

Master Thesis

Edit: Master Thesis Title

Jorge Eduardo FRÍAS NAVARRETE

Submitted in partial fulfillment of the requirement for the degree of:

Master of Science

Student ID: 012329686
Degree programme: Quantitative Finance
Supervisor: Univ.Prof. David PREINERSTORFER, Ph.D.
Date of Submission: September 04, 2025

*Department of Finance, Accounting and Statistics.
Vienna University of Economics and Business.
Welthandelsplatz 1, 1020 Vienna, Austria.*

Abstract

Here goes my abstract text. Here goes my abstract text. Here goes my abstract text. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Suspendisse eu dolor luctus, rhoncus leo in, commodo turpis. Aenean sed enim in sem euismod porta. Vivamus tempor lorem nec eros rhoncus, eu hendrerit libero tincidunt. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos himenaeos. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Quisque dapibus turpis quis nibh molestie dapibus. Aliquam erat volutpat. Integer et odio nec mauris sollicitudin mattis.

Table of contents

List of figures

List of tables

1 Introduction

The first cryptocurrency, Bitcoin, was created in 2009 by Satoshi Nakamoto, who presented it as a peer-to-peer electronic coin with secured and verified transactions through an encrypted proof-of-work mechanism (?). As originally proposed, Bitcoin was designed as an alternative, decentralized cash system offering low-cost and near-real-time transactions, while avoiding currency controls imposed by national governments or financial institutions¹ (?). These features quickly attracted widespread public attention. However, due to its high volatility, researchers have questioned its role as a purely digital currency and instead classified it as an investment or speculative asset (?, ?, ?).

Since then, the cryptocurrency market has expanded rapidly, giving rise to thousands of new coins. In the second quarter of 2025, the total cryptocurrency market capitalization amounted to nearly 3.5 trillion USD, according to data from CoinGecko (?). Despite this rapid growth, perceptions of cryptocurrencies remain divided. Some view them as investments tied to the underlying technologies, such as blockchain and smart contracts, or simply as a form of speculation (?, ?). Others, however, see them as bubbles, fraud schemes, or scams, often driven by internet and social media marketing—for example, rug pulls involving so-called “memecoins,” or, more recently, the LIBRA cryptocurrency scandal in February 2025, when the coin was promoted by Argentinian president Javier Milei, soared in value, and collapsed only a few hours later (?, ?, ?).

As mentioned earlier, a key characteristic of cryptocurrencies is their high volatility, which greatly exceeds that of other traditional assets such as equity indices, gold, silver, foreign exchange currencies, and commodities (?, ?). According to the standard asset pricing theory, investors should be compensated for bearing such risks. The principle that higher risk should be associated with higher expected returns is central in finance, beginning with the capital asset pricing model (CAPM) of Sharpe (?) and Lintner (?), and later extended by Merton (?), who introduced state variables to capture changes in investment and consumption decisions through the intertemporal CAPM, and by

¹Contrary to the common belief, Bitcoin is not anonymous. All Bitcoin transactions are publicly visible in the network and only the identity of the user behind a Bitcoin address is unknown, until their identity is revealed through a purchase or another action. See Meiklejohn et al. (?) and <https://bitcoin.org/en/you-need-to-know>.

1 Introduction

Ross (?), who formalized multi-factor risk pricing through the arbitrage pricing theory (APT). In particular, the APT shows that, in the absence of arbitrage opportunities, asset returns can be represented by a linear factor model, where returns are explained by their exposures to systematic risk factors. In empirical applications, this relation is often estimated through time-series regressions (?). Let $r_{i,t+1} \in \mathbb{R}$ denote the excess return on asset i from period $t-1$ to t , for $i = 1, \dots, N$ and $t = 1, \dots, T$. Let $f_{t+1} \in \mathbb{R}^K$ be a $K \times 1$ vector of risk factors. The model can then be written as

$$r_{i,t} = \alpha_{i,t-1} + \beta'_{i,t-1} f_t + \epsilon_{i,t},$$

where $\beta_{i,t-1} \in \mathbb{R}^K$ measures the exposure of asset i to the risk factors, $\alpha_{i,t-1}$ represents a pricing error (equal to zero under correct specification), and $\epsilon_{i,t}$ is the idiosyncratic component of returns.

A major challenge of the framework described above is identifying the set of factors that best capture asset returns, as these factors are not directly observable. This raises the question of whether they truly explain the cross-section of excess returns or whether such returns should instead be attributed to asset mispricing. This motivates the main questions addressed in this thesis:

- Which factors account for the variation in cryptocurrency returns?
- To what extent can the return cross-section be explained by systematic risk factors?
- Does allowing for dynamic factor loadings improve the prediction of cross-sectional excess returns?

The main goal of this thesis is to apply established factor models from the financial literature to a large panel of cryptocurrency data and to compare their predictive performance under static and dynamic loadings. In particular, I replicate the approaches of Kelly et al. (?) and Bianchi & Babiak (?) for the cryptocurrency market. The analysis relies on a model that allows factor loadings to vary over time through observable characteristics, using the Instrumented Principal Component Analysis (IPCA) methodology.

1.1 Literature review

Linear factor pricing models play a fundamental role in the field of finance. Building on the theoretical foundations of APT, a large body of academic research have worked to identify the sources of economic risks and the factors that explain the cross-section of asset returns. Broadly speaking, two main strands have emerged in the empirical literature (?).

One strand of the literature pre-specifies the factors f_{t+1} and represents them with long-short portfolios, often referred to as factor-mimicking portfolios or sorted portfolios. These long-short portfolios are based on well-established knowledge of the empirical behavior of asset returns and are therefore treated as fully observable (?). The main drawback of this approach is that it presumes a prior understanding of the cross-sectional dynamics of asset returns, even though such knowledge is incomplete or imperfect.

Although the construction of each factor varies across studies, the process typically involves sorting assets into quintiles (or deciles) based on a given characteristic and forming the factor return as the difference between the top and bottom groups. Fama & French (?) were the first to formalize this approach in the context of linear factor models, introducing a three-factor model (FF3) that included the market, size, and value factors to explain stocks and bond returns. Carhart (?) expanded the FF3 by adding a momentum factor, which captures the one-year asset momentum, forming in this way a 4-factor model. Later, Fama & French (?) extended the FF3 by incorporating profitability and investment factors, creating a 5-factor model to capture additional stock return variation beyond size and value.

The number of risk factors proposed in the literature is vast, with hundreds of them reported across different studies (?, ?). Feng et al. (?) developed a model selection framework to evaluate the contribution of newly proposed factors, finding that most are redundant relative to existing ones. Hou et al. (?) and A. Y. Chen & Zimmermann (?) replicated 452 and 319 long-short strategies from the literature, respectively. Hou et al. failed to reproduce the results of more than half of predictors in their set, finding most of them statistically insignificant and concluding that many published return predictors are not reliable. By contrast, Chen and Zimmerman showed that nearly all of the literature results can be successfully replicated.

A second strand of research views the factors as latent and applies data-compression techniques, such as Principal Component Analysis (PCA), to simultaneously extract common factors and estimate their betas directly from the panel of realized returns (?).

1 Introduction

This method derives factors purely from a statistical criteria and therefore requires no prior knowledge of the cross-sectional behavior of returns. Its main limitation, however, is that PCA can only estimate static loadings, implying that asset exposures to systematic risk are assumed constant over time. Moreover, PCA cannot incorporate additional information beyond returns, which restricts its ability to identify more appropriate asset pricing models (?).

The pioneers in this approach are Chamberlain & Rothschild (?) and Connor & Korajczyk (?). Chamberlain & Rothschild (?) defined the concept of approximate factor structure and showed that asset returns on large markets can be represented by a small number of common factors that can be extracted with PCA, as long as the covariance matrix of asset returns has K unbounded eigenvalues. Building on this, Connor & Korajczyk (?) developed an econometric method using asymptotic principal components that estimates latent factors and their loadings from large panels of returns, providing consistent APT-based performance measures and an application to portfolio evaluation.

More recently, Kelly et al. (?) introduced the Instrumented PCA (IPCA). Unlike PCA, which assumes static factor loadings, IPCA allows loadings to vary with observable asset characteristics such as size, volatility, or momentum. These characteristics serve as instruments for conditional loadings, enabling the method to incorporate more information than returns alone and to handle unbalanced panels of data. Bali et al. (?) extended the IPCA approach to a joint factor model that explains the risk-return trade-off across different asset classes –bonds, stocks, and options–. In a related work, Z. Chen et al. (?) proposed the Regressed PCA (RPCA), which extracts common latent factors across stocks, bonds, and options by combining cross-sectional Fama–MacBeth regressions (?) on asset characteristics with standard PCA.

While most of the literature has focused on understanding stock market returns, a growing body of research has examined the dynamics of cryptocurrency returns. Inspired by the FF3 model in equities, Y. Liu et al. (?) and W. Liu et al. (?) construct a similar three-factor model for cryptocurrency returns using market, size, and momentum factors. Using weekly data, they show that this model captures a large share of cryptocurrency returns and, in particular, reveals strong anomaly effects in the momentum and size factors. However, Jung & Park (?) show that the three-factor model of Y. Liu et al. (?) explains only about one-third of cryptocurrency return variation. They attribute the remaining variation to a common component outside the three-factor model, closely linked to the value of fiat money, highlighting the role of global macroeconomic variables in cryptocurrency pricing. Further work by Y. Liu & Tsyvinski (?) shows that cryptocurrency returns are also linked to network factors, which capture

user adoption. They also find strong momentum effects and show that investor attention can predict future returns. Building on these findings, Cong et al. (?) show that value and network adoption provide strong risk premia across more than 4,000 cryptocurrencies. They propose a five-factor “C-5” model –market, size, momentum, value, and network–that performs better than earlier models in- and out-of-sample, and also report market segmentation across different categories of cryptocurrencies.

Studies adopting a latent-factor structure include Bouri et al. (?) and Bianchi & Babiak (?). Bouri et al. (?) apply a regime-switching factor model, where the comovement of cryptocurrency returns depends on market states. They show that accounting for these state-dependent comovements improves the forecasting performance of major cryptocurrencies compared to standard PCA and a random-walk model. In contrast, Bianchi & Babiak (?) apply the IPCA model to the cryptocurrency market, constructing 32 characteristics to instrument the dynamic factor loadings. They show that this time-varying latent-factor framework measures the variation in realized returns more accurately than conventional observable-factor models or standard PCA, both at the daily and weekly frequency. They also find that characteristics related to speculative demand and liquidity are the most significant in capturing the systematic mispricing of returns.

1.2 Data concerns

One of the main challenges in this thesis was obtaining a large panel of cryptocurrency data. I extracted market data from the free [CoinCodex](#) API, which provides access to the full historical data of the cryptocurrencies listed on its platform. In contrast, most crypto market data providers –also called coin-ranking sites, such as CoinMarketCap, CoinGecko, CryptoCompare (CoinDesk)– offer limited access to historical data (usually one year) or none at all without a paid subscription. Some exchange platforms, such as Bybit, Binance, Coinbase, and Cex, allow users to extract market data for free through their public APIs. However, the number of cryptocurrencies (and thus, the cross-section) available from these sources was relatively small compared with CoinCodex, and the time span was shorter ².

The choice of which data source is appropriate for scientific research is subject to debate. For example, Alexander & Dakos (?) examine different cryptocurrency data

²For example, Bitcoin data started from late 2013 in CoinCodex, compared to November, 2022 in Bybit, January, 2019, in Binance, and June, 2021, in Coinbase. The available cryptocurrencies paired with Tether USD (USD) were 763 in Bybit, 623 in Binance, and 116 (USD) in Coinbase.

1 Introduction

providers and find inconsistencies in regression estimates, suggesting that the source of cryptocurrency data can influence empirical results. Moreover, they document distorted coin prices on coin-ranking sites, caused by inflated or artificial trading volumes³, emphasizing the importance of using traded data from crypto exchanges. By contrast, Vidal-Tomás (?) argue that coin-ranking sites use the same underlying process as crypto exchanges and other platforms to compute a cryptocurrency price, and they report no significant differences in empirical results when using alternative data sources. To address these concerns, I apply a series of pre-processing filters, described in Section ??, to mitigate the impact of potential inaccuracies in my dataset.

The remainder of the thesis is structured as follows. Section 2 summarizes the IPCA model, the estimation strategy and the performance measures applied in the analysis. Section 3 describes the data extraction and the sample construction process. Section 4 presents the empirical findings, and Section 5 concludes.

³Coin-ranking sites rank coins and exchanges by trading volume and market capitalization. As highlighted by Alexander & Dakos (?), the prices quoted on some of these sites are calculated by aggregating the prices from hundreds of exchanges using a volume-weighted average. Because many exchanges artificially inflate their volume to boost their position in the rankings, the resulting aggregated prices are influenced by fake volumes and therefore inconsistent with traded prices.

2 Methodology

In this section, I present the main method used in this thesis: Instrumented Principal Component Analysis (IPCA), introduced by Kelly et al. (?). IPCA estimates latent factors and dynamic factor loadings by linking them to observable asset-specific characteristics. Unlike standard PCA, which assumes static loadings and relies uniquely on return data, IPCA allows factors loadings to vary with asset characteristics, such as size, volatility, volume, or momentum, which act as instruments for the conditional loadings. Moreover, it enables the estimation of K factor loadings directly from the panel of asset characteristics. Another advantage is that IPCA can be applied to unbalanced panels, which is particularly useful in the cryptocurrency market where new coins are regularly introduced and others become inactive or unavailable, making missing data in the cross-section very common.

2.1 IPCA model and estimation

Consider a linear factor model. Let $r_{i,t+1} \in \mathbb{R}$ denote the excess return on cryptocurrency i from period t to $t+1$, for $i = 1, \dots, N$ and $t = 1, \dots, T$. The general IPCA model specification is defined as

$$r_{i,t+1} = \alpha_{i,t} + \beta'_{i,t} f_{t+1} + \epsilon_{i,t+1}, \quad (2.1)$$

with

$$\alpha_{i,t} = z'_{i,t} \Gamma_{\alpha} + \nu_{\alpha,i,t}, \quad \beta_{i,t} = z'_{i,t} \Gamma_{\beta} + \nu_{\beta,i,t},$$

where $f_{t+1} \in \mathbb{R}^K$ is the $K \times 1$ vector of latent factors. The $K \times 1$ vector $\beta_{i,t}$ captures the dynamic factor loadings, which may depend on observable cryptocurrency characteristics contained in the $L \times 1$ vector of instruments $z_{i,t}$. The main idea is that linking model parameters to observable characteristics allows expected returns to adjust more quickly to new information than when using parameter estimates from rolling window time-series regressions (?). This link is captured through the $L \times K$ matrix Γ_{β} , which

2 Methodology

maps a potentially large number of cryptocurrency characteristics L into a small number K of latent factor loadings. Similarly, the $L \times 1$ vector Γ_α maps characteristics to anomaly intercepts. Finally, the terms $\nu_{\alpha,i,t}$ and $\nu_{\beta,i,t}$ are residuals that capture variation in loadings orthogonal to the observable instruments.

In IPCA, two specifications can be considered. As discussed earlier, characteristics are used as instruments for the time-variation in conditional loadings, so that the mapping $z_{i,t} \mapsto \beta_{i,t}$ is determined by the low-dimensional matrix Γ_β . A distinction is then made between a restricted and an unrestricted specification. The restricted model imposes $\Gamma_\alpha = \mathbf{0}$ and assumes that characteristics affect expected returns only through risk exposures, which means there are no “anomaly” intercepts. In contrast, the unrestricted model sets $\Gamma_\alpha \neq \mathbf{0}$, with $\alpha_{i,t}$ capturing mean returns from characteristics that are not determined by risk exposures alone.

For the restricted model ($\Gamma_\alpha = \mathbf{0}$), Equation ?? can be rewritten in vector form as

$$r_{t+1} = Z_t \Gamma_\beta f_{t+1} + \epsilon_{t+1}^*, \quad (2.2)$$

where r_{t+1} is an $N \times 1$ vector of individual cryptocurrency returns, Z_t is the $N \times L$ matrix of stacked characteristics, and $\epsilon_{t+1}^* = \epsilon_{t+1} + \nu_{\alpha,t} + \nu_{\beta,t} f_{t+1}$ is a composite error vector stacking individual residuals. The estimation problem is to minimize the sum of squared composite model errors:

$$\min_{\Gamma_\beta, F} \sum_{t=1}^{T-1} (r_{t+1} - Z_t \Gamma_\beta f_{t+1})' (r_{t+1} - Z_t \Gamma_\beta f_{t+1})$$

The solution is obtained by alternating least squares, iterating the first-order conditions of f_{t+1} and Γ_β (?):

$$\hat{f}_{t+1} = (\hat{\Gamma}_\beta' Z_t' Z_t \hat{\Gamma}_\beta)^{-1} \hat{\Gamma}_\beta' Z_t' r_{t+1}, \quad \forall t \quad (2.3)$$

$$\text{vec}(\hat{\Gamma}_\beta) = \left(\sum_{t=1}^{T-1} Z_t' Z_t \otimes \hat{f}_{t+1} \hat{f}_{t+1}' \right)^{-1} \left(\sum_{t=1}^{T-1} [Z_t \otimes \hat{f}_{t+1}]' r_{t+1} \right) \quad (2.4)$$

In this sense, ALS alternates between estimating factor realizations via cross-sectional regressions on latent loadings (Equation ??) and updating Γ_β through regressions on factors interacted with characteristics (Equation ??).

Similarly, the unrestricted model ($\Gamma_\alpha \neq \mathbf{0}$) can be rewritten in vector form as

$$r_{t+1} = Z_t \tilde{\Gamma} \tilde{f}_{t+1} + \epsilon_{t+1}^*, \quad (2.5)$$

where $\tilde{\Gamma} = [\Gamma_\alpha, \Gamma_\beta]$ and $\tilde{f}_{t+1} = [1, f'_{t+1}]'$. Note that the unrestricted model simply augments the factor specification to include a constant. The first-order conditions slightly change to

$$f_{t+1} = \left(\Gamma'_\beta Z'_t Z_t \Gamma_\beta \right)^{-1} \Gamma'_\beta Z'_t (r_{t+1} - Z_t \Gamma_\alpha), \quad \forall t, \quad (2.6)$$

$$\text{vec}(\tilde{\Gamma}) = \left(\sum_{t=1}^{T-1} Z'_t Z_t \otimes \tilde{f}_{t+1} \tilde{f}'_{t+1} \right)^{-1} \left(\sum_{t=1}^{T-1} [Z_t \otimes \tilde{f}_{t+1}]' r_{t+1} \right) \quad (2.7)$$

In the unrestricted model, the intercept captures only the part of mean returns that is not already explained by factor loadings. In other words, it accounts for the residual variation in expected returns that characteristics cannot map into risk exposures.

2.1.1 Interpretation as a managed portfolio

As discussed in Kelly et al. (?), the asset pricing literature traditionally evaluates pricing factor performance using test portfolios, such as the value-sorted portfolios in the Fama-French data library¹, rather than individual assets. These portfolios reduce idiosyncratic variation by averaging across many securities. Kelly et al. (?) show that the IPCA framework provides an analogous representation through characteristic-managed portfolios. Each managed portfolio is constructed by a weighted average of asset returns, where the weights are given by their observable characteristics. For L asset-specific characteristics, the $L \times 1$ vector of managed portfolio returns is

$$x_{t+1} = \frac{1}{N_{t+1}} Z'_t r_{t+1},$$

where Z_t is the $N \times L$ matrix of characteristics at time t , r_{t+1} is the $N \times 1$ vector of realized asset returns, and N_{t+1} is the number of available assets.

Although the main focus of this thesis is on explaining the relationship between cryptocurrency returns and common risk factors using the panel of individual cryptocurren-

¹see https://mba.tuck.dartmouth.edu/pages/faculty/ken.french/data_library.html

cies, I also report results for models estimated with characteristic-managed portfolios.

2.2 Performance measures

Kelly et al. (?) propose two metrics to evaluate and compare the asset pricing performance of the IPCA model across different choices of K factors and between restricted and unrestricted specifications. These measures are referred to as the total R^2 and the predictive R^2 . However, since both statistics can take negative values —although “ R^2 ” suggests non-negative values—I refer to them here as simply the “total R” and the “predictive R”.

The total R measures the overall fit of the IPCA model by quantifying how much of the variation in realized returns can be explained by the estimated factors and conditional loadings. It is defined as

$$R_{\text{total}} = 1 - \frac{\sum_{i,t} \left(r_{i,t+1} - z'_{i,t} (\hat{\Gamma}_{\alpha} + \hat{\Gamma}_{\beta} \hat{f}_{t+1}) \right)^2}{\sum_{i,t} r_{i,t+1}^2}.$$

The predictive R measures how much of the variation in realized returns is explained by the model’s conditional expected returns, obtained by replacing realized factors with the estimated risk prices $\hat{\lambda}$. It is defined as

$$R_{\text{pred}} = 1 - \frac{\sum_{i,t} \left(r_{i,t+1} - z'_{i,t} (\hat{\Gamma}_{\alpha} + \hat{\Gamma}_{\beta} \hat{\lambda}) \right)^2}{\sum_{i,t} r_{i,t+1}^2},$$

In the restricted specification ($\Gamma_{\alpha} = 0$), the predictive R describes how well characteristics explain expected returns only through their effect on factor loadings, that is, through systematic risk exposures. In the unrestricted specification, the predictive R measures how well characteristics explain expected returns both through factor loadings and through anomaly intercepts.

2.3 Hypothesis tests

Kelly et al. (?) develop three hypothesis tests that help determine the whether one specification significantly improves the model description of asset returns.

Asset pricing test $\Gamma_\alpha = \mathbf{0}$

The first hypothesis test evaluates whether anomaly intercepts capture variation in returns beyond systematic risk exposures. In the unrestricted specification in Equation ??, expected returns are modeled as a linear function of both factor loadings and anomaly intercepts. The null hypothesis is

$$H_0 = \Gamma_\alpha = \mathbf{0}_{L \times 1}$$

against the alternative

$$H_1 = \Gamma_\alpha \neq \mathbf{0}_{L \times 1}$$

If the null is not rejected, characteristics influence expected returns only through factor loadings, and alphas are not associated with the characteristics in $z_{i,t}$. Rejecting the null indicates that characteristics help explain average returns directly through anomaly intercepts, in addition to their role in determining exposures to risk.

Following Kelly et al. (?), the null hypothesis is tested using a Wald-type statistic, which evaluates the distance between the restricted and unrestricted models as the sum of squared elements of the estimated Γ_α vector:

$$W_\alpha = \hat{\Gamma}'_\alpha \hat{\Gamma}_\alpha$$

Inference is carried out using a bootstrap procedure. After estimating the unrestricted model and retaining $\hat{\Gamma}_\alpha$, $\hat{\Gamma}_\beta$, and $\{\hat{f}_t\}_{t=1}^T$, the managed portfolio residuals are constructed as $d_{t+1} = Z'_t \epsilon_{t+1}^*$ from the managed portfolio definition

$$x_{t+1} = Z'_t r_{t+1} = (Z'_t Z_t) \Gamma_\alpha + (Z'_t Z_t) \Gamma_\beta f_{t+1} + Z'_t \epsilon_{t+1}^*$$

These residuals are resampled and the fitted values $\{\hat{d}_t\}_{t=1}^T$ stored. Then, for each bootstrap replication $b = 1, \dots, 1000$, a new sample of portfolio returns is generated as

$$\tilde{x}_{t+1}^b = (Z'_t Z_t) \hat{\Gamma}_\beta \hat{f}_{t+1} + \tilde{d}_{t+1}^b, \quad \tilde{d}_{t+1}^b = q_{1,t+1}^b \hat{d}_{q_{2,t+1}^b}$$

Here, $q_{2,t+1}^b$ is a random time index drawn uniformly from the set of all possible dates, and $q_{1,t+1}^b$ is a Student- t random variable with unit variance and five degrees of freedom. Using these bootstrap samples, the unrestricted model is re-estimated and the statistic recomputed as

$$\tilde{W}_\alpha^b = \tilde{\Gamma}_\alpha^{b'} \tilde{\Gamma}_\alpha^b$$

Finally, the empirical p -value is obtained as the fraction of bootstrap statistics \tilde{W}_α^b

2 Methodology

that exceed the observed value W_α from the actual data.

Testing instruments significance

This test evaluates whether a specific characteristic significantly contributes to factor loadings after controlling for all other characteristics. The analysis is based on the restricted model with $\Gamma_\alpha = 0$, where the goal is to assess whether the l^{th} characteristic helps explain the conditional loadings $\beta_{i,t}$. For this, first, the loading matrix is written as

$$\Gamma_\beta = [\gamma_{\beta,1}, \dots, \gamma_{\beta,L}]',$$

with $\gamma_{\beta,l}$ denoting the $K \times 1$ vector of coefficients linking characteristic l to the K latent factors. Under the null hypothesis, the l^{th} characteristic plays no role in determining exposures, so its entire row is set to zero:

$$H_0 : \Gamma_\beta = [\gamma_{\beta,1}, \dots, \gamma_{\beta,l-1}, \mathbf{0}_{K \times 1}, \gamma_{\beta,l+1}, \dots, \gamma_{\beta,L}]$$

against the alternative allowing for a non-zero contribution from characteristic l .

$$H_1 : \Gamma_\beta = [\gamma_{\beta,1}, \dots, \gamma_{\beta,L}]',$$

The Wald-type statistic used to evaluate this hypothesis is

$$W_{\beta,l} = \hat{\gamma}'_{\beta,l} \hat{\gamma}_{\beta,l}$$

Inference is based on the same residual bootstrap procedure as in the alpha test. One thousand bootstrap samples are generated under the null hypothesis that the l^{th} characteristic has no effect on factor loadings, the portfolio returns is re-estimated for each sample, and the corresponding statistics $\tilde{W}_{\beta,l}^b$ are computed. The p -value is obtained as the fraction of bootstrap statistics that exceed the observed $W_{\beta,l}$.

Testing pre-specified factors

In addition to estimating latent factors, the IPCA can nest pre-specified, common observable factors, to compare against a the general IPCA specification. Following Kelly et al. (?), the models can be implemented as (i) the traditional time-series approach with static loadings estimated asset-by-asset on the observable factors, and (ii) an instrumented version that keeps the factor returns fixed but parameterizes loadings as functions of characteristics. Therefore, the second specification is a combination

between pre-specifying observable factors in the IPCA model, and estimating its loadings dynamically, for each period t , or even generate latent factors additional to the pre-specified ones. The model is written as

In addition to estimating latent factors, IPCA can also incorporate pre-specified observable factors, allowing direct comparison with the general specification. Following Kelly et al. (?), two versions can be implemented: (i) the traditional time-series approach with static loadings estimated asset-by-asset on the observable factors, and (ii) an instrumented version that fixes the factor returns but models loadings as functions of characteristics. The latter combines pre-specified observable factors with the IPCA structure, since loadings are estimated dynamically each period t , while additional latent factors may also be generated alongside the pre-specified ones. The model takes the form

$$r_{i,t+1} = \beta_{i,t}f_{t+1} + \delta_{i,t}g_{t+1} + \epsilon_{i,t+1},$$

with

$$\delta_{i,t} = z'_{i,t}\Gamma_{\delta} + \nu_{\delta,i,t},$$

where the term $\delta_{i,t}g_{t+1}$ captures the contribution of the $M \times 1$ vector of observable factors g_{t+1} , and Γ_{δ} is the $L \times M$ mapping from characteristics to their loadings. Estimation proceeds as in the unrestricted case, but now with $\tilde{\Gamma} = [\Gamma_{\beta}, \Gamma_{\delta}]$ and $\tilde{f}t + 1 = [f't + 1, g'_{t+1}]'$. The first-order condition in Equation ?? remains the same, while Equation ?? becomes

$$f_{t+1} = \left(\Gamma'_{\beta}Z'_tZ_t\Gamma_{\beta}\right)^{-1}\Gamma'_{\beta}Z'_t(r_{t+1} - Z_t\Gamma_{\delta}g_{t+1}), \quad \forall t. \quad (2.8)$$

Kelly et al. (?) propose a test to assess the explanatory power of observable factors after controlling for the baseline IPCA specification. The null hypothesis states that observable factors add no additional explanatory power

$$H_0 : \Gamma_{\delta} = \mathbf{0}_{L \times M},$$

against the alternative

$$H_1 : \Gamma_{\delta} \neq \mathbf{0}_{L \times M}$$

The Wald-type statistic used to evaluate this hypothesis is

$$W_{\delta} = \text{vec}(\hat{\Gamma}_{\delta})'\text{vec}(\hat{\Gamma}_{\delta})$$

which measures the distance between the specification that includes observable factors and the restricted model that excludes them. A large W_{δ} suggests that observable

2 Methodology

factors provide incremental explanatory power for asset returns after accounting for the latent IPCA factors. Inference is based on the same residual bootstrap procedure as in the previous tests, using $b = 1, \dots, 1000$ bootstrap samples.

3 Data

In this section, I introduce the cryptocurrency data used in this thesis, and describe the series of filters applied to clean and prepare the dataset, and the summary statistics of the cryptocurrency excess returns. In addition, I present the set of asset-specific characteristics constructed from the cryptocurrency market data, which are used as instruments for latent factor exposures in the IPCA model. Finally, I construct a set of observable risk factors, or factor-mimicking portfolios, which are used as pre-specified factors in the analysis. Appendix ?? and -Section ?? provides a detailed description of the set of characteristics and factors, respectively.

The data extraction and pre-processing are primarily conducted in R 4.5.1 (?), using, among other packages¹, the `tidyverse` (v. 2.0.0; ?). Additional cleaning steps and visualizations are performed in Python 3.13.5 (?). The full reproducible code is available in Appendix -Section ??.

3.1 Data extraction and sample construction

I collect daily cryptocurrency data on open, high, close, and low (OHCL) prices, 24-hour volume, and market capitalization (calculated as the cryptocurrency’s USD price multiplied by its circulating supply) from [CoinCodex](#), a website-data provider that gathers and aggregates data from more than 400 exchanges. I extract the data, all expressed in US dollars, using the CoinCodex API as follows:

1. I retrieve the list of all available cryptocurrencies and extract each cryptocurrency shortname, also referred to as the “slug”. At the time of writing, there are 14,907 unique cryptocurrency shortnames listed in the API.
2. Using the slug, I construct an URL for each cryptocurrency to obtain the meta-data from the API. I parse the JSON API response into a dataframe and extract the OHCL prices, volume, and market capitalization daily data. I exclude those observations with non-zero or missing values in any of these fields.

¹See Appendix ?? for the full list of software used in the empirical study.

3 Data

Out of the 14,907 cryptocurrencies listed, only 7,272 entries contained available data. Next, following the methodology of Bianchi & Babiak (?) and Mercik et al. (?), I apply a series of cleaning and filtering steps in order to remove possible inaccuracies in the dataset:

1. Non-positive and missing values. As mentioned earlier, I remove observations where prices, volume, or market capitalization were non-positive or missing.
2. Small cryptocurrencies. Similar to Y. Liu et al. (?), I screen out small cryptocurrencies and consider only those with a market capitalization greater than one million USD. Therefore, I exclude observations for coins whose market capitalization falls below this minimum threshold, which allows for the possibility that a coin may become “small” after a certain period or event.
3. Cryptocurrency type. Based on the cryptocurrency classification from [CoinMarketCap](#) and CoinCodex, I exclude:
 - stablecoins. I include (i) centralized stablecoins, which are backed and pegged to fiat currency or physical assets by a third party, such as Tether (USDT), USD Coin (USDC), and Euro Coin (EURC), and (ii) algorithmically stabilized stablecoins, which use algorithms to adjust the circulating supply in response to changes in demand to maintain a stable value with the underlying asset, such as DAI and AMPL (FSB, ?).
 - wrapped cryptocurrency tokens, which mirror the value of another cryptocurrency from a different blockchain, e.g., Wrapped Bitcoin (wBTC) or Wrapped Ethereum (wETH) (?).
 - cryptocurrencies backed by or pegged to gold or precious metals, including Pax Gold (PAXG) or XAGx Silver Token (XAGX).
4. Erroneous trading volume. To filter out cryptocurrencies with “fake” or “erroneous” trading volume, I calculate the daily volume-to-market-capitalization ratio for each token and exclude observations where the ratio exceeds 1.
5. Extreme returns. To minimize the influence of extreme values in my results, I winsorize daily cryptocurrency returns to lie within the range of -90% to 500%.
6. Time period. Even though cryptocurrency data are available since 2014, I use data from June 1, 2018 for the empirical analysis due to the low amount of coins available before this date (see Figure ??).
7. Minimum observations. In order to maintain practical relevance, I keep crip-

to currencies that have at least 365 consecutive daily observations and those with at least 730 observations in the complete panel of coin characteristics (see Section ??), which is equivalent to 2 years of historical data. Therefore, I exclude very short-lived coins, but retain failed coins with this relatively large number of observations, which help to lessen the so called “survivorship bias”.

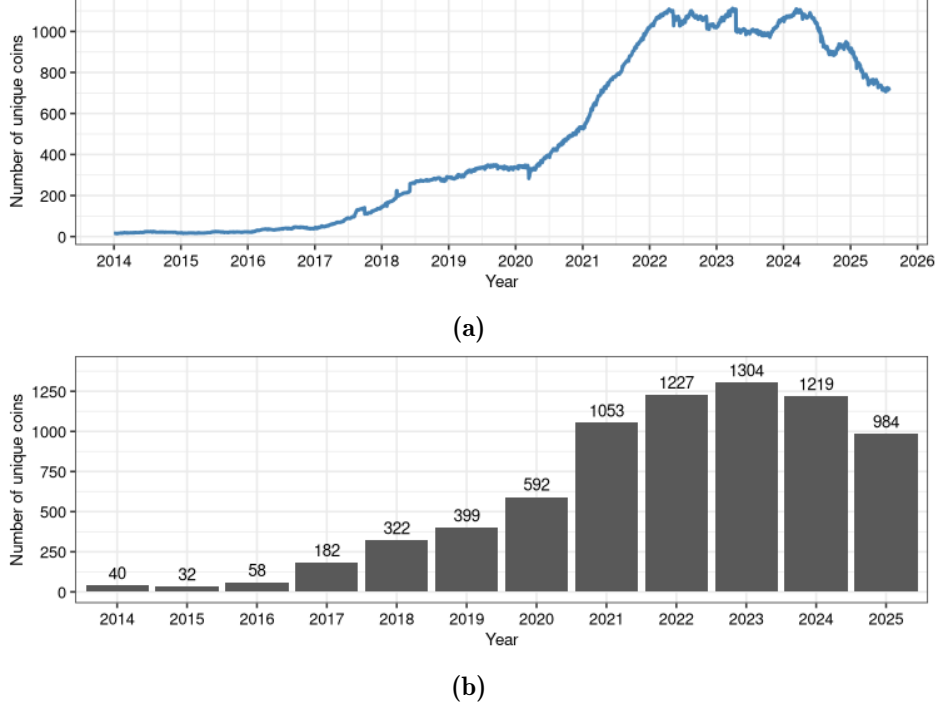


Figure 3.1: Number of cryptocurrencies over time. Panel A shows the daily time series of the number of unique cryptocurrencies. Panel B displays the number of unique cryptocurrencies recorded each year. Both panels correspond to the dataset after applying the filtering steps (1) to (5), covering the period from January 1, 2014, to July 31, 2025, and including 1,416 unique cryptocurrencies. Note that coins may enter or exit the market over time.

3.2 Sample overview

After applying all the filters, the resulting sample consists of 973 unique cryptocurrencies and 1,478,936 observations from June 1, 2018, to July 31, 2025, where a day starts at 00:00:00 UTC. It is important to mention that the number of cryptocurrencies fluctuates over the entire period, which results in an unbalanced panel of data. Table ?? provides a description of the yearly cross-sectional statistics: the sample starts with 254 different cryptocurrencies in 2018 and peaks in 2023 with 939 unique cryptocurrencies, before decreasing to 780 in 2025. The minimum daily cross-section is 121 in 2018, and then increases drastically up to 793 in 2023. For context, at the time of writing, CoinMarketCap tracks around 19 million cryptocurrencies, and CoinGecko around 19

3 Data

Table 3.1: Cross-section size of the sample. The table reports the number of unique coins per year, as well as the minimum daily cross-section size in the filtered sample.

Year	2018	2019	2020	2021	2022	2023	2024	2025
Unique coins	254	337	420	714	938	939	906	780
Min. daily cross-section	121	239	207	381	699	793	710	578

Table 3.2: Summary statistics of daily returns. The table reports summary statistics of daily returns for the filtered sample, the top 100 and top 10 cryptocurrencies ranked by market capitalization, and for Bitcoin, Ethereum, and Ripple individually. Reported statistics include the number of daily observations, the number of unique coins over the sample period, the mean and standard deviation of returns, and the 10th percentile, lower quartile, median, upper quartile, and 90th percentile of the distribution of the returns. The sample period is from June 1, 2018, to July 31, 2025.

==

thousands. When compared to these numbers, the size of the sample may seem small; however, it actually covers most of the whole cryptocurrency market capitalization (see Figure ??). The sample period includes important events in the market, such as

Table ?? summarizes the descriptive statistics for the cryptocurrency daily returns across different subsamples and Bitcoin, Ethereum, and Ripple, which are the three largest cryptocurrencies in the sample. Interestingly, the larger samples exhibit a larger volatility and more pronounced extreme returns, both positive and negative. Bitcoin shows the lowest mean return during the sample period (0.16% per day), though this value very close to that of Ethereum (0.17%) and Ripple (0.20%), and only slightly below other cryptocurrency subsamples.

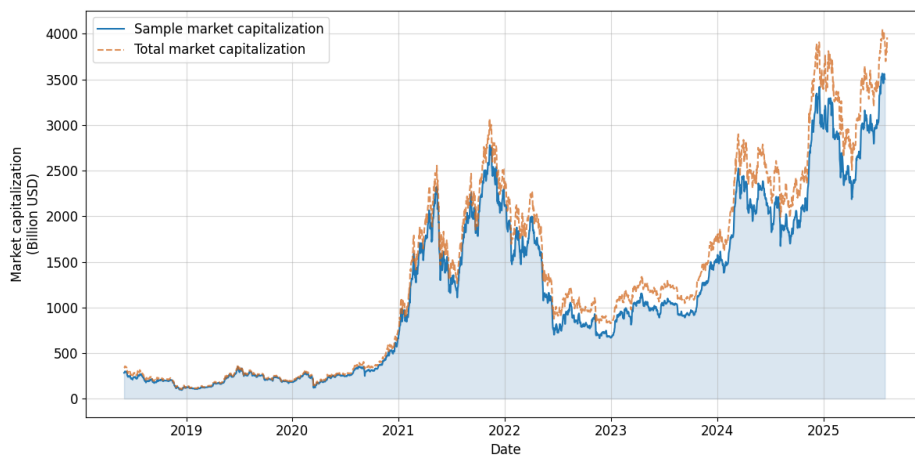


Figure 3.2: Cryptocurrency market capitalization. The figure compares the cryptocurrency market capitalization in the filtered sample (blue line) with the total market capitalization (yellow line) from June 1, 2018 to July 31, 2025. Source: total market capitalization from [CoinGecko](#).

The sample period spans several major market, economic, and political events, these

include: the start of the COVID-19 pandemic and the subsequent crypto bubble in 2020-2021, El Salvador adoption of Bitcoin as legal tender in September 2021, and China’s ban on cryptocurrency exchanges and mining in October 2021. The period also experienced multiple cryptocurrency exchange hacks², and geopolitical shocks such as the Russia-Ukraine war in February 2022, and the Palestine-Israel war in October 2023. More recently, in 2024, the U.S. Securities and Exchange Commission (SEC) approved the listing and trading of several crypto spot ETFs in January, and Donald Trump’s election as U.S. president, with Elon Musk playing an important role in his campaign (??; ??; ??; ??; ??).

3.3 Characteristic construction and description

For the analysis, I construct 41 asset-specific characteristics from the cross-section of 973 cryptocurrencies using data on prices, volume, and market capitalization. Specifically, I follow the methodology of Bianchi & Babiak (??), Y. Liu et al. (??), and Mercik et al. (??) to construct the set of characteristics widely used in the cryptocurrency and financial literature, which serve as return predictors in the empirical analysis. These characteristics are grouped into six categories: market and size, volatility and risk, trading activity, liquidity, past returns, and distribution. Table ?? summarizes the set of characteristics, while Appendix ?? provides detailed definitions and construction procedures.

3.4 Observable risk factors

In addition to the set of characteristics described above, I construct a set of observable risk factors. In the asset pricing literature, the convention is to analyze the risk compensation of asset returns using factor-mimicking portfolio (e.g. ??; ??, ?). This typically involves sorting assets cross-sectionally into quintiles based on a specific characteristic and forming a factor return, calculated as the difference in returns between the top and the bottom quintiles. This approach replicates a strategy that buys the portfolio of assets with high values of a particular characteristic (long), and sells the portfolio with the lowest values (short).

Building on this methodology, I construct a series of observable risk factors that prior literature have shown to explain the cross-section of cryptocurrency returns. Specifi-

²For example, Binance, largest crypto exchange in the world, was hacked in 2019, and KuCoin and Crypto.com were hacked in 2020 and 2022, respectively. (??)

3 Data

Table 3.3: Cryptocurrency characteristics. The table presents the 41 cryptocurrency characteristics used as return predictors in the empirical analysis. The characteristics are grouped in six categories: price and size, volatility and risk, trading activity, liquidity, past returns, and distribution.

No.	Characteristic	Symbol	Definition
Panel A: Price & size			
NA	NA	mcap	NA
(2)	Price	prc	Last day's logged closing price.
(3)	Closeness to the 90-day high	dh90	Last day's price over the maximum price in the previous 90 days.
Panel B: Volatility & risk			
(4)	Market beta	beta	CAPM market beta, estimated from 30 days of daily returns.
(5)	Idiosyncratic volatility	ivol	Volatility of CAPM residuals over 30 days of daily returns.
(6-7)	Realized volatility	rvol_*d	Realized volatility, calculated from 7 and 30 days of OHCL prices.
(8)	Return volatility	retvol	Standard deviation of daily returns over 7 days.
(9)	Value-at-Risk	var	The historical Value-at-Risk at 5% level over 90 days.
(10)	Expected Shortfall	es_5	The expected shortfall at the 5% level over 90 days.
(11)	Price delay	delay	Improvement in R^2 after adding lagged one-and two-day market excess return to the CAPM.
Panel C: Trading activity			
(12)	Trading volume	volume	Last day's daily trading volume in US dollars.
(13)	Average volume	volume_*d	Mean volume over the past 7 and 30 days.
(15)	Turnover	turn	The last day's trading volume over current market capitalization.
(16)	Average 7-day turnover	turn_7d	Mean turnover over the past 7 days.
(17)	Turnover volatility	std_turn	Turnover volatility over the past 30 days.
(18)	Trading volume volatility	std_vol	Volume's logged volatility over the past 30 days.
(19)	Volume's coefficient of variation	cv_vol	Volume's volatility over its mean in the previous 30 days.
Panel D: Liquidity			
(20)	Bid-ask spread	bidask	Mean estimated bid-ask spread calculated over the past 30 days.
(21)	Illiquidity	illiq	Mean absolute daily return over trading volume over the past 30 days.
(22)	Standardized abnormal turnover	sat	Last day's turnover minus its 30-day average, divided its volatility over 30 days.
(23)	De-trended turnover	dto	De-trended turnover minus the value-weighted daily market turnover.
(24)	Volume Shock 15-day	volsh_15d	Log deviation of trading volume from its rolling 15-day average.
(25)	Volume Shock 30-day	volsh_30d	Log deviation of trading volume from its rolling 30-day average.
Panel E: Past returns			
(26)	Daily reversal	r2_1	Return on the previous trading day.
(27-30)	Momentum	r*_1	7, 14, 21, and 30-day cumulative return ending 1 day before the prediction date.
(31)	Intermediate momentum	r30_14	Cumulative return from 30 to 14 days before the prediction date.
(32)	Long-term reversal	r180_60	Cumulative return from 180 to 60 days before the prediction date.
(33)	CAPM alpha	alpha	CAPM intercept, estimated from 30 days of daily returns.
Panel F: Distribution			
(34-35)	Skewness	skew_*d	Skewness of the daily return distribution over a 7-and 30-day period.
(36-37)	Kurtosis	kurt_*d	Kurtosis of the daily return distribution over a 7-and 30-day period.
(38-39)	Maximum daily return	maxret_*d	The maximum daily return in the past 7-and 30 days.
(40-41)	Minimum daily return	minret_*d	The minimum daily return in the past 7-and 30 days.

cally, I include the market, size, momentum, liquidity, and volatility factors, following Y. Liu et al. (?), Bianchi & Babiak (?), and Lan & Frömmel (?). Details on their construction are provided in Appendix ???. As described in Section ??, the IPCA allows for the inclusion of pre-specified factors within the more general model specification. I make use of this feature and pre-specify the observable factors in the IPCA model, with and without using asset-characteristics to instrument for dynamic loadings.