

Stable Portfolio Compression for Index Tracking: Tree-Based Basis Selection, Local Ridge Mapping, and Turnover-Aware Replication*

Angus Cheung

December 14, 2025

Abstract

Stable Portfolio Compression (SPC) is a pipeline for compressing a large index into a compact basis of representative assets while preserving index-like exposures and return behavior. The method combines noise-robust correlation estimation (Ledoit–Wolf and PCA denoising), graph-theoretic structure learning via a Minimum Spanning Tree (MST), basis selection heuristics on trees, local Ridge regressions over neighborhoods, and an analytic or turnover-penalized quadratic program to map index weights to basis weights.

1 Introduction

Managing indices with hundreds or thousands of constituents is costly to trade, monitor, and rebalance. Stable Portfolio Compression (SPC) reduces this complexity by selecting a small, representative set of basis assets that preserves index returns and exposures, enabling efficient tracking and replication with far fewer positions.

SPC converts correlations into distances, builds a Minimum Spanning Tree to reveal market structure, selects diverse basis assets, and models non-basis assets via expanding-window Ridge regressions on nearby basis names. The resulting coefficient maps translate index weights into basis weights that approximate the full index.

A key use case is passive exchange-traded fund (ETF) index tracking, where holding SPC-selected basis assets (or liquid proxies) achieves near-benchmark exposure with lower costs and operational burden.

2 Notation

Let:

- N be the number of assets, T the number of time periods.
- $P_{t,i} > 0$ be the price of asset i at time t ; $\mathbf{P}_t \in \mathbb{R}_+^N$ the price vector.
- Returns $R_{t,i} := \frac{P_{t,i} - P_{t-1,i}}{P_{t-1,i}}$; $\mathbf{R}_t \in \mathbb{R}^N$.

*Project repository: <https://github.com/angus4718/stable-portfolio-compression/>

- $\widehat{\Sigma} \in \mathbb{R}^{N \times N}$ an estimator of covariance; \widehat{C} the correlation matrix with entries $\widehat{C}_{ij} \in [-1, 1]$.
- Distance matrix D with entries $D_{ij} = \sqrt{2(1 - \widehat{C}_{ij})}$.
- Tree $\mathcal{T} = (V, E)$, $V = \{1, \dots, N\}$, $E \subset V \times V$ with weights $w_{ij} = D_{ij}$ on edges.
- Basis set $B \subset V$, $|B| = k \ll N$.
- Index weights at time t : $\mathbf{w}_t^{\text{idx}} \in \mathbb{R}^N$, row-normalized across active assets.
- Basis weights at time t : $\mathbf{w}_t^B \in \mathbb{R}^k$.

3 Returns, Correlation, and Distance Estimation

3.1 Return construction

Given price levels $\{P_{t,i}\}_{t=0}^T$, define simple returns

$$R_{t,i} = \frac{P_{t,i}}{P_{t-1,i}} - 1, \quad t = 1, \dots, T. \quad (1)$$

Collect $\mathbf{R} \in \mathbb{R}^{T \times N}$; for downstream linear models we drop rows that are all NaN (e.g., due to listing churn), and align series by pairwise overlap when forming second-order moments.

3.2 Pairwise-overlap correlation

For assets i, j , let $\mathcal{T}_{ij} := \{t : R_{t,i}, R_{t,j} \text{ both observed}\}$ and $n_{ij} = |\mathcal{T}_{ij}|$. The Pearson correlation under pairwise deletion is

$$\widehat{C}_{ij} = \begin{cases} \frac{\sum_{t \in \mathcal{T}_{ij}} (R_{t,i} - \bar{R}_i)(R_{t,j} - \bar{R}_j)}{\sqrt{\sum_{t \in \mathcal{T}_{ij}} (R_{t,i} - \bar{R}_i)^2} \sqrt{\sum_{t \in \mathcal{T}_{ij}} (R_{t,j} - \bar{R}_j)^2}}, & n_{ij} \geq m_{\min}, \\ \text{NaN (to be denoised/imputed)}, & \text{otherwise,} \end{cases} \quad (2)$$

with \bar{R}_i computed on $\{t \in \mathcal{T}_{ij}\}$. Diagonals are fixed to $\widehat{C}_{ii} = 1$ if column i has at least m_{\min} observations. This overlap-aware construction maximizes data usage under missingness patterns typical of dynamic universes.

3.3 Ledoit–Wolf (LW) shrinkage

Let $\mathbf{X} \in \mathbb{R}^{T \times N}$ be column-median-imputed returns (residual NaNs $\rightarrow 0$) with rows centered by column means when forming the sample covariance $\mathbf{S} = \frac{1}{T-1} \mathbf{X}^\top \mathbf{X}$. The LW estimator shrinks \mathbf{S} toward a structured, well-conditioned target \mathbf{F} :

$$\widehat{\Sigma}_{\text{LW}} = \delta^* \mathbf{F} + (1 - \delta^*) \mathbf{S}, \quad \delta^* \in [0, 1]. \quad (3)$$

We adopt the constant-correlation (a.k.a. equicorrelation) or diagonal-scalar target, both common in practice:

- Diagonal-scalar: $\mathbf{F} = \bar{\sigma}^2 \mathbf{I}$ with $\bar{\sigma}^2 = \frac{1}{N} \sum_i S_{ii}$.
- Constant-correlation: $F_{ii} = S_{ii}$, $F_{ij} = \bar{\rho} \sqrt{S_{ii} S_{jj}}$ for $i \neq j$, with $\bar{\rho}$ the average sample correlation.

The optimal shrinkage intensity δ^* minimizes $\mathbb{E}\|\widehat{\Sigma} - \Sigma\|_F^2$ and admits a closed form:

$$\delta^* = \frac{\text{Var}[\mathbf{S}] - \text{Cov}[\mathbf{S}, \mathbf{F}] + \|\mathbf{S} - \mathbf{F}\|_F^2/T}{\|\mathbf{S} - \mathbf{F}\|_F^2}, \quad (4)$$

estimated consistently from the data (the scikit-learn LedoitWolf implementation computes this). LW guarantees positive definiteness and reduces eigenvalue spread, improving conditioning for downstream inversion (e.g., regressions, QP). Convert to correlation via

$$\widehat{C}_{ij} = \frac{\widehat{\Sigma}_{ij}}{\sqrt{\widehat{\Sigma}_{ii}\widehat{\Sigma}_{jj}}}, \quad \widehat{C}_{ii} = 1, \quad (5)$$

with diagonal clipping to handle any numerical artifacts. LW is particularly beneficial when N is not negligible versus T (high-dimensional regime), where \mathbf{S} is noisy and often ill-conditioned.

3.4 PCA-based correlation reconstruction: factor truncation and stabilization

Let \mathbf{C} be the raw pairwise-overlap correlation with NaNs replaced by 0 for numerical routines. Compute its spectral decomposition (or SVD) and retain k principal components to reconstruct a low-rank approximation:

$$\mathbf{C} = \mathbf{U}\Lambda\mathbf{U}^\top, \quad \lambda_1 \geq \dots \geq \lambda_N, \quad (6)$$

$$k = \begin{cases} \text{user-specified,} & \text{fixed } k, \\ \min\{k' : \sum_{\ell=1}^{k'} \lambda_\ell / \sum_{\ell=1}^N \lambda_\ell \geq \tau\}, & \text{threshold } \tau \in (0, 1), \end{cases} \quad (7)$$

$$\widehat{\mathbf{C}}_{\text{PCA}} = \mathbf{U}_{[:,1:k]}\Lambda_{1:k,1:k}\mathbf{U}_{[:,1:k]}^\top. \quad (8)$$

We then enforce symmetry, clip to $[-1, 1]$, and set $\text{diag}(\widehat{\mathbf{C}}_{\text{PCA}}) = \mathbf{1}$. Truncation suppresses high-frequency (idiosyncratic) noise that inflates small negative eigenvalues and destabilizes graph geometry. Economically, retained components approximate systematic factors; discarded components capture measurement noise and transient co-movements, which would otherwise distort distances on the asset manifold.

3.5 Correlation-to-distance transform

We map correlations to a Euclidean chordal distance on the unit sphere:

$$D_{ij} = \sqrt{2(1 - \widehat{C}_{ij})}. \quad (9)$$

This choice is justified by the embedding argument: If returns are z-scored, the correlation \widehat{C}_{ij} is the cosine of the angle θ_{ij} between unit vectors $\hat{\mathbf{r}}_i, \hat{\mathbf{r}}_j$. The chordal distance on the unit sphere satisfies $D_{ij} = \|\hat{\mathbf{r}}_i - \hat{\mathbf{r}}_j\|_2 = \sqrt{2(1 - \cos \theta_{ij})}$, giving a proper metric obeying the triangle inequality.

4 Minimum Spanning Tree (MST)

Given a complete weighted graph with weights $w_{ij} = D_{ij}$, we construct an MST \mathcal{T} via Kruskal's algorithm.

4.1 Kruskal's algorithm

Let $\mathcal{E} = \{(i, j) : i < j\}$ and sort by weights. Maintain a disjoint-set union (*UnionFind*) with path compression and union-by-size.

Input: edge list $\{(i, j, w_{ij})\}$. (10)

Procedure: Initialize UF over V ; $E_T \leftarrow \emptyset$. (11)

For edges in ascending w_{ij} : (12)

if UF.find(i) \neq UF.find(j) : (13)

UF.union(i, j), $E_T \leftarrow E_T \cup \{(i, j, w_{ij})\}$. (14)

Output: $\mathcal{T} = (V, E_T)$. (15)

4.2 Tree distances

Because \mathcal{T} is a tree, the path between any u, v is unique. Let the adjacency be $A(u) = \{(v, w_{uv})\}$. Distances $d_{\mathcal{T}}(u, v)$ are computed by BFS/DFS accumulation of edge weights from u . All-pairs can be computed by BFS from each node.

4.3 Degrees and path betweenness

We define the degree as $\deg(u) = |A(u)|$. To compute the *path betweenness count* $\text{btw}(u)$ (the number of unordered pairs (s, t) whose unique path in a tree passes through u) efficiently for all nodes, we use a linear-time algorithm based on rooting the tree, computing parent/child relations, and aggregating subtree sizes.

Algorithm (all-nodes betweenness in a tree). Given a weighted or unweighted tree $\mathcal{T} = (V, E)$ with $n = |V|$:

1. Root selection and parent map:

- 1.1. Pick an arbitrary root $r \in V$.
- 1.2. Run a DFS from r to compute a parent map $\text{parent}(\cdot)$ and an order of visitation (children discovered after parent).

2. Subtree sizes:

- 2.1. Initialize $\text{subtree}[u] \leftarrow 1$ for all u .
- 2.2. Traverse nodes in reverse DFS order; for each node u , for every child v (i.e., $\text{parent}(v) = u$), accumulate

$$\text{subtree}[u] \leftarrow \text{subtree}[u] + \text{subtree}[v].$$

3. Component sizes after removing u :

- 3.1. Removing u partitions the tree into components formed by each child-subtree of u and, if $u \neq r$, the “above- u ” remainder of size $n - \text{subtree}[u]$.
- 3.2. Thus the multiset of component sizes is

$$\mathcal{S}(u) = \{\text{subtree}[v] : \text{parent}(v) = u\} \cup (\{n - \text{subtree}[u]\} \text{ if } \text{parent}(u) \neq \emptyset).$$

4. Betweenness from component sizes:

- 4.1. Let $s_\ell \in \mathcal{S}(u)$ and denote $S = \sum_\ell s_\ell = n - 1$. The number of unordered pairs (s, t) whose unique path passes through u equals the number of pairs drawn from different components:

$$\text{betw}(u) = \sum_{\ell < m} s_\ell s_m = \frac{S^2 - \sum_\ell s_\ell^2}{2} = \frac{(n-1)^2 - \sum_\ell s_\ell^2}{2}.$$

5 Basis Selection on Trees

We seek $B \subset V$, $|B| = k$, that balances diversity (spread) and importance (weights). Let $w : V \rightarrow \mathbb{R}_{\geq 0}$ be a node weight signal (e.g., market cap) and \tilde{w} its min–max normalization to $[0, 1]$. Stickiness parameter $\beta \in [0, 1]$ biases selection toward nodes in a previous basis `prev_basis` by adding a scaled bonus.

5.1 Max-spread selection

Define the average tree distance of node v to all nodes by $\bar{d}(v) = \frac{1}{|V|} \sum_{u \in V} d_T(v, u)$. The initial score combines spread and weight; stickiness adds a bonus equal to a fraction of the maximum baseline score:

$$s_0^{\text{base}}(v) = \alpha \bar{d}(v) + (1 - \alpha) \tilde{w}(v), \quad (16)$$

$$s_0(v) = s_0^{\text{base}}(v) + \beta \mathbf{1}\{v \in \text{prev_basis}\} \cdot \max_{u \in V} s_0^{\text{base}}(u). \quad (17)$$

Seed with $b_1 = \arg \max_v s_0(v)$ and maintain the distance `nearest_dist`[v] = $\min_{b \in B} d_T(v, b)$ to the nearest basis for every node v .

Greedy fill until $|B| = k$: for each candidate $v \notin B$,

$$s^{\text{base}}(v) = \alpha \cdot \text{nearest_dist}[v] + (1 - \alpha) \cdot \tilde{w}(v), \quad (18)$$

$$s(v) = s^{\text{base}}(v) + \beta \mathbf{1}\{v \in \text{prev_basis}\} \cdot \max_{u \notin B} s^{\text{base}}(u). \quad (19)$$

Select $v^* = \arg \max_{v \notin B} s(v)$, update $B \leftarrow B \cup \{v^*\}$, and refresh `nearest_dist`[u] $\leftarrow \min\{\text{nearest_dist}[u], d_T(u, v^*)\}$ for all u .

5.2 Hub-and-branch selection

Compute a mixed centrality signal from degree and path betweenness, each min–max normalized to $[0, 1]$:

$$c(v) = \alpha \cdot \widetilde{\deg}(v) + (1 - \alpha) \cdot \widetilde{\text{betweenness}}(v), \quad (20)$$

$$h^{\text{base}}(v) = (1 - \gamma) \cdot \widetilde{c}(v) + \gamma \cdot \tilde{w}(v), \quad (21)$$

$$h(v) = h^{\text{base}}(v) + \beta \mathbf{1}\{v \in \text{prev_basis}\} \cdot \max_{u \in V} h^{\text{base}}(u), \quad (22)$$

with $\gamma \in [0, 1]$ trading off centrality vs weights. Choose hubs H as the top- h nodes by $h(v)$ where

$$h = \max\{1, \min(k - 1, \lfloor 0.2k \rfloor)\},$$

and add H to B .

Remove hubs from the tree and compute connected components (branches) $\{C\}$. For each branch, score candidates by distance to the nearest hub u^* mixed with weights and stickiness:

$$r^{\text{base}}(v) = \eta \cdot d_{\mathcal{T}}(v, u^*) + (1 - \eta) \cdot \tilde{w}(v), \quad (23)$$

$$r(v) = r^{\text{base}}(v) + \beta \mathbf{1}\{v \in \text{prev_basis}\} \cdot \max_{u \in C} r^{\text{base}}(u), \quad (24)$$

with $\eta \in [0, 1]$. Add one representative per branch in descending $r(v)$ until $|B| = k$. If short after exhausting branches, fill the remainder via a farthest-first greedy pass over residual nodes using `nearest_dist` updates as in max-spread selection.

6 Local Ridge Regressions

For each non-basis asset $i \in V \setminus B$, select q nearest basis neighbors by the full-sample distance matrix \mathbf{D} :

$$\mathcal{N}_i = \operatorname{argmin}_{B' \subset B, |B'|=q} \sum_{b \in B'} D_{ib} \quad (\text{implemented by sorting } \{D_{ib}\}_{b \in B} \text{ and taking } q \text{ smallest}). \quad (25)$$

Fit Ridge regression without intercept on the overlapping sample:

$$\min_{\boldsymbol{\beta}_i \in \mathbb{R}^{|\mathcal{N}_i|}} \sum_{t \in \mathcal{T}_i} \left(R_{t,i} - \sum_{b \in \mathcal{N}_i} \beta_{i,b} R_{t,b} \right)^2 + \lambda \|\boldsymbol{\beta}_i\|_2^2, \quad (26)$$

with closed-form solution $\hat{\boldsymbol{\beta}}_i = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$ where \mathbf{y} is $R_{t,i}$ and \mathbf{X} columns are $\{R_{t,b}\}_{b \in \mathcal{N}_i}$.

Reconstruction and diagnostics. For all times with non-missing neighbor returns, define reconstructed returns

$$\hat{R}_{t,i} = \sum_{b \in \mathcal{N}_i} \hat{\beta}_{i,b} R_{t,b}. \quad (27)$$

Compute RMSE and R^2 per asset on the fit sample:

$$\text{RMSE}_i = \sqrt{\frac{1}{|\mathcal{T}_i|} \sum_{t \in \mathcal{T}_i} (R_{t,i} - \hat{R}_{t,i})^2}, \quad (28)$$

$$R_i^2 = 1 - \frac{\sum_t (R_{t,i} - \hat{R}_{t,i})^2}{\sum_t (R_{t,i} - \bar{R}_i)^2}. \quad (29)$$

7 Weight Mapping: Assets to Basis

Let $\mathbf{A}_t \in \mathbb{R}^{N \times k}$ be the asset-to-basis mapping at time t formed by:

1. Identity rows for basis assets present in the universe: $(\mathbf{A}_t)_{b,b} = 1$.
2. For non-basis i , coefficient row $(\mathbf{A}_t)_{i,b} = \hat{\beta}_{i,b}$ from the selected coefficient pivot date (static or most recent $\leq t$).

Given index weights $\mathbf{w}_t^{\text{idx}} \in \mathbb{R}^N$, the analytic mapping is

$$\mathbf{w}_t^B = \mathbf{A}_t^\top \mathbf{w}_t^{\text{idx}} \quad (\text{or equivalently } \mathbf{w}_t^{B\top} = (\mathbf{w}_t^{\text{idx}})^\top \mathbf{A}_t). \quad (30)$$

This does not constrain $\mathbf{1}^\top \mathbf{w}_t^B = 1$, leaving residual mass as cash if the mapping under-fills.

7.1 Turnover-penalized quadratic program

To stabilize rebalancing, solve

$$\min_{\mathbf{w} \in \mathbb{R}^k} \underbrace{\left\| \mathbf{A}_t \mathbf{w} - \mathbf{w}_t^{\text{idx}} \right\|_2^2}_{\text{fit}} + \lambda \underbrace{\left\| \mathbf{w} - \mathbf{w}_{t-1}^B \right\|_2^2}_{\text{turnover}}, \quad (31)$$

with unconstrained \mathbf{w} (negative basis weights allowed). The objective is strictly convex with solution

$$\mathbf{w}_t^B = \left(\mathbf{A}_t^\top \mathbf{A}_t + \lambda \mathbf{I} \right)^{-1} \left(\mathbf{A}_t^\top \mathbf{w}_t^{\text{idx}} + \lambda \mathbf{w}_{t-1}^B \right), \quad (32)$$

used directly or via a numerical solver (e.g., SLSQP). In the implementation, when no previous \mathbf{w}_{t-1}^B exists, the mapping defaults to the analytic $\mathbf{A}_t^\top \mathbf{w}_t^{\text{idx}}$.

8 Index and Replicate Returns; Backtesting

Given prices \mathbf{P}_t and returns \mathbf{R}_t , index return uses lagged weights

$$r_t^{\text{idx}} = \sum_{i=1}^N w_{t-1,i}^{\text{idx}} R_{t,i}. \quad (33)$$

Replicate return uses basis returns and lagged basis weights

$$r_t^{\text{rep}} = \sum_{b \in B} w_{t-1,b}^B R_{t,b}. \quad (34)$$

Active return $a_t = r_t^{\text{rep}} - r_t^{\text{idx}}$.

8.1 Tracking error and diagnostics

$$\text{TE}_{\text{RMSE}} = \sqrt{\frac{1}{T'} \sum_t a_t^2}, \quad (35)$$

$$\text{running_RMSE}_t = \sqrt{\frac{1}{t} \sum_{s=1}^t a_s^2}, \quad (36)$$

$$\text{ann. TE} \approx \text{TE}_{\text{RMSE}} \cdot \sqrt{12}, \quad (37)$$

$$R^2 := 1 - \frac{\sum_t (r_t^{\text{idx}} - r_t^{\text{rep}})^2}{\sum_t (r_t^{\text{idx}} - \bar{r}^{\text{idx}})^2}. \quad (38)$$

Turnover at date t for replicate weights \mathbf{w}_t^B :

$$\text{turnover}_t = \frac{1}{2} \sum_{b \in B} |w_{t,b}^B - w_{t-1,b}^B|. \quad (39)$$

9 Synthetic Market Simulator

We present a configurable stochastic generator for monthly equity prices, shares outstanding, market capitalizations, and index membership. The model couples sector structure, heavy-tailed dispersion in shares outstanding, sector-correlated returns with optional drift, and Poisson-driven universe churn (additions/removals). It produces monthly panels of prices and market caps, market-cap-normalized index weights, and per-month constituent sets for a top- K index.

9.1 Universe, sectors, and shares outstanding

Let \mathcal{S} be the sector set with $|\mathcal{S}| = S$. Sector $s \in \mathcal{S}$ initially contains n_s assets. Each asset i is assigned a sector $s(i) \in \mathcal{S}$. Shares outstanding Q_i are heavy-tailed to induce skewed market caps:

$$Q_i \sim \text{Shifted-Pareto}(\alpha, \text{scale}) \quad \text{with} \quad Q_i \leftarrow \max((X + 1) \cdot \text{scale}, 1), \quad X \sim \text{Pareto}(\alpha). \quad (40)$$

Initial prices are fixed at $P_{0,i} = 100$, yielding initial market caps $M_{0,i} = Q_i P_{0,i}$. Tickers are uniquely identified and tagged with sector metadata and creation/removal months.

9.2 Within-sector correlation and monthly returns

For each month and sector s , we generate correlated Gaussian shocks with constant within-sector correlation. Let n_s be the number of active assets in sector s . Define the correlation matrix $\mathbf{C}_s \in \mathbb{R}^{n_s \times n_s}$ by

$$(\mathbf{C}_s)_{ij} = \begin{cases} 1, & i = j, \\ \rho_s, & i \neq j, \end{cases} \quad \rho_s \in [0, 1), \quad (41)$$

and, if needed, stabilize to positive-definiteness via a small diagonal jitter. Draw standard normals $\mathbf{z}_s \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n_s})$ and set

$$\boldsymbol{\epsilon}_s = \mathbf{L}_s \mathbf{z}_s, \quad \mathbf{L}_s \mathbf{L}_s^\top = \mathbf{C}_s, \quad (42)$$

$$\sigma_s = \frac{\sigma_s^{\text{ann}}}{\sqrt{12}}, \quad \mu_s \sim \mathcal{U}[\mu_{\min}, \mu_{\max}], \quad (43)$$

$$\mathbf{r}_s = \mu_s \mathbf{1} + \sigma_s \boldsymbol{\epsilon}_s. \quad (44)$$

Each asset i in sector s takes $R_{t,i} = (\mathbf{r}_s)_i$. Prices evolve with clipping to avoid non-positive values:

$$P_{t,i} = \max(P_{t-1,i}(1 + R_{t,i}), \varepsilon), \quad \varepsilon > 0 \text{ small.} \quad (45)$$

Market caps update as $M_{t,i} = Q_i P_{t,i}$.

9.3 Additions and removals (churn)

At month t , the active universe is U_t . The number of removals and additions scale with $|U_t|$ via Poisson draws:

$$N_{\text{rem}} \sim \text{Poisson}(\lambda_{\text{rem}} |U_t|), \quad N_{\text{add}} \sim \text{Poisson}(\lambda_{\text{add}} |U_t|). \quad (46)$$

Removals are sampled without replacement from the bottom- p fraction (by current market cap) of removable names in U_t . Additions are created with sector drawn uniformly from \mathcal{S} ; each new asset receives fresh Q_i from the shifted-Pareto law, is initialized at $P_{t,i} = 100$, and begins contributing to prices and market caps from its creation month onward.

9.4 Index formation and weights

Let $I_t \subseteq U_t$ be the set of top- K tickers by market cap at time t (with K configurable). The market-cap-weighted index uses weights normalized over active names:

$$w_{t,i}^{\text{mkt}} = \frac{M_{t,i}}{\sum_{j \in U_t} M_{t,j}}, \quad i \in U_t. \quad (47)$$

The simulator records, per month, the price panel $\{P_{t,i}\}$, market caps $\{M_{t,i}\}$, weights $\{w_{t,i}^{\text{mkt}}\}$, and constituent sets I_t .

9.5 Per-constituent weight capping algorithm

We optionally enforce a per-constituent maximum weight $\bar{w} \in (0, 1]$ on market-cap-normalized weights to control concentration. For each month t , let $\mathbf{w}_t \in \mathbb{R}^N$ denote the row-normalized weights over the set of active tickers $A_t \subseteq U_t$ (inactive tickers are treated as missing). The goal is to produce adjusted weights $\tilde{\mathbf{w}}_t$ that satisfy:

$$(i) \quad \tilde{w}_{t,i} \geq 0 \quad \forall i \in A_t, \quad (ii) \quad \sum_{i \in A_t} \tilde{w}_{t,i} = 1, \quad (48)$$

$$(iii) \quad \tilde{w}_{t,i} \leq \bar{w} \quad \forall i \in A_t, \quad (iv) \quad \tilde{w}_{t,i} = \text{NaN} \text{ for } i \notin A_t. \quad (49)$$

A necessary and sufficient feasibility condition is

$$\bar{w} \cdot |A_t| \geq 1, \quad (50)$$

otherwise the cap is infeasible for period t .

We use an iterative clip-and-redistribute scheme that preserves the relative proportions among uncapped names at each step:

Initialization. Set $\tilde{\mathbf{w}}_t^{(0)} = \mathbf{w}_t$ on A_t and define the active receiver set $R_t^{(0)} = \{i \in A_t : \tilde{w}_{t,i}^{(0)} < \bar{w}\}$.

Iteration. For $k = 0, 1, 2, \dots$:

1. Identify violators $V_t^{(k)} = \{i \in A_t : \tilde{w}_{t,i}^{(k)} > \bar{w}\}$. If $V_t^{(k)} = \emptyset$, stop and return $\tilde{\mathbf{w}}_t^{(k)}$.

2. Compute total excess

$$E_t^{(k)} = \sum_{i \in V_t^{(k)}} (\tilde{w}_{t,i}^{(k)} - \bar{w}). \quad (51)$$

Clip violators: set $\tilde{w}_{t,i}^{(k+1)} = \bar{w}$ for $i \in V_t^{(k)}$ and $\tilde{w}_{t,i}^{(k+1)} = \tilde{w}_{t,i}^{(k)}$ for $i \notin V_t^{(k)}$.

3. Determine receivers (currently uncapped and positive):

$$R_t^{(k)} = \{i \in A_t : \tilde{w}_{t,i}^{(k+1)} < \bar{w}, \tilde{w}_{t,i}^{(k+1)} > 0\}. \quad (52)$$

If $R_t^{(k)} = \emptyset$ and $E_t^{(k)} > 0$, feasibility is violated; abort.

4. Redistribute the excess proportionally to receivers:

$$S_t^{(k)} = \sum_{j \in R_t^{(k)}} \tilde{w}_{t,j}^{(k+1)}, \quad \tilde{w}_{t,i}^{(k+1)} \leftarrow \tilde{w}_{t,i}^{(k+1)} + E_t^{(k)} \cdot \frac{\tilde{w}_{t,i}^{(k+1)}}{S_t^{(k)}} \quad \forall i \in R_t^{(k)}. \quad (53)$$

10 Backtesting Protocols

10.1 Static scripts

1. Basis selection from all data to date: compute \mathbf{D} , \mathcal{T} , and B via max-spread or hub-branch.
2. Local Ridge: fit $\hat{\beta}_{i,b}$ using neighbors from \mathbf{D} .
3. Weight mapping: compute \mathbf{A} and \mathbf{w}^B (optionally with turnover penalty).

10.2 Time-series analysis

At each date t :

1. Compute \mathbf{D}_t from data up to t and select basis B_t ($|B_t| = k$).
2. Fit local Ridge mapping coefficients $\hat{\beta}^{(t)}$ using data up to t .
3. Map $\mathbf{w}_t^{\text{idx}} \mapsto \mathbf{w}_t^B$ using the pivot at date t ; optionally solve (31).
4. Evaluate with next-period returns using lagged weights.

10.3 Baseline strategy

Top- N baseline replicator: at each t , take the N_0 largest index weights and renormalize:

$$\tilde{w}_{t,i} = \begin{cases} \frac{w_{t,i}^{\text{idx}}}{\sum_{j \in \text{Top-}N_0} w_{t,j}^{\text{idx}}}, & i \in \text{Top-}N_0, \\ 0, & \text{otherwise.} \end{cases} \quad (54)$$

Replicate returns use $\tilde{\mathbf{w}}_{t-1}$ and full-universe returns.

11 Conclusions

SPC integrates robust covariance estimation, tree geometry, greedy max-spread and hub-branch selections, local regressions, and turnover-aware mapping to compress index exposures into a compact, tradable basis with controlled tracking error and turnover.