

Question 1: Application for lending department of a bank

Feng Qiu (feng.qiu.1101@gmail.com)
Ying(Alice) Xia (alice.ying.xia@gmail.com)

November 23, 2025

Contents

1 Literature Review	5
1.1 Constructive Accounting Approach	5
1.2 Structural Models Based on Double-Entry Accounting	6
1.3 Identity-Constrained Structural State-Space Models	7
2 A Simple Model of the Balance Sheet Based on the Tools of Velez-Pareja(09) and Velez-Pareja(10).	9
2.1 Model Structure	9
2.1.1 Definitions of Variables	9
2.1.2 Inputs From Previous Period	10
2.1.3 Policy Inputs Estimated From Historical Data	10
2.1.4 Inputs Required for Forecasting Period t	10
2.1.5 Variables Forecasted Directly for Period t	11
2.1.6 Quantities Derived From the Forecast	11
2.1.7 Final Outputs	12
2.2 Evaluation	12
2.3 S&P 500 Forecast Evaluation Results for 2024	13
2.3.1 Overall Forecast Accuracy	13
2.3.2 Most Accurate Forecasts	14
2.3.3 Least Accurate Forecasts	14
2.3.4 Data Coverage and Failures	16
3 Potential Improvement: Predicting Revenue Growth with MLP and Firm Fixed Effects via Embeddings	17
3.1 Overview	17
3.2 Motivation for Firm Embeddings	17
3.3 Model Inputs	18
3.3.1 Firm Embedding	18
3.3.2 Lagged Firm-Level Predictors	18
3.4 Model Architecture: MLP with Firm Embeddings	19
3.4.1 Network Structure	19
3.5 Interpretation of Learned Embeddings	19
3.6 Comparison with Classical Fixed-Effects Models	20

Chapter 1

Literature Review

Forecasting a firm's balance sheet over a one-year horizon involves two core requirements. The projections must be economically meaningful—consistent with the firm's operations, investment needs, and financing choices—while also respecting the accounting identities and subtotals that hold with certainty each period. Most fundamentally,

$$\text{Assets} = \text{Liabilities} + \text{Equity},$$

and related leverage and working-capital relationships must be satisfied exactly. A forecasting framework that meets only one of these requirements risks producing either mechanically inconsistent or financially implausible statements. The balance sheet of a firm is a multivariate time series governed by strict linear constraints. Traditional forecasting models—ARIMA, VAR models, machine learning regressors, and neural networks—do not guarantee that their predictions satisfy these accounting identities.

1.1 Constructive Accounting Approach

A natural starting point is the constructive accounting approach proposed by Vélez-Pareja (2007). His working paper demonstrates that financial statements can be forecasted *without* relying on plug variables and *without* generating circular references, provided the modeling process is anchored in explicit accounting mechanics. In practice, this requires (i) well-defined operating assumptions, (ii) detailed schedules for depreciation, working capital, and operating expenses, (iii) a comprehensive Cash Budget, and (iv) explicit modeling of short- and long-term financing flows. Under this structure, the Income Statement and Balance Sheet are derived only after financing needs and cash movements are fully determined. Any residual imbalance is treated as a genuine modeling error. This approach reasserts the fundamental role of double-entry bookkeeping rather than allowing ad-hoc spreadsheet adjustments.

Vélez-Pareja (2009) extends this logic to a full financial planning and valuation framework. He shows how a structured parameter table and Cash Budget can be used to generate consistent forecasts of operational behavior, financing surpluses and deficits, and valuation through the Capital Cash Flow (CCF) method. The framework underscores the interplay between short-term operating gaps, long-term investment needs, and excess cash accumulation. A central insight is that forecasting and valuation should both derive from a single, coherent set of accounting flows rather than from independently specified modules.

A related circularity in discounted cash-flow valuation—namely, the fact that leverage affects the discount rate while the discount rate simultaneously affects firm value—has been addressed analytically by Mejía-Peláez and Vélez-Pareja (2011). Their working paper provides closed-form expressions for levered firm value, the cost of equity, and WACC under various tax-shield assumptions, eliminating the need for iterative spreadsheet routines. The contribution is important because it demonstrates that internal consistency in valuation is achievable through analytical structure rather than numerical iteration.

Collectively, these constructive accounting contributions offer a rigorous foundation for producing internally coherent financial statements. However, they remain essentially static: they focus on building consistent financial statements, not on generating multi-period, probabilistic forecasts for large panels of firms.

1.2 Structural Models Based on Double-Entry Accounting

A central challenge in forecasting a firm’s balance sheet is that its line items are jointly determined. Double-entry bookkeeping ensures that every transaction affects at least two accounts, which creates deterministic articulation identities across the balance sheet, income statement, and cash-flow statement. Standard univariate forecasting or separate regressions for individual accounts often violate these identities and ignore the inherent simultaneity among accounting variables. Two influential papers—Christodoulou and McLeay (2014) and Christodoulou and McLeay (2019)—develop econometric frameworks that explicitly incorporate double entry into the estimation of accounting systems. These studies provide the conceptual foundation for constructing coherent forecasting models for financial statements.

Christodoulou and McLeay (2014) emphasize that accounting variables are inherently endogenous because they are linked by rank-deficient linear identities such as the clean-surplus relation, fixed-asset articulation, and the fundamental balance-sheet identity ($\text{assets} = \text{liabilities} + \text{equity}$). Regressing one accounting variable on others using OLS leads to biased or uninterpretable coefficients, because the regressors themselves are codetermined through the same identities. To address this, the authors propose a structural regression system in which:

- each accounting variable is modeled by its own behavioral equation;
- the canonical accounting identities enter the system as separate equations with no error terms;
- these identity equations impose linear restrictions on the coefficients of the behavioral equations;
- estimation proceeds via constrained least squares, SUR, or 2SLS/3SLS using lagged accounting variables and external covariates as instruments.

They show that imposing double-entry constraints yields economically meaningful parameter estimates even when unconstrained OLS produces superficially reasonable fitted

values. The principal contribution is not improved point forecasts per se, but ensuring that the estimated relationships and the resulting forecasts are internally coherent and compliant with accounting logic.

[Christodoulou and McLeay \(2019\)](#) extend the structural identity approach to changes in net operating assets (ΔNOA) and net financial claims (ΔNFC), focusing on the high-level identity:

$$\Delta\text{NOA}_t = -\Delta\text{NFC}_t.$$

They show how this relation arises from the joint articulation of the earnings identity and the cash-flow identity when all debit/credit flows are decomposed into operational, accrual, financing, and tax components. The paper develops a constrained SUR system in which:

- components of ΔNOA (e.g., sales, receipts, purchases, payments, provisions, depreciation, capex) follow autoregressive equations;
- components of ΔNFC are modeled analogously;
- the double-entry identity $\Delta\text{NOA} + \Delta\text{NFC} = 0$ is imposed as a parameter restriction across equations.

Unconstrained OLS produces unstable or sign-inconsistent estimates, whereas constrained estimation yields coherent and interpretable dynamics for both operational and financing flows. This framework demonstrates how a high-level accounting identity can anchor a dynamic system that produces internally consistent forecasts for the entire articulated financial structure.

1.3 Identity-Constrained Structural State-Space Models

A broader stream of research examines how time-series models can be adapted so that forecasts automatically satisfy accounting identities. [Pandher \(2007\)](#) offers an early and influential contribution in this direction. Working with monetary aggregates, he develops a state-space model in which the underlying components—such as net foreign and net domestic assets—are treated as latent processes, and the accounting identity linking them is imposed directly within the measurement equations. One of the advantages of this setup is that the Kalman filter produces filtered and forecasted values that *cannot* violate the identity, because consistency is built into the model rather than imposed as an external adjustment. Empirical results for Germany, the UK, and the US show that this constrained approach not only preserves internal coherence but also improves forecast accuracy relative to simple autoregressive models and unconstrained state-space formulations. Although his application focuses on monetary aggregates, the same modeling strategy extends naturally to corporate balance sheets: when assets and liabilities evolve jointly over time, embedding the balance-sheet identity within the state-space structure provides a disciplined way to maintain coherence across all forecasted items.

[Angelini et al. \(2008\)](#) extend similar ideas to a broader macroeconomic setting through a dynamic factor model designed to estimate and forecast euro-area national accounts. Their

model uses a large panel of monthly indicators, extracted through a small number of latent factors, to jointly produce monthly interpolations and quarterly forecasts of GDP and its major components. Importantly, the state-space representation encodes both the temporal aggregation from months to quarters and the cross-equation accounting relationships inherent in the national accounts. Using real-time data vintages, they show that this constrained factor-model approach outperforms traditional tools such as bridge equations, quarterly autoregressions, and VAR models—not only by improving predictive accuracy but also by ensuring that the reconstructed monthly paths remain consistent with the national accounting framework.

Taken together, these two studies demonstrate that state-space methods provide a flexible and principled framework for forecasting systems in which multiple variables must evolve jointly while obeying strict accounting relations. For balance-sheet forecasting, this line of research offers a clear and practical template: model the asset and liability components as jointly driven latent processes, encode the accounting constraints directly in the structure of the state space, and use the Kalman filter to generate forecasts that are statistically efficient and internally coherent.

Chapter 2

A Simple Model of the Balance Sheet Based on the Tools of Velez-Pareja(09) and Velez-Pareja(10).

2.1 Model Structure

This section describes the full structure of the forecasting model for earnings and the balance sheet, following the transaction-based and non-circular framework of Vélez-Pareja (2007, 2009); Mejía-Peláez and Vélez-Pareja (2011). All variables, inputs, and dependencies are clearly defined. The goal is to ensure that every year's financial statements satisfy accounting identities without using plugs.

2.1.1 Definitions of Variables

Income statement items

- S_t : Revenue,
- COGS_t : Cost of goods sold,
- SGA_t : Selling, general and administrative expenses,
- Dep_t : Depreciation expense,
- EBIT_t : Earnings before interest and taxes,
- IntExp_t : Interest expense,
- EBT_t : Earnings before taxes,
- Tax_t : Taxes,
- NI_t : Net income.

Working capital items

- AR_t : Accounts receivable, Inv_t : Inventory, AP_t : Accounts payable,

$$NWC_t = AR_t + Inv_t - AP_t.$$

Other assets and liabilities

C_t : Cash, PPE_t : Property, plant, and equipment, OA_t : Other assets,
 STD_t : Short-term debt, LTD_t : Long-term debt, OL_t : Other liabilities,
 E_t : Equity, Div_t : Dividends.

Cash flow quantities

OCF_t : Operating cash flow, $Capex_t$: Capital expenditure, $FCFF_t$: Free cash flow,
 C_t^{pre} : Cash before financing, Deficit_t , Surplus_t : Financing imbalance.

2.1.2 Inputs From Previous Period

These values come directly from the historical year $t - 1$ (for $t = 1$) or from the prior forecasted year (for $t \geq 2$):

$$S_{t-1}, NI_{t-1}, AR_{t-1}, Inv_{t-1}, AP_{t-1}, NWC_{t-1}, \\ C_{t-1}, PPE_{t-1}, STD_{t-1}, LTD_{t-1}, E_{t-1}, OA_{t-1}, OL_{t-1}.$$

These are required because accounting variables are inherently recursive.

2.1.3 Policy Inputs Estimated From Historical Data

Using 3–5 years of Yahoo Finance statements, we estimate stable firm policies:

$$gs : \text{Revenue growth rate}, m_{COGS}, m_{SGA} : \text{Cost ratios}, \\ \delta : \text{Depreciation rate}, r_d : \text{Interest rate on debt}, \tau : \text{Effective tax rate}, \\ p : \text{Dividend payout ratio}, \\ DSO, DIH, DPO : \text{Working-capital policies}, \\ \kappa : \text{Capex-to-sales ratio}, \chi : \text{Minimum cash as \% of revenue}.$$

These do *not* change during the forecast horizon.

2.1.4 Inputs Required for Forecasting Period t

To compute results for year t , the model needs:

1. Previous-period financial position.
2. Policy parameters.
3. Any user-specified adjustments (e.g. changes in capex or leverage).

2.1.5 Variables Forecasted Directly for Period t

Revenue

$$S_t = S_{t-1}(1 + g_S).$$

Operating costs

$$\text{COGS}_t = m_{\text{COGS}}S_t, \quad \text{SGA}_t = m_{\text{SGA}}S_t.$$

Depreciation

$$\text{Dep}_t = \delta \cdot PPE_{t-1}.$$

Earnings

$$\text{EBIT}_t = S_t - \text{COGS}_t - \text{SGA}_t - \text{Dep}_t.$$

$$\text{IntExp}_t = r_d(\text{STD}_{t-1} + \text{LTD}_{t-1}).$$

$$\text{EBT}_t = \text{EBIT}_t - \text{IntExp}_t.$$

$$\text{Tax}_t = \tau \cdot \max(\text{EBT}_t, 0).$$

$$\text{NI}_t = \text{EBT}_t - \text{Tax}_t.$$

Working capital

$$AR_t = DSO \cdot \frac{S_t}{365}, \quad Inv_t = DIH \cdot \frac{\text{COGS}_t}{365}, \quad AP_t = DPO \cdot \frac{\text{COGS}_t}{365}.$$

Capital expenditure and PPE

$$\text{Capex}_t = \kappa S_t,$$

$$PPE_t = PPE_{t-1} + \text{Capex}_t - \text{Dep}_t.$$

Other assets and liabilities

$$OA_t = OA_{t-1}(1 + g_{OA}), \quad OL_t = OL_{t-1}(1 + g_{OL}).$$

2.1.6 Quantities Derived From the Forecast

Operating cash flow

$$OCF_t = \text{EBIT}_t(1 - \tau) + \text{Dep}_t - (NWC_t - NWC_{t-1}).$$

Free cash flow

$$FCFF_t = OCF_t - \text{Capex}_t.$$

Cash before financing

$$C_t^{\text{pre}} = C_{t-1} + FCFF_t - \text{Div}_t - \text{IntExp}_t.$$

Minimum cash policy

$$C_t^{\min} = \chi S_t.$$

Financing requirements

$$\text{Deficit}_t = \max(0, C_t^{\min} - C_t^{\text{pre}}), \quad \text{Surplus}_t = \max(0, C_t^{\text{pre}} - C_t^{\min}).$$

Debt adjustments

$$STD_t = STD_{t-1} + \Delta STD_t^+ - \Delta STD_t^-,$$

where increases occur when financing deficits and reductions occur when surplus cash is used for repayment.

Cash after financing

$$C_t = C_t^{\text{pre}} + \Delta STD_t^+ - \Delta STD_t^-.$$

Equity update

$$\Delta RE_t = NI_t - Div_t, \quad E_t = E_{t-1} + \Delta RE_t.$$

2.1.7 Final Outputs

For each forecast year t , we obtain:

- Full income statement (all items from S_t to NI_t).
- Cash flow statement (including OCF_t , $FCFF_t$, cash budget).
- Balance sheet (cash, working capital, PPE, debt, equity).

By construction,

$$\text{Assets}_t = \text{Liabilities}_t + E_t$$

holds automatically, without the use of balancing plugs.

2.2 Evaluation

Financial-statement line items vary substantially in scale (for example, revenue versus cash, accounts receivable, or property, plant and equipment). Using raw errors such as MAE or RMSE would cause large-magnitude items to dominate the evaluation. Therefore, we adopt a scale-free error measure: the *Symmetric Mean Absolute Percentage Error (sMAPE)*.

For a series of true values y_i and forecasts \hat{y}_i , the sMAPE is defined as:

$$\text{sMAPE}(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n \frac{|\hat{y}_i - y_i|}{(|y_i| + |\hat{y}_i|)/2}.$$

This formulation normalizes the error by the average magnitude of the actual and forecasted values, making the metric symmetric with respect to over- and under-prediction and

robust when true values are small. Because sMAPE is unitless, it is directly comparable across items with different scales.

To evaluate forecasting accuracy for a company, we compute sMAPE for each of the K income-statement and balance-sheet items (e.g. revenue, net income, cash, accounts receivable, inventory, PPE, accounts payable, short-term debt, long-term debt, equity). Let sMAPE_k denote the error for item k . The overall company-level forecasting performance is then obtained by taking the simple average across all items:

$$\text{sMAPE}_{\text{company}} = \frac{1}{K} \sum_{k=1}^K \text{sMAPE}_k.$$

This provides a unified, scale-independent measure of how well the model forecasts the full set of financial statements. Lower values indicate better performance (for example, 5–10% is excellent, 10–20% good, 20–30% moderate, and values above 30% typically indicate poor or unstable forecasting assumptions).

2.3 S&P 500 Forecast Evaluation Results for 2024

This section summarizes the out-of-sample forecast evaluation for the S&P 500 constituents for fiscal year 2024. The evaluation covers all 500 index members. Forecasts were successfully generated and evaluated for 494 firms (98.8%), while 6 firms (1.2%) could not be evaluated due to missing financial statements on Yahoo Finance.

2.3.1 Overall Forecast Accuracy

Across the 494 evaluated stocks, the overall average symmetric mean absolute percentage error (sMAPE) is 31.62%, with a median of 30.77%. The model performs slightly better on the balance sheet than on the income statement: the average sMAPE is 33.77% for income-statement items and 29.73% for balance-sheet items.

The cross-sectional distribution of overall sMAPE is summarized below:

- Minimum: 4.24% (ticker V),
- 25th percentile: 21.02%,
- Median: 30.77%,
- 75th percentile: 40.61%,
- Maximum: 88.94% (ticker SW).

These statistics indicate that, while a subset of firms is forecast with very high accuracy, there is substantial dispersion across the cross-section. In particular, firms with volatile earnings, large investment cycles, or major changes in capital structure are harder to capture with the simple one-year-ahead model.

2.3.2 Most Accurate Forecasts

Table 2.1 reports the 20 firms with the lowest overall sMAPE. For each firm, we decompose the error into income-statement (IS) and balance-sheet (BS) components.

Table 2.1: Top 20 Most Accurate Forecasts (Lowest Overall sMAPE, 2024)

Rank	Ticker	Overall sMAPE	IS sMAPE	BS sMAPE
1	V	4.24%	3.13%	4.98%
2	TRV	4.53%	1.94%	5.83%
3	PNC	5.42%	3.34%	6.47%
4	AMP	6.69%	5.27%	7.54%
5	PG	6.92%	8.50%	5.82%
6	AMZN	7.06%	4.54%	8.50%
7	SYF	8.14%	8.68%	7.78%
8	PGR	8.16%	9.37%	7.68%
9	REGN	8.33%	2.99%	11.39%
10	PEP	8.41%	8.34%	8.46%
11	MTCH	8.44%	7.20%	9.26%
12	CHTR	9.19%	12.74%	6.09%
13	ITW	9.39%	7.98%	10.38%
14	NOC	9.42%	5.77%	11.51%
15	CTSH	9.49%	11.25%	7.95%
16	CPT	9.51%	6.68%	11.21%
17	EQIX	9.73%	6.19%	12.82%
18	MKC	9.89%	11.84%	8.53%
19	KMB	9.90%	12.61%	8.00%
20	HIG	10.07%	15.01%	6.55%

A number of mature, large-cap firms (e.g., V, PG, AMZN, PEP) feature among the best-performing stocks, suggesting that the simple model can capture stable earnings and balance-sheet patterns more reliably for firms with relatively smooth growth and conservative financial policies.

2.3.3 Least Accurate Forecasts

Table 2.2 lists the 20 firms with the highest overall sMAPE. These cases highlight where the current specification struggles most, often due to volatile earnings, substantial investment cycles, or structural changes in leverage and funding.

In many of these firms (e.g., TTWO, MU, EQT, MRNA), the income-statement errors are particularly large, reflecting substantial short-term earnings volatility that is difficult to extrapolate from historical data. In other cases (e.g., SW, SMCI, PEG), the balance sheet is the main source of error, indicating that the simple balance-sheet dynamics may be too rigid to capture rapid changes in leverage, working capital, or investment intensity.

Table 2.2: Bottom 20 Least Accurate Forecasts (Highest Overall sMAPE, 2024)

Rank	Ticker	Overall sMAPE	IS sMAPE	BS sMAPE
1	SW	88.94%	68.08%	103.55%
2	MS	85.61%	104.14%	71.71%
3	TTWO	84.27%	111.63%	60.33%
4	MU	82.94%	99.37%	71.45%
5	EQT	80.43%	103.17%	60.54%
6	AVGO	75.53%	79.43%	72.80%
7	BK	71.76%	92.69%	56.06%
8	MMM	70.13%	85.99%	59.03%
9	MRNA	69.18%	100.16%	47.49%
10	LYV	67.95%	82.45%	57.79%
11	NCLH	66.48%	71.69%	62.84%
12	SCHW	66.21%	80.81%	55.26%
13	GS	63.43%	78.35%	52.24%
14	NVDA	62.75%	74.21%	54.73%
15	SMCI	62.59%	48.04%	72.78%
16	AMTM	61.78%	69.48%	55.03%
17	PFE	61.56%	68.56%	56.65%
18	WYNN	61.38%	56.32%	64.92%
19	PEG	60.10%	55.57%	64.63%
20	ALB	60.00%	72.12%	51.52%

2.3.4 Data Coverage and Failures

A small subset of stocks (DFS, WBA, PARA, ANSS, HES, JNPR) could not be evaluated because of missing financial statements on Yahoo Finance. These data gaps reduce the effective coverage of the S&P 500 universe but do not materially affect the aggregate patterns reported above.

Chapter 3

Potential Improvement: Predicting Revenue Growth with MLP and Firm Fixed Effects via Embeddings

3.1 Overview

This chapter develops a modern machine learning framework for forecasting revenue growth using panel data across firms and time. Traditional panel econometric models incorporate firm fixed effects through additive intercepts, typically expressed as

$$g_{i,t} = \alpha_i + f(X_{i,t-1}) + \varepsilon_{i,t},$$

where $g_{i,t}$ denotes revenue growth for firm i at time t and $X_{i,t-1}$ denotes lagged predictors. In contrast, the approach developed here replaces the scalar firm effect α_i with a *learned dense embedding vector* $e_i \in \mathbb{R}^k$, capturing persistent and latent firm-specific characteristics. This embeds firm fixed effects directly into the machine learning architecture, allowing for flexible nonlinear mappings from firm traits and financial features to future revenue growth.

3.2 Motivation for Firm Embeddings

Firm embeddings generalize the notion of fixed effects. Instead of modeling a single intercept α_i for each firm, we learn a vector e_i that captures richer, multidimensional structure. This embedding may encode:

- long-run growth regimes,
- business model features,
- persistent profitability characteristics,
- industry membership,
- structural differences in financial ratios,

- typical revenue volatility or momentum patterns.

These embeddings are learned end-to-end within the model, along with the weights of the prediction network. As a result, they adapt optimally to the revenue dynamics observed across the entire panel dataset.

3.3 Model Inputs

For each firm i at time t , the model takes two types of inputs:

3.3.1 Firm Embedding

Each firm is assigned a unique integer identifier, which is mapped through an embedding layer to produce a dense vector:

$$e_i = \text{Embed}(i), \quad e_i \in \mathbb{R}^k.$$

This vector replaces the traditional fixed effect.

3.3.2 Lagged Firm-Level Predictors

Lagged financial features form the second component of the input. Let

$$X_{i,t-1} \in \mathbb{R}^p$$

contain predictors such as:

- lagged revenue growth: $g_{i,t-1}, g_{i,t-2}$,
- profitability ratios: ROA, ROE, gross margin,
- investment ratios: CAPEX/Sales, R&D/Sales,
- leverage measures: Debt/Assets, interest coverage,
- liquidity ratios: Cash/Assets,
- working-capital deltas: $\Delta\text{Receivables}$, $\Delta\text{Inventory}$,
- firm size: $\log(\text{Revenue})$, $\log(\text{Assets})$,
- industry median revenue growth,
- macroeconomic indicators.

The model input vector is the concatenation:

$$\text{Input}_{i,t} = [e_i \parallel X_{i,t-1}].$$

3.4 Model Architecture: MLP with Firm Embeddings

The prediction model consists of two main components:

1. an **embedding layer** mapping firm IDs to e_i ;
2. a **multi-layer perceptron (MLP)** that models nonlinear relationships between the concatenated inputs and revenue growth.

Formally, the prediction for revenue growth is:

$$\hat{g}_{i,t} = h_\theta([e_i \parallel X_{i,t-1}]),$$

where $h_\theta(\cdot)$ denotes the neural network with parameters θ .

The model is trained by minimizing the squared error objective:

$$\min_{\theta, \{e_i\}} \sum_{i,t} (g_{i,t} - \hat{g}_{i,t})^2.$$

3.4.1 Network Structure

A typical architecture is:

- Embedding layer: maps firm ID to e_i .
- Dense layer(s): optionally transform the lagged features $X_{i,t-1}$.
- Concatenation: $[e_i \parallel X_{i,t-1}]$.
- MLP with ReLU or GELU activation.
- Output layer predicting $\hat{g}_{i,t}$.

3.5 Interpretation of Learned Embeddings

Once trained, the embedding vectors $\{e_i\}$ provide a latent representation of long-run firm characteristics. Empirically, these embeddings tend to capture:

- industry clusters,
- differences in scale, profitability, and leverage,
- typical revenue growth regimes,
- structural risk exposures,
- persistent temporal patterns in revenue dynamics.

This multidimensional structure allows the model to generalize more effectively than classical fixed-effect regressions.

3.6 Comparison with Classical Fixed-Effects Models

The embedding-based approach offers several advantages over a traditional linear fixed-effects specification:

1. **Richer representation.** Embeddings encode many latent traits rather than a single intercept.
2. **Flexible nonlinear mapping.** The MLP can capture complex interactions among financial variables.
3. **Better generalization.** Firms with limited historical data borrow statistical strength from similar firms via embedding clustering.
4. **Scalability.** Embeddings avoid the dimensionality explosion of including thousands of firm dummies.

Bibliography

- Vélez-Pareja, I. (2007). *Forecasting Financial Statements with No Plugs and No Circularity*.
- Vélez-Pareja, I. (2009). *Constructing Consistent Financial Planning Models for Valuation*.
- Mejía-Peláez, J., and Vélez-Pareja, I. (2011). *Analytical Solution to the Circularity Problem in the Discounted Cash Flow Valuation Framework*.
- Christodoulou, D., and McLeay, S. (2014). Double-Entry Constraint, Structural Modeling and Econometric Estimation. *Contemporary Accounting Research*.
- Christodoulou, D., and McLeay, S. (2019). The double entry structural constraint on the econometric estimation of accounting variables. *The European Journal of Finance*, 25(10), 853–875.
- Pandher, G. (2007). Modelling and controlling monetary and economic identities with constrained state space models. *International Statistical Review*.
- Angelini, E., Camba-Méndez, G., Giannone, D., Rünstler, G., and Smets, F. (2008). Estimating and forecasting the euro area monthly national accounts from a dynamic factor model. *European Economics: Macroeconomics & Monetary Economics eJournal*.