# Can Large Language Models Mine Interpretable Financial Factors More Effectively? A Neural-Symbolic Factor Mining Agent Model

Anonymous ACL submission

#### Abstract

Finding interpretable factors for stock returns is the most vital issue in the empirical asset pricing domain. As data-driven methods, existing factor mining models can be categorized into symbol-based and neural-based models. Symbol-based models are interpretable but inefficient, while neural-based approaches are efficient but lack interpretability. Hence, mining interpretable factors effectively presents a significant challenge. Inspired by the success of Large Language Models (LLMs) in 011 various tasks, we propose a FActor Mining Agent (FAMA) model that enables LLMs to integrate the strengths of both neural and sym-014 015 bolic models for factor mining. In this paper, FAMA consists of two main components: 017 Cross-Sample Selection (CSS) and Chain-of-Experience (CoE). CSS addresses the homogeneity challenges in LLMs during factor mining by assimilating diverse factors as in-context samples, whereas CoE enables LLMs to leverage past successful mining experiences, expediting the mining of effective factors. Experimental evaluations on real-world stock market data demonstrate the effectiveness of our approach by surpassing the SOTA RankIC by 0.006 and RankICIR by 0.105 in predicting 027 S&P 500 returns. Furthermore, the investment simulation shows that our model can achieve superior performance with an annualized return of 39.0% and a Sharpe ratio of 667.6%.

#### 1 Introduction

The task of predicting market trends in finance presents a formidable challenge, given the intricate interplay of various factors (Hou et al., 2011), such as the dynamics of demand and supply (Hendricks and Singhal, 2009), market sentiment (Verma and Soydemir, 2009) and government regulations (Ali Imran et al., 2020). In the field of quantitative trading, the conventional approaches often extract factors as indicative signals for market trends from raw historical stock data first, then



Figure 1: An illustration of three distinct factor mining approaches: (a) symbolic factor model, (b) neural factor model, and (c) our proposed neural-symbolic model.

serve them as input features for machine learning models (Sharpe, 1964; Ross, 2013; Duan et al., 2022). A pivotal step in this process entails discerning and extracting effective factors that demonstrate robust predictive capabilities for market trends (Ng et al., 1992). As an illustrative example, the Capital Asset Pricing Model (CAPM) (Sharpe, 1964) employed the market's excess return as a predictive factor for the return of a financial asset, thereby contributing a seminal factor to finance.

Hence discovering factors with high returns has been a trendy topic among investors and researchers. The prevailing methods for mining factors can be in general divided into two groups, namely symbolic factor and neural factor models. As illustrated in Figure 1(a), in symbolic factor models, factors are represented as symbolic expressions, then symbolic regression (Makke and Chawla, 2022) serves as a common technique

for factor mining (Jin et al., 2019; Zhang et al., 062 2020; Chen et al., 2021; Cui et al., 2021). For 063 instance, considering two factors, Factor1 =064 close/open and Factor2 = log(close), the factor values are calculated by the opening and closing price, then the two factors are inputted into 067 a symbolic regression model to generate a novel factor,  $NewFactor = \log(close/open)$ . The interpretability of the symbolic factor model arises from the explicit representation of the calculation process for the factors. However, due to the vast search 072 space of symbolic factors, mining with symbolic factor models often proves inefficient. Conversely, neural factor approaches, gaining recent popularity, transform factors into numerical features to opti-076 mize factor extraction. As depicted in Figure 1(b), neural factor models predict market trends by extracting numerical factor features from stock data through feature extractors (Kelly et al., 2019; Gu et al., 2021; Duan et al., 2022). Compared with symbolic factor models, neural factor models exhibit proficiency in extracting effective numerical factors. However, the financial interpretability in neural factor models struggles with implicit features. The question we are facing is: Can an effective approach be devised for mining financially interpretable factors conducive to predicting market trends?

Recent advancements in LLMs have demonstrated success across financial tasks, including sentiment analysis (Guo et al., 2023) and financial text generation (Yang et al., 2023). Thanks to its powerful In Context Learning (ICL) ability (Brown et al., 2020), we conceptualize LLMs as a neuro-symbolic model illustrated in As depicted in Figure 1(c), that bridges two distinct representations, i.e. numerical and symbolic factors, aiming to achieve efficient mining of interpretable ones. It is facile to consider utilizing the factors disclosed (Kakushadze, 2016) as contextual examples to generate new factors through In-context Learning. Since the disclosed factors are often limited in number, high correlation, and low complexity, direct mining factors using ICL encounter challenges. These issues can be summarized in two aspects: (1) The heightened homogeneity observed among factors, characterized by the uniform structure, culminates in the generation of the singular factor form through ICL. (2) The presence of a noteworthy proportion of ineffective factors acts as an impediment, hindering ICL from effectively

094

098

100

101

102

104

105

106

108

109

110

111

112

exploring novel patterns. Therefore, the efficacy of mining effective factors using LLMs is contingent upon selecting diversity-guiding factors as contextual samples to mitigate homogeneity. Additionally, encouraging ICL to explore new patterns is key to increasing the proportion of effective factors. 113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

153

154

156

157

158

159

160

In this paper, we present the FActor Mining Agent (FAMA), consisting of two main parts: Cross-Sample Selection (CSS) and Chainof-Experience (CoE) methods. CSS is designed to ensure the diversification of factor mining by amalgamating low correlation classes of factors as contextual samples, which empowers LLMs to incorporate diversity-guiding factors and mitigate the homogeneity of mined factors. CoE efficiently encourages ICL to explore new paradigms by incorporating the paths of mining effective factors as experiential prompts, which contributes to the further optimization of factor mining in LLMs. Our experimental results show better performance of our model in predicting stock market returns compared to previous approaches. Moreover, our model also demonstrates a superior annualized return and Sharpe ratio in the investment simulations.

Our contributions can be summarized as follows:

- To the best of our knowledge, we are the first ones to use LLMs as a bridge between symbolic and neural representations in the task of factor mining.
- We propose a factor mining agent (FAMA) to facilitate LLMs as factor miners, in which its components CSS and CoE are designed to tackle homogeneity issues and encourage ICL in exploring new directions respectively.
- We expand the capabilities of LLMs to perform factor mining tasks and present a series of experiments to demonstrate the effectiveness of our proposed model.

## 2 **Problem Formulation**

## 2.1 Financial Factor

Consider a stock dataset for *n* stocks over *T* trading days. The features of all stocks are denoted as  $\mathbf{X} = [x_1, x_2, ..., x_n]$ . Consider the *m* features, such as open and close prices, pertaining to each stock *j*, denoted as  $x_j \in \mathbb{R}^{m \times T}$ . We define the factor space as  $\mathcal{F}$ , where each factor  $f_i \in \mathcal{F}$  is defined as  $f_i : \mathbb{R}^{m \times T} \to \mathbb{R}^T$ . The value of factor  $f_i$  on stock *j* is defined as  $f_i(x_i) \in \mathbb{R}^T$ . To 161 162

163

164

166

167

169

170

171

172

173

174

175

176

177

178

179

181

183

185

186

187

190

191

192

193

194

195

196

198

199

200

201

204

conveniently represent the symbolic form of factors, we employ the symbol function  $s(f_i)$  to denote the symbolic text of factor  $f_i$ . For example,  $s(f_{101}) = "((close - open)/(high - low))"$ .

### 2.2 Factor Distance and Correlation

In practice, factor categorization has traditionally depended on artificial classification rooted in financial principles, such as momentum (Carhart, 1997) and trend (Han et al., 2016) factors. Despite the demonstrated high accuracy associated with this approach, it involves a labor-intensive process. To enhance the efficiency of factor classification, we advocate for a quantitative exploration of correlations among factors. We consider the factor space  $\mathcal{F}$  is equipped with a distance mapping  $d: \mathcal{F} \times \mathcal{F} \to \mathbb{R}$ , thereby establishing it as a complete distance space  $(\mathcal{F}, d)$ , then correlations between factors can be defined within this space  $(\mathcal{F}, d)$  as  $r : \mathcal{F} \times \mathcal{F} \rightarrow [-1, 1]$ . This approach enables a more efficient analysis of factor correlations without a labor-intensive process.

#### 2.3 Factor Mining

The goal of factor mining is to produce a new set of factors  $F \subset \mathcal{F}$  that will lead to better predictive performance of stocks in their portfolios. To evaluate the predictive performance of factors, we employ the Rank Information Coefficient (RankIC) (Chuan and Wu, 2019). RankIC measures the correlation between a factor's ranking in equity exposure and its subsequent return ranking. The RankIC on period t and average RankIC  $\gamma$  is defined as follows:

$$RankIC_{t}(f) = Corr(order_{t-1}^{f}, order_{t}^{r}),$$
  

$$\gamma(f) = \frac{1}{T} \sum_{t=1}^{T} RankIC_{t},$$
(1)

where  $order_{t-1}^{f}$  signifies the factor value ranking at time t-1, and  $order_{t}^{r}$  represents the return ranking at time t, with Corr(x, y) denoting the correlation coefficient between vectors x and y. Given the initial factor set  $F = \{f_1, ..., f_l\}$ , its effectiveness is assessed by computing the average RankIC of the factors within the set, as described below:

$$\gamma(F) = \mathbb{E}_i[\gamma(f_i)], f_i \in F.$$
(2)

We denote the combined model as  $g(\mathbf{X}; \mathbf{F}; \mathbf{I})$ , where  $\mathbf{X}$  is the stock feature matrix and  $\mathbf{I}$  is the prompt entered into the LLMs. Our goal is that the new set of factors mined by the model *g* achieves the optimal average RankIC, defined as follows:

$$g^{*}(\mathbf{X}; \mathbf{F}) = g(\mathbf{X}; \mathbf{F}; \mathbf{I}^{*})$$
$$\mathbf{I}^{*} = \operatorname{argmax}_{\mathbf{I}} \gamma(g(\mathbf{X}; \mathbf{F}; \mathbf{I})).$$
(3)

205

206

207

209

210

211

212

213

214

215

216

217

218

219

220

221

223

224

225

227

229

230

232

234

235

236

237

238

239

240

241

## **3** Factor Mining Agent

As illustrated in Figure 2, our proposed FActor Mining Agent (FAMA) consists of two main parts: (1) Cross-Sample Selection (CSS) and (2) Chainof-Experience (CoE). FAMA improves the mining factor effectiveness through iterative mining. In each iteration, FAMA generates diversity guiding factors via CSS and empirical paths through CoE as prompts fed into LLMs for mining factors.

#### 3.1 Definitions

To measure the distance and correlation between factors quantitatively mentioned in Section 2.2, we start with calculating the weighted average price of the stock pool. It is defined as:

$$\mathbf{p} = \mathbf{w}\mathbf{X},\tag{4}$$

where  $\mathbf{w} \in \mathbb{R}^n$  denotes the total market value weight corresponding to the company's stock. Subsequently, we calculate the factor exposure  $\mathbf{v}_i$  of factor  $f_i$  at the weighted average price p and employ z-score normalization as:

$$\mathbf{v}_{\mathbf{i}} = \frac{f_i(\mathbf{p}) - mean(f_i(\mathbf{p}))}{std(f_i(\mathbf{p}))}.$$
 (5)

Consequently, we define the distance between two factors as:

$$d(f_i, f_j) = \|\mathbf{v_i} - \mathbf{v_j}\|_2.$$
(6)

Then, the correlation coefficient between the factors is defined as:

$$r(f_i, f_j) = \operatorname{Corr}(\mathbf{v}_i, \mathbf{v}_j)$$
  
= 
$$\frac{\sum_{t=1}^{T} (\mathbf{v}_{it} - \overline{\mathbf{v}}_i) (\mathbf{v}_{jt} - \overline{\mathbf{v}}_j)}{\sum_{t=1}^{T} (\mathbf{v}_{it} - \overline{\mathbf{v}}_i)^2 (\mathbf{v}_{jt} - \overline{\mathbf{v}}_j)^2}.$$
 (7)

#### 3.2 Cross-Sample Selection

The CSS selects low-correlation guiding factors as contexts thereby avoiding homogeneity of the generated factors. It categorizes the factors into different classes, sampling from the classes to get a context sample of diversity factors. Here, we propose a clustering algorithm based on KMeans (Krishna and Murty, 1999) for factor clustering. The



Figure 2: An overview of the FAMA model. FAMA(CI-n) denotes the nth iteration of the FAMA model. Initially, (a) the input factors, stock data, and experience chain data are fed into the FAMA model. Subsequently, (b) the CoE module utilizes the outcomes of FAMA(CI-(k - 1)) to produce a novel CoE<sup>k</sup>, and incorporates the diverse guidance factors generated by the (c) CSS module to formulate a prompt. Lastly, the prompt is fed into the LLMs to mine a new factor of FAMA(CI-k) as illustrated in (d), which is then stored in the factor database.

243factor value  $\mathbf{v}_i$  of factor  $f_i$  obtained from Equa-244tion 5 is used for clustering. Initially, we ran-245domly select k factor values as clustering centers246{ $\mu_1, \mu_2, \cdots, \mu_k$ }. For each factor value  $\mathbf{v}_i$ , its247class is calculated as  $c(f_i) = argmin_j \|\mathbf{v}_i - \mu_j\|^2$ .248Subsequently, we update the clustering center using249the formula:

 $\boldsymbol{\mu}_j = \frac{1}{\sum\limits_{c(f_i)=j} 1} \sum\limits_{c(f_i)=j} \mathbf{v}_i.$ (8)

We define the loss of the factor cluster model as:

$$J = \sum_{i=1}^{k} \sum_{c(f_j)=i} \|\mathbf{v}_j - \boldsymbol{\mu}_i\|^2.$$
 (9)

253 The optimal classification is defined as:

251

254

$$c^* = \underset{\mathbf{c}}{\operatorname{argmin}} J. \tag{10}$$

255 We denote the set of factors  $C_i$  belonging to the 256 same class *i* as:

$$C_i = \{f_j | c^*(f_j) = i\}.$$
 (11)

Subsequently, we randomly draw a sample  $f^i$  from each category  $C_i$  to get a factor combination:

$$FC = [f^1, f^2, \cdots, f^k], \ f^i \in C_i.$$
 (12)

Finally,  $l(l \le k)$  factors in the factor combination FC are selected as context samples:

$$S = [s(f^{i_1}), s(f^{i_2}), \cdots, s(f^{i_l})], f^{i_j} \in FC.$$
(13)

#### 3.3 Chain-of-Experience

This part aims to involve past successful mining experiences in ICL to facilitate factor mining effectiveness. The generation of experience chains is divided into two phases: the initial generation phase and the enhanced generation phase. In the initial phase, we employ the initial set of factors for generation. Following the acquisition of the previous factor clustering results  $C_i$  through Equation 11 and with the size of  $C_i$  denoted as  $p_i$ , the initial experience chain for category  $C_i$  is generated. This generation process relies on the ranking result of  $\gamma$ as defined in Equation 1, which can be described as follows:

$$CoE_{i}^{0} = s(f_{1}^{(i)}) \to s(f_{2}^{(i)}) \to \dots \to s(f_{p_{i}}^{(i)}),$$
  
$$\gamma(f_{1}^{(i)}) \le \gamma(f_{2}^{(i)}) \le \dots \le \gamma(f_{p_{i}}^{(i)}).$$
  
(14)

In the enhanced phase, the experience chain used in the previous step is denoted as  $CoE_i^{(k-1)}$ . We choose ICL-generated factor  $f^{*(i)}$  with  $CoE_i^{(k-1)}$ having a higher  $\gamma$  than all chain factors. Then, we 261 262

265

267

268

269

270

272

273

274

275

276

277

278

279

compute the correlation r defined in Equation 7 for the new factor  $f^{*(i)}$  and factors on  $CoE_i^{(k-1)}$  to get 284 the highest correlation factor  $f_h^{(i)}$ . If the matched factor  $f_h^{(i)}$  is at the end of the chain  $C_i$ , the new factor is treated as an extension of the experience. Otherwise, the new factor  $f^{*(i)}$  represents a new experience, and it is introduced into the chain subsequent to a split triggered by the matching factor 290  $f_{h}^{(i)}$ . This process can be defined uniformly as:

$$CoE_{i}^{k} = s(f_{1}^{(i)}) \to \dots \to s(f_{h}^{(i)}) \to s(f^{*(i)}),$$
  
$$r(f_{h}^{(i)}, f^{*(i)}) \ge r(f_{j}^{(i)}, f^{*(i)}), \forall 0 \le j \le p_{i}.$$
  
(15)

Our proposed FAMA integrates Cross-Sample Selection (CSS) outlined in Section 3.2 and Chainof-Experience (CoE) detailed in Section 3.3 together to automatically generate diversity-guiding factor samples and experience chains for iterative factor mining. In each iteration, we utilize the sample generated through CSS as an in-context example and select the corresponding experience chain, then feed them into the LLMs. The output factor is added to the factor set and also contributes to a new experience chain. The specific algorithm is presented in Algorithm 1.

296

302

Algorithm 1: Factor Mining Agent **Data:** Initial factor set  $F = \{f_1, \dots, f_n\},\$ number of mining m. **Result:** Final factor set *F*, experience chain set  $CoE^m$ . 1 Generate the initial experience chain set  $CoE^0 = \{e_1, \cdots, e_k\};$ // Equation 14 **2** for  $i \leftarrow 1$  to m do  $C \leftarrow \text{Cluster}(F); // \text{Equation 11}$ 3  $S \leftarrow \text{SelectSamples}(C, CoE^{(i-1)});$ 4 // Equation 13 foreach  $(s, e) \in S$  do 5  $prompt \leftarrow s + e;$ 6  $f' \leftarrow \text{LLM}(prompt);$ 7 if  $\gamma(f') > max(\gamma(f)), \forall f \in e$  then 8  $e' \leftarrow \text{GenChain}(e, f');$ // Equation 15  $F \leftarrow F \cup \{f'\};$ 10  $CoE^i \leftarrow E \cup \{e'\};$ 11 12 return  $F, CoE^m$ ;

#### 4 **Experiments**

305 Our experimental investigation revolves around addressing three key questions: 307 • Q1: How does our proposed model compare 308 to prior factor mining methods? 309 • **O2:** Which factors within the experience 310 chain contribute to the enhancement of the 311 RankIC&RankICIR? 312 • Q3: How does our model perform under a 313 more realistic investment situation? 314 4.1 **Experiment Settings** 315 We use 38 factors from Alpha101 (Kakushadze, 316 2016) as our initial factor set F, the number of 317 clusters m is chosen to be 7, and the number of 318 randomly sampled factors l is set to 2. We choose 319 text-davinci-002<sup>1</sup> as the LLM for factor mining. 320 The full factors and prompt examples are listed in 321 Appendix A and Appendix B. 322 4.2 Datasets 323 Given that these factors are specifically crafted for 324 the U.S. stock market, we opt for the correspond-325 ing U.S. stock index, namely the S&P500 as the 326 stock set. Our dataset comprises all stocks from 327 the S&P500 index, with a focus on key fields including closing price, opening price, low price, 329 high price, adjusted closing price, and total vol-330 ume. The temporal scope of the stock data spans 331 from 2015/01/01 to 2022/01/01. The dataset is divided into a training set (2015/01/01-2020/01/01), 333

a validation set (2020/01/01-2021/01/01) and a test set (2021/01/01-2022/01/01). In our model, we only use stock data for the time period 2020/06/01-2021/01/01 as the training set, which is 10%amount of the provided training set.

334

335

336

337

339

340

341

342

343

344

345

346

347

#### Baselines 4.3

We explored SOTA models in recent years for comparison, encompassing both symbolic factor models and neural factor models as follows:

• Alpha101 (Kakushadze, 2016) publicly disclosed by WorldQuant LLC<sup>2</sup>, accompanied by precise code-based definitions. It serves as our initial set of factors from which our factors are derived.

```
<sup>1</sup>https://platform.openai.com/docs/model-index-for-
researchers
```

<sup>2</sup>https://www.worldquant.com/

Category	Model	Interpretability	Training data usage	Rank IC	Rank ICIR
Symbolic	Alpha101	$\checkmark$	-	0.025(0.000)	0.365(0.000)
	GP	$\checkmark$	100%	0.027(0.005)	0.149(0.034)
	LLM	$\checkmark$	10%	0.015(0.008)	0.139(0.011)
Neural	DTransformer	Х	100%	0.025(0.005)	0.124(0.015)
	ALSTM	×	100%	0.028(0.006)	0.167(0.021)
	FactorVAE	×	100%	0.048(0.008)	0.379(0.042)
Neural Symbolic	FAMA(C)	$\checkmark$	10%	0.023(0.006)	0.204(0.019)
	FAMA(I-1)	$\checkmark$	10%	0.016(0.006)	0.149(0.017)
	FAMA(CI-3)	$\checkmark$	10%	0.030(0.008)	0.372(0.031)
	FAMA(CI-7)	$\checkmark$	10%	0.054(0.010)	0.485(0.051)

Table 1: The performance of the compared models in returns prediction on the test dataset. Higher values for Rank IC and Rank ICIR indicate superior performance. *Interpretability* indicates that the mined factors are financially interpretable. *LLM* is the result of directly mining factors using LLMs. The term FAMA(C) corresponds to the CSS model. Additionally, FAMA(I-n) signifies the application of the COE iteration n. The **bold** part highlights the best-performing model in the evaluation. The mean and standard deviation of results from 10 experiments are reported.

• **GP** (Chen et al., 2021) Genetic programming algorithms create new factors through the mutation of factor expression trees, a widely cited model in factor mining.

348

351

352

360

361

362

- ALSTM (Qin et al., 2017) proposes a framework based on attentional mechanisms and long and short-term memory to predict stock trends.
- **DTransformer** (Wang et al., 2022) forecasts market indices by leveraging fundamental rules characterizing stock market dynamics through an encoder-decoder architecture and a full attention mechanism.
- FactorVAE (Duan et al., 2022) generates a prior risk factor return rate within the Variational Autoencoder (VAE) framework. It refines the prior factor return rate to approximate the posterior factor return rate.

#### 4.4 Cross-Sectional Returns Prediction

In this experiment, we employ both the neural and symbolic factor models to forecast future stock returns for answering Q1. The Average Rank IC is calculated between the forecasted and actual stock returns, as defined in Equation 2. To better illustrate the relationship between prediction effectiveness and risk, we introduce the Rank ICIR, defined as the ratio of the mean value of the Rank IC to the standard deviation:

Rank ICIR = 
$$\mathbb{E}_f[\frac{\gamma(f)}{\sigma_{RankIC_t(f)}}]$$
 (16)

376

377

378

379

381

382

384

385

386

As evidenced in Table 1, FAMA demonstrates superior performance compared to the most recent benchmark, FactorVAE. FAMA exhibits improvements of 0.006 on RankIC and 0.106 on RankICIR.

In addition, it can be observed from Table 1, that both CSS and CoE exhibit improvement in factor mining effects. Achieving satisfactory prediction results using CSS or CoE individually faces challenges. When CSS and CoE are employed together, the predictive performance of the model improves with an increasing number of mining iterations.



Figure 3: The results of parameter effects. Subfigure (a) illustrates RankIC and RankICIR in relation to the number of CoE iterations. Meanwhile, Subfigure (b) portrays the plot of RankIC and RankICIR with respect to the number of CSS samples.

To explore the impact of the number of CoE iterations on the model, we set the CoE iterations from 1 to 7 and verify the effect of the corresponding iterations. Results in Figure 3(a) show that the model's prediction effectiveness gradually improves with an increase in CoE iterations. The improvement effect of CoE largely depends on the generation effect of the previous round of factors.

390

394

400 401

402

403

404

405

406

407

To explore the impact of sample number selection on the model, we changed the number of crosssample selections and conducted experiments. As shown in Figure 3(b), until the number of samples is 3, increasing the number of samples improves the performance of the model. When the quantity of samples surpasses a threshold of three, the efficacy of the model shows a decrement. This observation signifies that an excessive abundance of samples fails to enhance the performance.

## 4.5 Randomized Modification of Chain-of-Experience



Figure 4: Impact of randomly deleting CoE nodes at different locations on model prediction. *Initial* is the performance of factors generated by retaining the complete experience chain of factors. *Head*, *Middle*, *Tail* are the performance of factors generated after randomly deleting the factors located at the head, middle, and tail of the experience chain.

In the pursuit of unraveling the fundamental com-408 ponents of the Chain-of-Experience (CoE) func-409 tion, we conducted an experiment that entailed the 410 random deletion of nodes within the CoE. The ob-411 jective of this endeavor is to address the inquiry 412 encapsulated in Q2. Nodes are categorized into 413 head nodes, middle nodes, and tail nodes. Given 414 that intermediate nodes may consist of multiple 415 nodes, we randomly select one among them as the 416 middle node. In each round of CoEs, we systemat-417 ically delete the head node, middle node, and tail 418 node, utilizing the modified CoEs for factor mining. 419 The results, averaged over multiple rounds, are de-420 picted in Figure 4. We observed that the removal 421 of initial nodes enhances the performance of factor 422 mining. This observation suggests that the inclu-423 sion of an excessive number of low-performing 424

nodes compromises the efficacy of factor mining in the LLM. Thus, it becomes imperative to adjust the length of the chain over time for optimal results. 425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

#### 4.6 Portfolio Investment Simulation

We intend to answer Q3 by designing an investment simulation of the stock market. For our model, we implement a multi-factor strategy to predict factors using the following approach. We select factors with positive average RankIC values during the valid period from 2020/01/01-2021/01/01. Funds for each factor are allocated based on weights given by:

$$w_i = \frac{RankIC_i^{past}}{\sum_{i=0}^{n} RankIC_i^{past}},$$
 (17)

where  $RankIC_i^{past}$  represents the mean RankIC value during the valid period. We choose stocks with the top 20% factor value to buy and sell them in next day.

We evaluate the portfolio investment performance using standard metrics, including Annualized Return (AR), Volatility (Vol), and Sharpe Ratios (SR):

$$AR = (1+R)^{252/N} - 1,$$
 (18)

$$Vol = \sigma_p * \sqrt{252},\tag{19}$$

$$SR = \frac{(R_p - R_f)}{\sigma_p} * \sqrt{252},$$
 (20)

where R represents the cumulative return rate, N is the total number of trading days,  $\sigma_p$  is the daily standard deviation of the portfolio,  $R_p$  is the expected daily return rate of the portfolio,  $R_f$  is the risk-free rate <sup>3</sup>.

Models	$AR(\uparrow)$	$Vol(\downarrow)$	$SR(\uparrow)$
S&P500	28.7%	13.0%	201.5%
GP	11.2%	6.8%	159.2%
Alpha101	13.2%	3.6%	340.8%
ALSTM	18.5%	22.3%	87.8%
DTransformer	18.6%	25.0%	80.8%
FactorVAE	31.8%	22.8%	132.2%
FAMA	39.0%	4.9%	667.6%

Table 2: Portfolio performance of the compared models on the test datasets;  $\uparrow$  indicates a larger value is better,  $\downarrow$  indicates a smaller value is better. The *S&P500* represents a portfolio comprising all *S&P500* stocks.

As depicted in Figure 5, the symbol-based approach exhibits lower volatility but yields comparatively lower returns. Conversely, neuro-based

<sup>&</sup>lt;sup>3</sup>For simplicity, we set the risk-free rate to zero.



Figure 5: Portfolio performance of factor mining models. *Cumulative Return* is defined as the ratio of the model's total return to the initial principal, calculated from the first day of the testing period to the end of the testing period.

485

486

487

457

approaches show higher returns, albeit accompanied by elevated volatility. It is noteworthy that our approach adeptly strikes a balance between returns and volatility, demonstrating a consistent performance throughout the investment simulation without experiencing significant fluctuations. This delicate equilibrium is achieved while concurrently realizing a commendable return, highlighting the robustness and stability inherent in our model.

It is evident from Table 2 that FAMA surpasses current SOTA models, in the context of portfolio investment simulation. Specifically, there is a notable increase of 7.2% in AR and a substantial improvement of 326.8% in the SR.

### 5 Related Work

Financial Factor Mining. The initial phase of factor mining involves the manual mining of factors. The Capital Asset Pricing Model (CAPM) (Sharpe, 1964), posits that the expected return of a financial asset primarily depends on the market's excess return. This contributed a groundbreaking factor to the financial field. To refine this conceptual framework, the Fama-French 3-factor model (Fama and French, 1993) extends the CAPM by introducing size and value risk factors alongside market risk factors. However, manual factor mining is considered labor-intensive. To address this limitation and efficiently mine effective factors in the market, various symbolic factor-based models have been proposed. AutoAlpha (Zhang et al., 2020) expedites the identification of promising factor search

spaces through the utilization of genetic algorithms. 488 Furthermore, AlphaEvolve (Cui et al., 2021) has 489 developed a factor mining framework grounded in 490 AutoML, facilitating the evolution of initial factors 491 into new factors characterized by excess returns 492 and correlations. Factors derived through symbolic 493 factor models exhibit clear factor calculation steps, 494 making them easily interpretable. However, con-495 strained by the vast symbolic factor target space, 496 these models are generally challenging to optimize. 497 This has prompted increased interest in the easy-to-498 optimize neural factor models. In a recent study, 499 AE (Gu et al., 2021) introduces a novel latent fac-500 tor conditional asset pricing model employing an 501 autoencoder. Additionally, FactorVAE (Duan et al., 502 2022) integrates a dynamic factor model with a 503 variational autoencoder to approximate the optimal 504 factor model. The neural factor model, a method 505 for extracting numerical characteristic factors from 506 stock data through feature extraction, is known for 507 its heightened optimization efficiency. Despite this 508 advantage, factors constrained by implicit features 509 present challenges in terms of artificial identifica-510 tion, resulting in a lack of interpretability in neural 511 factor models. In response to this, our proposed 512 model takes a strategic approach by combining 513 symbolic factors and leveraging neural factors for 514 feature extraction, achieving both financial inter-515 pretability and high efficiency in the realm of factor 516 mining. 517

## 6 Conclusion

In this paper, we consider Large Language Models (LLMs) as a neural symbolic model for financial factor mining. To facilitate LLMs to pursue our task, we proposed a model called Factor Mining Agent (FAMA), which comprises two integral components: Cross-Sample Selection (CSS) and Chain-of-Experience (CoE). CSS mitigates the homogeneity in the mined factors by amalgamating diverse guidance factors. CoE encourages In-Context Learning (ICL) to explore novel factor paradigms by leveraging the paths leading to the mining of effective factors as experiential prompts. Both CSS and CoE components are integrated into our factor mining agent to effectively mine financially interpretable factors. Experimental results demonstrate the effectiveness of our proposed approach. Our future work includes exploring more avenues to enhance the optimization of factor mining and addressing the illusionary effect of LLMs.

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

# 589 591 592 593 594 595 596 597 598 599 600 601 602 603 604 605 606 607 608 609 610 611 612 613 614 615 616 617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639

640

588

## 538 Limitations

When employing LLMs for factor mining, we observed the illusionary phenomenon of LLMs within
the financial domain that introduces interference in
the factor mining process. In future endeavors, our
emphasis will be directed towards mitigating the
illusionary effects of LLMs in the context of factor
mining.

## 546 Ethics Statement

547We utilize the OpenAI API in strict adherence to548the OpenAI User Rules for the generation of finan-549cial factors, ensuring the absence of harmful and550unethical information. Our approach has under-551gone validation in historical market scenarios and552expressly does not offer any form of investment553advice.

## References

555

561

566

568

571

573

574

575

576

577

578

579

580

581

- Zulfiqar Ali Imran, Abdullah Ejaz, Cristi Spulbar, Ramona Birau, and Periyapatna Sathyanarayana Rao Nethravathi. 2020. Measuring the impact of governance quality on stock market performance in developed countries. *Economic Research-Ekonomska Istraživanja*, 33(1):3406–3426.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Mark M Carhart. 1997. On persistence in mutual fund performance. *The Journal of finance*, 52(1):57–82.
- Tianxiang Chen, Wei Chen, and Luyao Du. 2021. An empirical study of financial factor mining based on gene expression programming. In 2021 4th International Conference on Advanced Electronic Materials, Computers and Software Engineering (AEMCSE), pages 1113–1117. IEEE.
- Yijian Chuan and Lan Wu. 2019. Information coefficient: From a new perspective. *Available at SSRN* 3387744.
- Can Cui, Wei Wang, Meihui Zhang, Gang Chen, Zhaojing Luo, and Beng Chin Ooi. 2021. Alphaevolve:
  A learning framework to discover novel alphas in quantitative investment. In *Proceedings of the 2021 International Conference on Management of Data*, pages 2208–2216.
- Yitong Duan, Lei Wang, Qizhong Zhang, and Jian Li. 2022. Factorvae: A probabilistic dynamic factor model based on variational autoencoder for predicting cross-sectional stock returns.

- Eugene F Fama and Kenneth R French. 1993. Common risk factors in the returns on stocks and bonds. *Journal of financial economics*, 33(1):3–56.
- Shihao Gu, Bryan Kelly, and Dacheng Xiu. 2021. Autoencoder asset pricing models. *Journal of Econometrics*, 222(1):429–450.
- Yue Guo, Zian Xu, and Yi Yang. 2023. Is chatgpt a financial expert? evaluating language models on financial natural language processing. *arXiv preprint arXiv:2310.12664*.
- Yufeng Han, Guofu Zhou, and Yingzi Zhu. 2016. A trend factor: Any economic gains from using information over investment horizons? *Journal of Financial Economics*, 122(2):352–375.
- Kevin B Hendricks and Vinod R Singhal. 2009. Demand-supply mismatches and stock market reaction: Evidence from excess inventory announcements. *Manufacturing & Service Operations Management*, 11(3):509–524.
- Kewei Hou, G Andrew Karolyi, and Bong-Chan Kho. 2011. What factors drive global stock returns? *The Review of Financial Studies*, 24(8):2527–2574.
- Ying Jin, Weilin Fu, Jian Kang, Jiadong Guo, and Jian Guo. 2019. Bayesian symbolic regression. *arXiv* preprint arXiv:1910.08892.
- Zura Kakushadze. 2016. 101 formulaic alphas. *Wilmott*, 2016(84):72–81.
- Bryan T Kelly, Seth Pruitt, and Yinan Su. 2019. Characteristics are covariances: A unified model of risk and return. *Journal of Financial Economics*, 134(3):501– 524.
- K Krishna and M Narasimha Murty. 1999. Genetic k-means algorithm. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 29(3):433–439.
- Nour Makke and Sanjay Chawla. 2022. Interpretable scientific discovery with symbolic regression: a review. *arXiv preprint arXiv:2211.10873*.
- Victor Ng, Robert F Engle, and Michael Rothschild. 1992. A multi-dynamic-factor model for stock returns. *Journal of Econometrics*, 52(1-2):245–266.
- Yao Qin, Dongjin Song, Haifeng Chen, Wei Cheng, Guofei Jiang, and Garrison Cottrell. 2017. A dual-stage attention-based recurrent neural network for time series prediction. *arXiv preprint arXiv:1704.02971*.
- Stephen A Ross. 2013. The arbitrage theory of capital asset pricing. In *Handbook of the fundamentals of financial decision making: Part I*, pages 11–30. World Scientific.
- William F Sharpe. 1964. Capital asset prices: A theory of market equilibrium under conditions of risk. *The journal of finance*, 19(3):425–442.

Rahul Verma and GökÇe Soydemir. 2009. The impact of individual and institutional investor sentiment on the market price of risk. *The Quarterly Review of Economics and Finance*, 49(3):1129–1145.

645

646

647

648

649

650 651

652

653

654 655

- Chaojie Wang, Yuanyuan Chen, Shuqi Zhang, and Qiuhui Zhang. 2022. Stock market index prediction using deep transformer model. *Expert Systems with Applications*, 208:118128.
  - Yi Yang, Yixuan Tang, and Kar Yan Tam. 2023. Investlm: A large language model for investment using financial domain instruction tuning. *arXiv preprint arXiv:2309.13064*.
- Tianping Zhang, Yuanqi Li, Yifei Jin, and Jian Li. 2020. Autoalpha: An efficient hierarchical evolutionary algorithm for mining alpha factors in quantitative investment. *arXiv preprint arXiv:2002.08245*.

# A Factor

	Factor				
0	" (-1 * correlation(rank(delta(log(volume), 2)), rank(((close - open) / open)), 6))"				
1	" (-1 * correlation(rank(open), rank(volume), 10))"				
2	" (-1 * Ts_Rank(rank(low), 9))"				
3	" (rank((open - (sum(vwap, 10) / 10))) * (-1 * abs(rank((close - vwap)))))"				
4	" (-1 * correlation(open, volume, 10))"				
5	"(-1 * rank(((sum(open, 5) * sum(returns, 5)) - delay((sum(open, 5) * sum(returns, 5)),10))))"				
6	" ((rank(ts_max((vwap - close), 3)) + rank(ts_min((vwap - close), 3))) *rank(delta(volume, 3)))"				
7	" (sign(delta(volume, 1)) * (-1 * delta(close, 1)))"				
8	" (-1 * rank(covariance(rank(close), rank(volume), 5)))"				
9	" ((-1 * rank(delta(returns, 3))) * correlation(open, volume, 10))"				
10	" (-1 * sum(rank(correlation(rank(high), rank(volume), 3)), 3))"				
11	" (-1 * rank(covariance(rank(high), rank(volume), 5)))"				
12	" (((-1 * rank(ts rank(close, 10))) * rank(delta(delta(close, 1), 1))) * rank(ts rank((volume / adv20),				
	5)))"				
13	" (-1 * rank(((stddev(abs((close - open)), 5) + (close - open)) + correlation(close, open, 10))))"				
14	" (((-1 * rank((open - delay(high, 1)))) * rank((open - delay(close, 1)))) * rank((open - delay(low,				
15	" (-1 * (delta(correlation(high, volume, 5), 5) * rank(stddev(close, 20))))"				
16	rank(((((-1 * returns) * adv20) * vwap) * (high - close)))				
17	" (-1 * ts_max(correlation(ts_rank(volume, 5), ts_rank(high, 5), 5), 3))"				
18	$\frac{1}{2} = \frac{1}{2} $				
19	" (((1.0 - rank(((sign((close - delay(close, 1))) + sign((delay(close, 1) - delay(close, 2))))				
	+sign((delay(close, 2) - delay(close, 3))))) * sum(volume, 5)) / sum(volume, 20))"				
20	$rank((-1 * ((1 - (open / close))\hat{1})))$				
21	" rank(((1 - rank((stddev(returns, 2) / stddev(returns, 5)))) + (1 - rank(delta(close, 1)))))"				
22	" ((Ts_Rank(volume, 32) * (1 - Ts_Rank(((close + high) - low), 16))) * (1 -Ts_Rank(returns, 32)))"				
23	" ((-1 * rank(Ts_Rank(close, 10))) * rank((close / open)))"				
24	" ((-1 * rank(stddev(high, 10))) * correlation(high, volume, 10))"				
25	$(((high * low)\hat{0}.5) - vwap)$				
26	(rank((vwap - close)) / rank((vwap + close)))				
27	" (ts_rank((volume / adv20), 20) * ts_rank((-1 * delta(close, 7)), 8))"				
28	" (-1 * correlation(high, rank(volume), 5))"				
29	" (-1 * ((rank((sum(delay(close, 5), 20) / 20)) * correlation(close, volume, 2))				
	*rank(correlation(sum(close, 5), sum(close, 20), 2))))"				
30	" ((((rank((1 / close)) * volume) / adv20) * ((high * rank((high - close))) / (sum(high, 5) /5))) -				
	rank((vwap - delay(vwap, 5))))"				
31	" (-1 * ts_max(rank(correlation(rank(volume), rank(vwap), 5)), 5))"				
32	" (-1 * delta((((close - low) - (high - close)) / (close - low)), 9))"				
33	$((-1 * ((low - close) * (open\hat{5}))) / ((low - high) * (close\hat{5})))$				
34	" (-1 * correlation(rank(((close - ts_min(low, 12)) / (ts_max(high, 12) - ts_min(low, 12)))),				
	rank(volume), 6))"				
35	" (0 - (1 * ((2 * scale(rank(((((close - low) - (high - close)) / (high - low)) * volume)))) -				
	scale(rank(ts_argmax(close, 10)))))"				
36	" ((rank(delay(((high - low) / (sum(close, 5) / 5)), 2)) * rank(rank(volume))) / (((high -low) /				
	(sum(close, 5) / 5)) / (vwap - close)))"				
37	((close - open) / ((high - low) + .001))				

## B Prompt

#### 59 B.1 Factor Mining

**Instruction:** ### Instruction 661 You are an alpha generator. You should follow the following codes: 1. The inputs are the alpha factors that are currently performing well, and you are required to output a new alpha factor that is generated from the fusion of these factors, and your factor must be different from the input factor. 2. Complete <fill\_alpha\_formula> with new alpha's formula. 3. Do not repeat example answer. 4. The specific operator is defined as follows: rank(x) = cross-sectional rank delay(x, d) = value of x d days ago670 correlation(x, y, d) = time-serial correlation of x and y for the past d days 671 covariance(x, y, d) = time-serial covariance of x and y for the past d daysscale(x, a) = rescaled x such that sum(abs(x)) = a (the default is a = 1)673 delta(x, d) = todays value of x minus the value of x d days ago signedpower(x, a) =  $x^a$  $decay_linear(x, d) = weighted moving average over the past d days with linearly$ decaying weights d, d 1, ..., 1 (rescaled to sum up to 1) 677 indneutralize(x, g) = x cross-sectionally neutralized against groups g ( 678 subindustries, industries, sectors, etc.), i.e., x is cross-sectionally demeaned within each group g  $ts_{0}(x, d) = operator 0$  applied across the time-series for the past d days; noninteger number of days d is converted to floor(d)  $ts_min(x, d) = time-series min over the past d days$  $ts_max(x, d) = time-series max over the past d days$ 684  $max(x, d) = ts_max(x, d)$ sum(x, d) = time-series sum over the past d days product(x, d) = time-series product over the past d daysstddev(x, d) = moving time-series standard deviation over the past d days5. Follow the path in "improve\_path". -> Indicates that the following factors have better performance than the previous factors. You should refer it to build new alpha. **Input Example:** ### Input alphas: (-1 \* correlation(open, volume, 10)) generate\_factor\_num: 1 improve\_path: close/open -> rank(close)/rank(open) **Output Example:** 

598 ### Answer: 599 ["(-1 \* correlation(rank(open), rank(volume), 10))"]