

PCA based construction of cryptocurrency index

Process Report

Zeyan Huang

April 30, 2025

Abstract

1 Introduction

After reviewing the paper "Principal Component Analysis-Based Construction and Evaluation of a Cryptocurrency Index"[1], I decided to replicate its results and explore potential improvements for analyzing the cryptocurrency market.

2 Data

To construct the dataset for this study, I collected daily market capitalization data for the top 1,200 cryptocurrencies from CoinGecko.com, covering the period from March 29, 2024, to March 28, 2025. The selection of these cryptocurrencies was based on their rankings as of March 28, 2025, ensuring that the dataset reflects the market structure as observed at the end of the analysis period.

One notable issue with the dataset is the presence of missing data points. Specifically, some cryptocurrencies had not yet launched at the start of the data period, resulting in incomplete time series that could complicate further analysis.

To address this issue, I adopted different approaches for handling missing data when calculating the static index and the dynamic index.

3 Static Index

The methodology presented in the paper uses the first principal component (PC1) as a set of weights, which are multiplied by the market capitalizations of the corresponding cryptocurrencies on a given day. The weighted values are then summed to compute the index value for each specific day.

In this replication, I applied principal component analysis to the full dataset—consisting of 1,200 cryptocurrencies over 365 days—and used PC1 as the weighting scheme to construct the daily index. However, the variance explained by PC1 was only 36%, which is relatively low and suggests that the resulting index may not be as reliable or informative as implied in

the original paper. Here I used the sum of the 1200 cryptocurrencies' market capitalization as the total market capitalization.

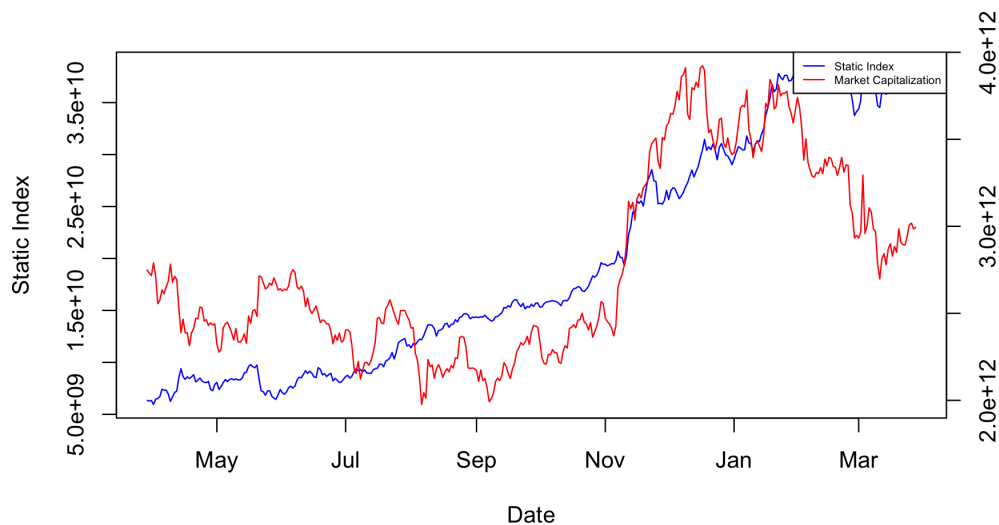


Figure 1: Static Index using PC1

As shown in the plot above, the calculated static index does not fully capture the overall movement of the cryptocurrency market. However, there are several periods where the index exhibits similar short-term trends to the total market capitalization, suggesting that the method may still partially reflect market dynamics.

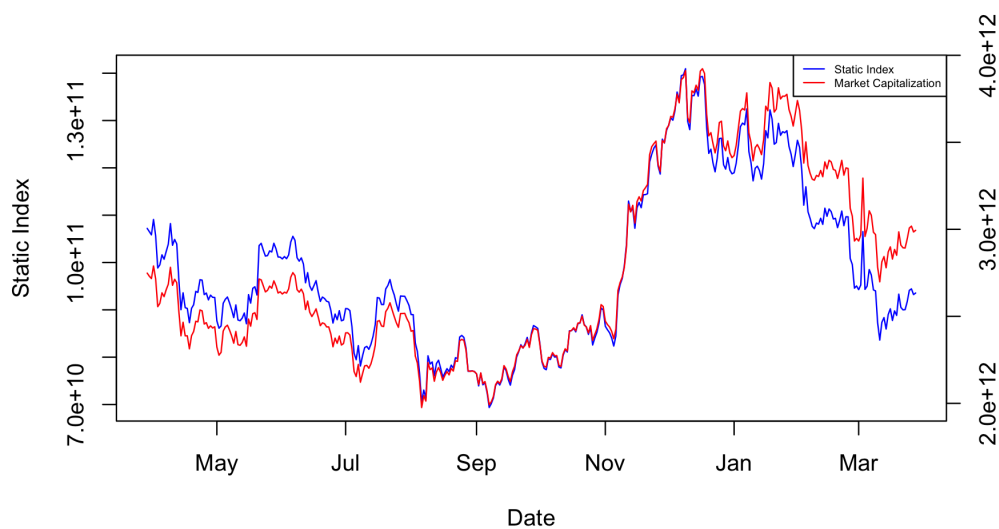


Figure 2: Static Index using absolute value of PC1 and PC2

Interestingly, when using the absolute value of the sum of PC1 and PC2 loadings as weights, the resulting static index appears to effectively capture the overall movement of the cryptocurrency market, even though the combined variance explained is only 61.5%.

The static index without using absolute value provided below.

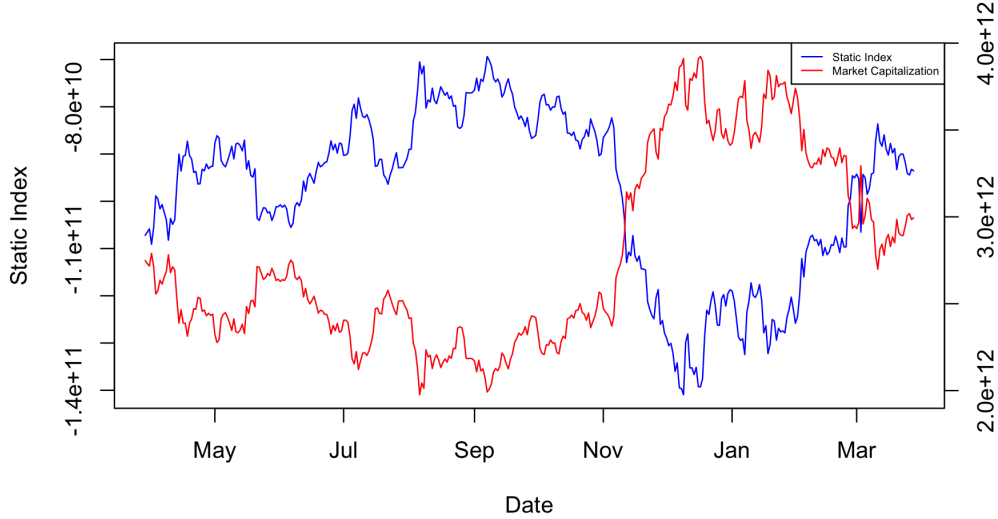


Figure 3: Static Index using PC2

Question!!! Why the absolute method here work well?

4 Dynamic Index Construction Methodology

Based on the PCA-based dynamic index methodology presented in the original paper, the calculation of the dynamic index consists of two main steps.

In the **first step**, a rolling window of T_n days (90 days in our implementation) is used to conduct PCA and correlation analysis to determine the number of constituent cryptocurrencies to include in the index. We define an update interval of T_{n-gap} days (30 days in our test). Every 30 days, we extract the preceding 90-day window of data, apply PCA, and identify the subset of cryptocurrencies whose individual indices (constructed using PC1) have a correlation higher than 0.99 with the index derived from the full set of cryptocurrencies. This subset of cryptocurrencies is saved for use in the next step, and the number of selected assets is denoted as n_c .

In the **second step**, the index is computed using a forward-looking window. Specifically, for each update period, we take the previous T_w days (also 30 in our implementation) of data for the n_c selected cryptocurrencies and apply PCA. The first principal component (PC1) is used as the weighting vector. The index value a_t for the subsequent T_w days is calculated as:

$$a_t = \sum_{i=1}^{n_c} PC_{1,i} \cdot M_{t,i}$$

where $PC_{1,i}$ is the loading of the i -th asset on the first principal component, and $M_{t,i}$ is the market capitalization of asset i on day t .

Once the daily values of a_t are obtained, the dynamic PCA-based index I_{PCA} is computed as:

$$I_{PCA} = \frac{a_t}{a_{base}} \times m$$

Here, m is a multiplier (set to 1000 for this test), and a_{base} is the value of a_t on the base day t_0 , which is the first day of each new T_w -day sub-period. This ensures that the index resets to m at the beginning of each sub-period, maintaining consistency with the methodology described in the paper.

The plot below compares the dynamic PCA-based index with the total market capitalization, using the first principal component (PC1) as the weighting vector. While the dynamic index generally tracks the overall movement of the market quite well, a noticeable divergence appears after December, where the index underperforms relative to the market capitalization. I believed this gap is because our m value is set as 1000 and not updated follows the market, so the next step should focus on the calculation methodology of choosing m .

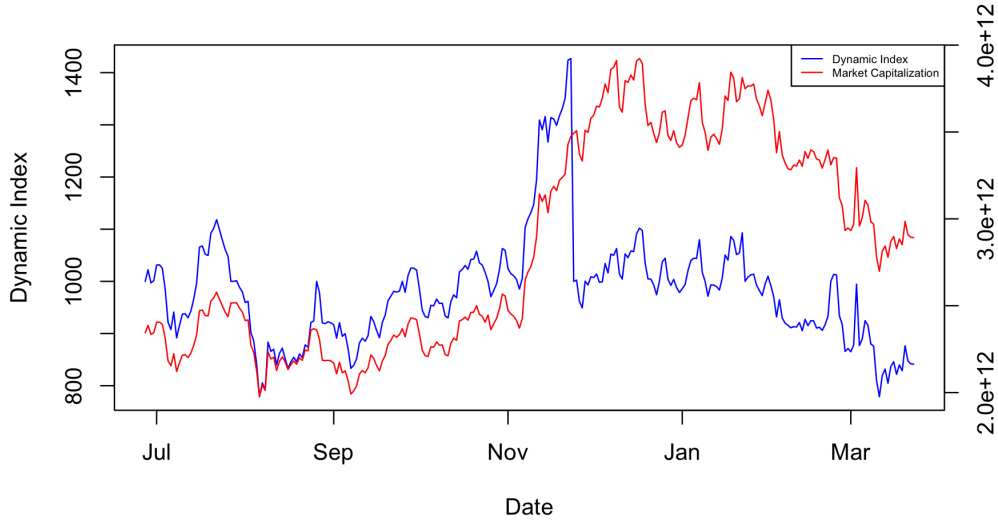


Figure 4: Dynamic index using PC1

References

- [1] Principal component analysis based construction and evaluation of cryptocurrency index.
Retrieved from <https://doi.org/10.1016/j.eswa.2020.113796>

Appendix

R Code for Data Preprocessing

Listing 1: Data Preprocessing

```
1 # Data
2 # Data preparation
3 # data from Coingecko
4 folder_path <- "market_cap_data/"
5 all_files <- list.files(folder_path)
6
7
8 dataset <- lapply(all_files, function(x){
9   data <- read.csv(paste0(folder_path, x), header = TRUE)
10  data <- data[-nrow(data), ]
11  data$Date <- as.Date(data$date)
12  data <- data[order(data$Date),]
13  data <- data[!duplicated(data$Date), ]
14
15
16  price <- as.numeric(data$market_cap)
17  names(price) <- data$Date
18  return(price)
19 }
20 )
21
22 # Name the data list
23 file_names <- gsub("\\.csv", "", all_files)
24 names(dataset) <- file_names
25
26 # Step 1: Convert each price series to a dataframe with Date + market cap
27 dataset_aligned <- lapply(names(dataset), function(name) {
28   df <- data.frame(Date = as.Date(names(dataset)[[name]]),
29                     Cap = dataset[[name]])
30   colnames(df)[2] <- name
31   return(df)
32 })
33
34 # Step 2: Merge all dataframes by Date, using full outer join
35 data <- Reduce(function(x, y) merge(x, y, by = "Date", all = TRUE),
36               dataset_aligned)
37
38 # Step 3: set all NA value as 0
39 # data[is.na(data)] <- 0
40
41 # Set Date as row index
42 rownames(data) <- data$Date
43 data$Date <- NULL
```

R Code for Static Index

Listing 2: SI

```

1 # Static Index
2 # Data prep
3 data_static_index <- data
4
5 data_static_index[is.na(data_static_index)] <- 0
6
7 (pca_static_index$sdev[1:7])^2
8 sum((pca_static_index$sdev)^2)
9 sum((pca_static_index$sdev[1:2])^2)/sum((pca_static_index$sdev)^2)*100
10
11 static_index <- as.matrix(data_static_index) %*% as.vector(pca_static_
    index$rotation[,1])
12
13 static_index <- as.data.frame(static_index)
14
15 static_index$Date <- as.Date(rownames(static_index))
16
17 plot(static_index$Date, static_index$V1, type = "l", xlab = "Date", ylab =
    "PCA Static Index")
18
19 # Get the daily market capitalization
20 data_mc <- data.frame(RowSum = rowSums(data_static_index))
21
22 rownames(data_mc) <- rownames(data_static_index)
23
24 # Combine the data_mc with static_index
25 data_mc_si <- merge(static_index, data_mc, by = "row.names")
26 rownames(data_mc_si) <- data_mc_si$Row.names
27 data_mc_si$Row.names <- NULL
28
29 (data_mc_si$Date, data_mc_si$V1, type = "l", col = "blue", xlab = "Date",
    ylab = "Static Index")
30
31 par(new = TRUE)
32 plot(data_mc_si$Date, data_mc_si$RowSum, type = "l", col = "red", axes =
    FALSE, xlab = "", ylab = "")
33 axis(side = 4) # Add right-side axis
34 mtext("Market Capitalization", side = 4, line = 3) # Label for right Y-
    axis
35
36 # Add legend
37 legend("topright", legend = c("Static Index", "Market Capitalization"),
38       col = c("blue", "red"), lty = 1, cex = 0.5)
39
40 static_index_pc2 <- abs(as.matrix(data_static_index) \%*\% as.vector(
    rowSums(pca_static_index$rotation[,1:2])))
41
42 static_index_pc2 <- as.data.frame(static_index_pc2 )
43
44 static_index_pc2$Date <- as.Date(rownames(static_index_pc2))
45
46 # Combine the data_mc with static_index
47 data_mc_si_pc2 <- merge(static_index_pc2, data_mc, by = "row.names")
48 rownames(data_mc_si_pc2) <- data_mc_si_pc2$Row.names

```

```

49 data_mc_si_pc2$Row.names <- NULL
50
51 plot(data_mc_si_pc2$Date, data_mc_si_pc2$V1, type = "l", col = "blue",
      xlab = "Date", ylab = "Static_Index")
52
53 par(new = TRUE)
54 plot(data_mc_si_pc2$Date, data_mc_si_pc2$RowSum, type = "l", col = "red",
      axes = FALSE, xlab = "", ylab = "")
55 axis(side = 4) # Add right-side axis
56 mtext("Market_Capitalization", side = 4, line = 3) # Label for right Y-
      axis
57
58 # Add legend
59 legend("topright", legend = c("Static_Index", "Market_Capitalization"),
60       col = c("blue", "red"), lty = 1, cex = 0.5)

```

R Code for Nc constituents Update function

Listing 3: NC

```

1 Nc_cal <- function(dataframe, threshold = 0.99){
2   # Step 1: Calculate PC1 for all columns
3   pca_full_temp <- prcomp(dataframe, scale. = TRUE)
4
5   # Initialize a list to store PC1 for each subset of columns
6   pc1_list_temp <- list()
7   index_list_temp <- list()
8   variance_explained_list <- list()
9
10  # Loop through subsets of columns (from 1:2, 1:3, ..., 1:ncol)
11  for (i in 2:ncol(dataframe)) {
12    subset_data <- dataframe[, 1:i]
13    pca_subset <- prcomp(subset_data, scale. = TRUE)
14    pc1_subset <- pca_subset$rotation[, 1]
15
16    pc_1_v <- (pca_subset$sdev[1])^2
17    pc_all_v <- sum((pca_subset$sdev)^2)
18    pc_1_explained_v_pct <- pc_1_v/pc_all_v
19    variance_explained_list[[i]] <- list(pc_1_v, pc_all_v, pc_1_explained_
      v_pct)
20
21    temp_index <- as.matrix(subset_data) %*% as.vector(pc1_subset)
22    pc1_list_temp[[i]] <- pc1_subset # Store in list (index starts at 1)
23    index_list_temp[[i]] <- temp_index
24  }
25
26  # Initialize a vector to store the correlations
27  correlations_temp <- numeric(length(pc1_list_temp) - 1)
28
29  # Index with all crypto
30  index_all_temp <- index_list_temp[[length(index_list_temp)]]
31
32  # Calculate the Correlation
33  for (i in 2: (length(index_list_temp))){

```



```

34     correlations_temp[i] <- cor(index_list_temp[[i]], index_all_temp)
35 }
36
37 # Find the first index where correlation >= threshold
38 first_reach_index_temp <- which(correlations_temp >= threshold)[1]
39 chosen_crypto_list <- c(colnames(dataframe)[1:first_reach_index_temp])
40
41 output_list <- list()
42 output_list[['Index']] <- index_list_temp
43 output_list[['Correlation']] <- correlations_temp
44 output_list[['Nc']] <- first_reach_index_temp
45 output_list[['Chosen_Crypto']] <- chosen_crypto_list
46 output_list[['Variance']] <- variance_explained_list
47 return(output_list)
48 }

```

R Code for Basic dynamic index calculation function for single period

Listing 4: DI_{base}

```

1 dynamic_index_base <- function(dataframe, Tw = 30){
2   a_list <- c()
3   a_base <- as.numeric(dataframe[Tw+1,])%%as.vector(prcomp(dataframe[1:Tw
4     ],, scale.= TRUE)$rotation[,1])
5   a_list <- append(a_list, 1)
6   # for (i in 2:Tw){
7     #   a_t <- (as.numeric(dataframe[Tw+i,])%%as.vector(prcomp(dataframe[i
8       : (Tw+i-1),], scale.= TRUE)$rotation[,1]))/a_base
9     #   a_list <- c(a_list, a_t)
10    # }
11    # a_t_df <- data.frame(Value = a_list, row.names = rownames(dataframe[(
12      Tw+1):(Tw+Tw),]))
13    a_t <- as.matrix(dataframe[(Tw+2):(Tw+Tw),])%%as.vector(prcomp(
14      dataframe[1:Tw,], scale.= TRUE)$rotation[,1])
15
16    a_list <- c(a_list, as.vector(a_t)/c(a_base))
17    a_t_df <- data.frame(Value = a_list, row.names = rownames(dataframe[(Tw
18      +1):(Tw+Tw),]))
19
20    return(a_t_df)
21 }

```

R Code for Dynamic Index function

Listing 5: DI_{base}

```

1 dynamic_index_cal <- function(dataframe, Tn = 90, Tn_gap = 30, Tw = 30,
2   threshold = 0.99){
3   # Nc Choosing process
4   nc_result <- list()

```

```

5  n_iter <- floor((nrow(dataframe) - Tn) / Tn_gap)
6  for (i in 0:n_iter){
7    data_90 <- dataframe %>%
8    slice((1+i*30):(90+i*30)) %>%          # Select the first 90 rows (time
      window)
9    select(where(~ all(!is.na(.)))) %>% # Drop columns that contain any NA
      values
10   select(where(~ all(. != 0))) %>%      # Drop columns that contain any 0
      values
11   select(where(~ {                      # Drop columns with 0 or NA
      standard deviation
12     sd_x <- sd(., na.rm = TRUE)        # Calculate standard deviation
      ignoring NA
13     is.finite(sd_x) && sd_x > 0        # Keep only columns with valid,
      non-zero SD
14   })) %>%
15   { .[, order(-as.numeric(.[1, ]))] } # Sort the remaining columns by
      the first row, descending
16
17   nc_result[[i+1]] <- Nc_cal(data_90, threshold = threshold)
18 }
19
20 # Dynamic Index Calculation
21 dynamic_index_results <- list()
22 for (i in 1:length(nc_result)){
23   data_60 <- dataframe %>%
24   slice((31+i*30):(90+i*30)) %>%          # Select the required 60 rows (
      time window)
25   select(where(~ all(!is.na(.)))) %>% # Drop columns that contain any NA
      values
26   select(where(~ all(. != 0))) %>%
27   select(all_of(nc_result[[i]][["Choosen_Crypto"]])) # Filter by selected
      crypto
28   if (nrow(data_60) < 60){
29     break
30   }
31   dynamic_index_results[[i]] <- dynamic_index_base(data_60)
32 }
33
34 final_results <- list(nc_result, dynamic_index_results)
35 return(final_results)
36 }

```