

Box A.2: DE-MEANING AND STANDARDIZATION OF FEATURES

Most raw features in the dataset – including prices, volumes, and macroeconomic levels – introduce significant risks if incorporated directly into predictive modelling. They exhibit (a) ***strong drifts*** (long-run systematic upward or downward movements) – which can lead to *spurious trend learning* (i.e., the model interpreting shared trends as causal rather than focusing on deviations that contain predictive signals); (b) ***seasonality*** – which may be misinterpreted as a predictive driver of returns instead of highlighting when values are unusually high or low relative to their seasonally typical baseline; and (c) ***wide scale differences*** across variables and tickers, creating risks of *non-comparability* due to inconsistent units, and *volatility of coefficients* (if features are on very different scales, coefficients will be on wildly different magnitudes, leading to unstable estimation and unfair penalization under regularization).

To mitigate these issues, a two-step transformation was applied:

1. De-meaning (centring): most features were de-meansed using trailing rolling or exponentially weighted averages, shifted by one day to prevent look-ahead. This removes slow-moving drifts and regular seasonal components, forcing the model to focus on deviations from typical values rather than absolute levels.

For these features, de-meaning subtracts a trailing mean, centring the series around zero:

$$\tilde{x_t} = x_t - \mu_{t,w}, \quad \mu_{t,w} = \text{rolling/EWM mean over window } w \text{ (shifted by one day)}$$

Slow-moving macroeconomic and earnings-level variables were normalized using global statistics (full-sample mean and standard deviation) without a rolling de-mean step.

2. Standardization (z-scoring): After centring, most features were standardized by dividing their de-meansed values by trailing volatility (standard deviation), computed on a rolling or exponentially weighted basis and shifted by one day.

$$z_t = \frac{x_t - \mu_{t,w}}{\sigma_{t,w}}, \quad \mu_{t,w} = \text{rolling/EWM mean}, \quad \sigma_{t,w} = \text{rolling/EWM std, (all shifted by one day)}$$

This expresses each observation in units of standard deviations from its recent norm, which makes features directly comparable across time and across tickers, regardless of their original scale (price, ratios, macro levels, etc.). For families not de-meansed (macroeconomic levels and earnings levels), standardization was applied using global statistics. The result is that all features are placed on a common, volatility-adjusted scale (typically between -2 and +2 for most observations).

Detailed Methodological Documentation: targets (future returns), event flags (e.g., event, earnings_date_*), calendar/timing variables, and binary indicators were auto excluded from transformation. All remaining numeric variables were assigned to families based on their economic and statistical characteristics and for each family, tailored normalization methods were applied:

Rolling windows (=60 trading days): Rolling normalization was applied to Price levels, Distance-to-trend/Momentum, and Market return features. A 60-day (~3 months) trailing window provides a stable local baseline while remaining responsive to medium-term changes, ensuring deviations are measured relative to a recent but robust benchmark.

Exponentially weighted moving averages (EWM): EWM normalization was applied to Volume/Turnover and Volatility/Beta features. Unlike rolling windows, EWM assigns greater weight to more recent observations, enabling faster adaptation to regime shifts in trading activity or volatility. For volume features, a half-life of 25 days was used, meaning the effective weight of an observation decreases by 50% every ~25 trading days. This choice ensures that the scale reflects the most recent trading patterns while still retaining a meaningful history. For volatility/beta features, a half-life of 40 days was chosen, providing a more gradual adaptation. This reflects the fact that volatility and beta tend to evolve more slowly than trading volumes, requiring a longer memory to establish a stable baseline.

Global statistics (full sample mean and standard deviation). Global normalization was applied to the slow-moving, low-frequency variables of Macroeconomic and Earnings level features, using a constant mean and standard deviation computed once over the full sample. This approach ensures stability for slow-moving, low-frequency series where short windows would generate noisy or artificial fluctuations. For Macroeconomic indicators, global statistics provide a consistent baseline for variables that change slowly and are typically released monthly or quarterly (e.g., CPI, GDP, unemployment). For Earnings levels, global normalization avoids overstating volatility between sparse quarterly updates while still putting different firms and earnings metrics on comparable scales.

Finally, to prevent look-ahead bias, all rolling and EWM statistics were computed on a trailing-only basis and shifted by one trading day. This ensures that the standardized value at time t uses information available only up to $t-1$, preserving the integrity of out-of-sample evaluation.

Contribution to modelling: by transforming features into rolling or exponentially weighted demeaned and standardized variants, the dataset was rendered more stationary and better aligned for modelling. This process produced cleaner signals, enhanced robustness, and ultimately supported more consistent and reliable predictive performance.

Normalization & De-Meaning Methods by Feature Family			
De-Meaning & Normalization Method	Feature Family	Features Included	Outputs
Rolling mean / std (60d)	Price levels	<code>open_</code> , <code>high_</code> , <code>low_</code> , <code>close_</code>	Raw + *_dm, *_z
	Distance-to-trend / Momentum	<code>bollinger_position_</code> , <code>macd_</code> , <code>rsl_</code> , <code>sma_ratio_20_</code> , <code>ret_1d_</code> , <code>ret_5d_</code> , <code>perf_vs_sp_ret_5d_</code> , <code>rolling_avg_10d_</code> , <code>mom_20d_</code> , <code>mom_63d_</code> , <code>sharp_20d_</code> , <code>daily_pct_change_</code> , <code>sp_ret_1d</code> , <code>sp_ret_5d</code>	Raw + *_dm, *_z
EWM mean / std (half-life 25d)	Market returns	<code>sentiment_</code> , <code>oi_score_</code>	Raw + *_dm, *_z
	Sentiment		Raw + *_dm, *_z
	Volume / Turnover	<code>volume_</code> , <code>vol_z_60d_</code>	Raw + *_dm, *_z
	Volatility / Beta	<code>atr_</code> , <code>beta_sp_60d_</code> , <code>bollinger_upper_</code> , <code>bollinger_lower_</code> , <code>rv_20d_</code> , <code>semivol_20d_</code> , <code>idio_vol_20d_</code> , <code>dd_60d_</code>	Raw + *_dm, *_z
Global mean / std (full sample)	Macroeconomic levels	<code>cpi_raw</code> , <code>cpi_mom</code> , <code>cpi_yoy</code> , <code>exports_raw</code> , <code>exports_mom</code> , <code>exports_yoy</code> , <code>imports_raw</code> , <code>imports_mom</code> , <code>imports_yoy</code> , <code>trade_balance_raw</code> , <code>trade_balance_mom</code> , <code>trade_balance_yoy</code> , <code>unemp_rate</code> , <code>gdp</code> , <code>debt_service_ratio</code> , <code>interest_rate_daily</code> , <code>reported_eps_</code> , <code>estimated_eps_</code> , <code>revenue_million_</code> , <code>surprise_percent_</code> , <code>earning_date_</code> , <code>earnings_date_</code>	Raw + *_gz
	Earnings levels		Raw + *_gz