
NOWCASTING R&D EXPENDITURES: *A Machine Learning Approach*

 **Atin Aboutorabi**
EPFL

atin.aboutorabi@epfl.ch

 **Gaétan de Rassenfosse**
EPFL

gaetan.derassenfosse@epfl.ch

July 17, 2024

ABSTRACT

Macroeconomic data are crucial for monitoring countries' performance and driving policy. However, traditional data acquisition processes are slow, subject to delays, and performed at a low frequency. We address this 'ragged-edge' problem with a two-step framework. The first step is a supervised learning model predicting observed low-frequency figures. We propose a neural-network-based nowcasting model that exploits mixed-frequency, high-dimensional data. The second step uses the elasticities derived from the previous step to interpolate unobserved high-frequency figures. We apply our method to nowcast countries' yearly research and development (R&D) expenditure series. These series are collected through infrequent surveys, making them ideal candidates for this task. We exploit a range of predictors, chiefly Internet search volume data, and document the relevance of these data in improving out-of-sample predictions. Furthermore, we leverage the high frequency of our data to derive monthly estimates of R&D expenditures, which are currently unobserved. We compare our results with those obtained from the classical regression-based and the sparse temporal disaggregation methods. Finally, we validate our results by reporting a strong correlation with monthly R&D employment data.

Keywords Economic Nowcasting · Machine Learning · Google trends · R&D Expenditures.

1 Introduction

Three megatrends are transforming the forecasting literature: the burgeoning of machine learning (ML) methods, the expanding volume of available data, and the increasing computing ability. These trends give rise to new classes of prediction models that can process more and new types of data. Notable examples of new data include satellite data (Henderson et al., 2012; Diebold et al., 2021), textual data (Bollen et al., 2011; De Caigny et al., 2020), and price data from online merchants to forecast inflation (Cavallo and Rigobon, 2016). Combined, these trends lead to more accurate predictions or enable the prediction of new attributes (Zhao and Yang, 2023; Ashtiani and Raahemi, 2023; Shu and Ye, 2023).

Our paper contributes to this research line. We introduce a framework, illustrated in Figure 1, in response to the so-called 'ragged-edge' problem—publication delays of headline variables in official statistics (Mosley et al., 2022), and their low-frequency. This two-step framework includes one step for a supervised learning task to predict observed low-frequency figures (*step A*) and an interpolation step for an unsupervised learning task to estimate unobserved high-frequency figures (*Step B*). The *step A* is a neural network-based model that incorporates low-frequency data and high-frequency data—in the form of web-search data—to nowcast countries' R&D expenditures. As innovation is the prime engine of economic growth (Romer, 1990; Aghion and Howitt, 2008), R&D investments are a central policy metric. However, they are poorly monitored, being based on expensive surveys and released yearly—sometimes even biennially—with a publication lag of two to three years. This delay hinders effective policy-making. As government policies seek to stimulate R&D investments (Edler and Fagerberg, 2017) and institutions track countries' innovation performances (WIPO, 2023), our inability to produce accurate estimates of recent R&D investments is particularly

prejudicial (OECD, 2009). How can countries pilot the EU’s target of reaching R&D investments amounting to 3 percent of GDP by 2030 (European Commission, 2020) if the data have such shortcomings?

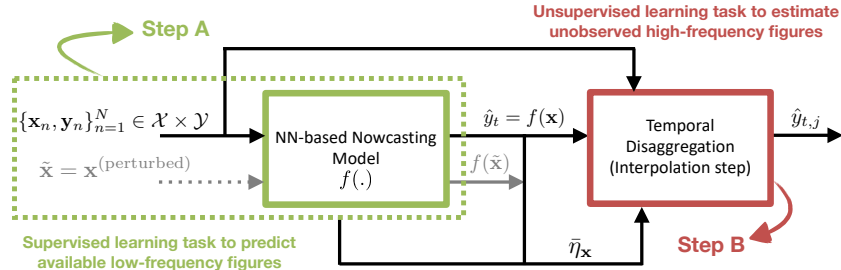
The *step B* involves interpolating the R&D series at a higher frequency than currently available. We exploit the high granularity of our input data to allocate yearly investments into monthly figures. Such data can be used to provide a better understanding of how R&D investments react to economic shocks (OECD, 2009) or policy stimuli (Guellec and van Pottelsberghe, 2003; González and Pazó, 2008; Tajaddini and Gholipour, 2021). The method that we propose is also an important first step towards the nowcasting of monthly series on R&D expenditures.

Our work builds on a handful of studies that have sought to nowcast economic activities using ML methods, including Sokolov-Mladenović et al. (2016), Dai et al. (2017), and Tümer and Akkuş (2018). It also relates to Preis et al. (2013), Borup et al. (2023), and Wołoszko (2020) that exploit web-search data, especially Google trends data, for prediction purposes. Finally, it relates to the literature that has sought to model R&D dynamics (Bloom, 2007) or forecast R&D and innovation activities (Cheng et al., 2005) using traditional econometric techniques. To the best of our knowledge, this study is the first to introduce a neural network-based nowcasting model for R&D expenditures. It also sets a new foundation for advancing the nowcasting of R&D expenditures at a higher frequency. Our method exhibits performance at least as good as existing regression-based temporal disaggregation methods (Chow and Lin, 1971; Mosley et al., 2022).

As the model leverages the web-search data, we first need to identify the relevant search terms. For that purpose, we conscientiously map the actors in the innovation ecosystem by surveying the ‘systems of innovation’ literature (Lundvall, 2010; Edquist, 2013). Next, we run horse race models of nowcasting for yearly R&D series, exploiting both traditional and web-based data. To our surprise, the neural-network-based model exploiting solely web-based data outperforms similar models exploiting solely traditional data, in terms of expressive power (prediction accuracy) and generalization (out-of-sample performance). We then move to high-frequency estimation relying on the output of the neural network model. To be more specific, we propose an interpolation approach, that involves corrupting the input and computing an elasticity value for each input feature based on output distribution. We compare our estimates to the well-known regression-based temporal disaggregation methods. Also, we validate them by using monthly figures for employment in scientific R&D services, reporting significant correlation with our R&D series.

The rest of the paper is organized as follows. Section 2 discusses the background literature. Section 3 outlines the neural network-based nowcasting model for yearly frequency. Section 4 describes the data used for the empirical analysis, and Section 5 discusses the results. Section 6 extends the nowcasting model using temporal interpolation techniques, and reports the results and its respective validations. Finally, Section 7 concludes and points out possible extensions.

Figure 1: Proposed Framework for Nowcasting R&D Expenditures.



2 Background

Macroeconomic variables play a crucial role in guiding global economic decisions. Covering key facets of economic activity, they offer invaluable insights to policymakers and businesses about an economy’s health and direction. Despite their importance, traditional methods of data acquisition have limitations. They rely on expensive surveys and low-frequency updates, potentially obscuring rapidly changing economic conditions.

Acknowledging these limitations, scholars have sought to ‘nowcast’ economic activity. Nowcasting aims to predict the present, the near future, or the recent past, thereby addressing the growing demand for real-time economic insights (Banbura et al., 2010). The first generation of nowcasting models were classical statistical models exploiting macroeconomic variables (Evans, 2005; Giannone et al., 2008; Banbura et al., 2010). For instance, the seminal work by Giannone et al. (2008) proposes a factor model predicting economic activity as captured by GDP. While these models

have pioneered this line of work, their reliance on macroeconomic data, to some extent, defeats their purpose by failing to capture the swiftly changing economic landscape, let alone potential economic shocks.

A subsequent line of inquiry has addressed this limitation by utilizing high-frequency data sources to complement the traditional ones and developing newer, more agile prediction models. Concerning data, scholars have shown that the use of high-frequency Internet search volume, such as Google trends, improves economic forecasting (Choi and Varian, 2009, 2012; Wu and Brynjolfsson, 2015; Ferrara and Simoni, 2019; Götz and Knetsch, 2019).¹

Advances in econometric modeling have seen the rise of (restricted) Mixed Data Sampling (MIDAS) approaches (Ghysels et al., 2004) and unrestricted MIDAS (U-MIDAS) approaches (Foroni et al., 2015), which offer novel ways of incorporating data of varying frequencies into macroeconomic analysis. MIDAS models employ distributed lag-polynomials to integrate high-frequency predictors into a low-frequency modeling framework (Ghysels et al., 2004; Clements and Galvão, 2008; Borup et al., 2023). U-MIDAS models, on the other hand, do not impose a lag-polynomial structure and do not require the alignment of the frequency of predictors and target variables, thereby allowing for more direct inclusion of high-frequency data (Foroni et al., 2015; Borup et al., 2023).

Finally, recent advances in artificial intelligence (AI) have reinforced the shift from classical econometric models to machine learning techniques for forecasting (Borup et al., 2023; Borup and Schütte, 2022; Woloszko, 2023). Building on these developments, we propose a neural-network-based nowcasting model that leverages high-frequency (search volume) data and low-frequency (macroeconomic) data. By constructing different configurations for the input space, we blend the same-frequency data sampling and U-MIDAS approaches.

3 A mixed-frequency neural network-based nowcasting model

3.1 General set-up

We consider a target ($y \in \mathbb{R}$) sampled at a yearly frequency and a vector of predictors ($\mathbf{x} \in \mathbb{R}^d$) containing mixed frequency input, *i.e.*, some variables are sampled monthly and others yearly. The pair (\mathbf{x}, y) represents any given data point (observation) for which the general prediction model is given by:

$$y = f(\mathbf{x}) + \varepsilon$$

in which ε is a zero-mean error term. The target value denotes the R&D expenditures for each country i at year t , and the input vector (\mathbf{x}) consists of four main components. We categorize these components into the following vectors: autoregressive (AR) terms of the target ($\mathbf{Y}_{t-\tau,i} \in \mathbb{R}^{2\tau}$), search volume data, *i.e.* Google trends data ($\mathbf{S}_{t,j,i} \in \mathbb{R}^{(1+\tau) \cdot k_s}$), macroeconomic variables ($\mathbf{Z}_{t-\tau,i} \in \mathbb{R}^{\tau \cdot k_z}$), and the general categorical features ($\mathbf{C}_{j,i} \in \mathbb{R}^{k_c}$), as defined further below. We can re-write the prediction model as:

$$y_{t,i} = f^{(m)}\left(\mathbf{Y}_{t-\tau,i}, \mathbf{S}_{t,j,i}^{(m)}, \mathbf{Z}_{t-\tau,i}, \mathbf{C}_{j,i}^{(m)}; \boldsymbol{\theta}^{(m)}\right) + \varepsilon_{t,j,i}^{(m)}. \quad (1)$$

This model accounts for variations that occur on a monthly basis thanks to the inclusion of monthly-level features in the predictor set, as indicated by superscript (m) . Note also the presence of the superscript on the model $f(\cdot)$ itself, highlighting the fact that the model contains features with a different sampling frequency than the target value. Finally, $\boldsymbol{\theta}$ represents the corresponding vector of model parameters.

We now discuss each of the four components of the model. The vector $\mathbf{Y}_{t-\tau,i}$ includes $\tau > 0$ lagged values of the target to account for potential serial correlation in $y_{t,i}$, without adding any forward-looking (look-ahead) bias:²

$$\mathbf{Y}_{t-\tau,i} = \begin{bmatrix} y_{t-1,i} & \dots & y_{t-\tau,i} & y'_{t-1,i} & \dots & y'_{t-\tau,i} \end{bmatrix}^\top.$$

In addition to the AR terms $y_{t,i}$, we incorporate a binary variable associated with each AR term $y'_{t,i}$, that indicates missing values on the $y_{t,i}$ variable. In case of missing value for any of the AR terms, we estimate it using linear interpolation (see Section 4.1).

¹Google trends is an analytical tool that quantifies search intensities for any given search term(s) by geographical area and over any selected time period. See Section 4.2 for details.

²After having tested various lags, the empirical analysis will consider $\tau = 3$. Note that even with less than $\tau = 3$ and only considering the final state ($\tau = 1$), the out-of-sample performance of the model is not significantly different, showing only minor improvement from $\tau = 1$ to $\tau = 3$ (not reported). However, due to the greater comprehensiveness of considering more lags and in order to evaluate properly the ability of neural networks to handle high-dimensional data, we present the empirical results for $\tau = 3$.

The second vector of predictors, $\mathbf{S}_{t,j,i}^{(m)}$, consists of monthly data obtained from Google trends. The data contain k_s number of different topics associated with a selection of search terms, as we elaborate in Section 4.³ Note the subscript j , which indicates values corresponding to the $(12 - j)^{th}$ month of the year, with $j = 0, \dots, 11$.⁴ Formally, we define $\mathbf{S}_{t,j,i}^{(m)}$ as follows:

$$\mathbf{S}_{t,j,i}^{(m)} = \begin{bmatrix} \bar{\mathbf{S}}_{t-\tau,i} & \bar{\mathbf{S}}_{t,j,i}^{\text{YTD}^{(m)}} \end{bmatrix}^\top.$$

The first sub-vector, $\bar{\mathbf{S}}_{t-\tau,i}$, contains the lagged monthly Google trends values averaged on a yearly frequency for the past τ years and for k_s number of predictors representing various topics associated with the selected search terms. We define:

$$\bar{\mathbf{S}}_{t-\tau,i} = [\bar{s}_{t-1,i} \dots \bar{s}_{t-\tau,i}],$$

$$\bar{s}_{t-\tau,i} = [\bar{s}_{1,t-\tau,i} \dots \bar{s}_{k_s,t-\tau,i}],$$

$$\bar{s}_{k_s,t-\tau,i} = \frac{1}{12} \sum_{j=0}^{11} s_{k_s,t-\tau,j,i}^{(m)}.$$

By aggregating the Google trends data into annual average for each topic $\bar{s}_{k_s,t-\tau,i}$ and constituting $\bar{\mathbf{S}}_{t-\tau,i}$ for all topics, we reduce the monthly variability and emphasize long-term trends. Also, by constructing the time-lagged vector $\bar{\mathbf{S}}_{t-\tau,i}$ for τ lagged values of these averages, we enable the model to capture temporal dynamics across several years and test their associations with the future outcomes.

The second sub-vector of $\mathbf{S}_{t,j,i}^{(m)}$, which is $\bar{\mathbf{S}}_{t,j,i}^{\text{YTD}^{(m)}}$, captures the year-to-date (YTD) Google trends values averaged over the past months of the current year (excluding the current month):

$$\bar{\mathbf{S}}_{t,j,i}^{\text{YTD}^{(m)}} = \begin{bmatrix} \bar{s}_{1,t,j,i}^{\text{YTD}^{(m)}} & \dots & \bar{s}_{k_s,t,j,i}^{\text{YTD}^{(m)}} \end{bmatrix},$$

$$\bar{s}_{k_s,t,j,i}^{\text{YTD}^{(m)}} = \begin{cases} \frac{1}{11-j} \sum_{j=j+1}^{11} s_{k_s,t,j,i}^{(m)} & \forall j \in [0, 10], \\ 0 & \forall j = 11. \end{cases}$$

In short, while the subvector $\bar{\mathbf{S}}_{t-\tau,i}$ captures the information from previous year(s), the subvector $\bar{\mathbf{S}}_{t,j,i}^{\text{YTD}^{(m)}}$ captures the information available in the current year (t) up to the previous month that is being estimated. Indeed, this definition implies that for January ($j = 11$), we have $\bar{s}_{k_s,t,j,i}^{\text{YTD}^{(m)}} = 0$.

The third vector of predictors, $\mathbf{Z}_{t-\tau,i}$, contains the lagged values of k_z different macroeconomic variables. We formally define it as:

$$\mathbf{Z}_{t-\tau,i} = [\mathbf{z}_{t-1,i} \dots \mathbf{z}_{t-\tau,i}]^\top,$$

$$\mathbf{z}_{t-\tau,i} = [z_{1,t-\tau,i} \dots z_{k_z,t-\tau,i}].$$

Finally, the fourth vector of predictors, $\mathbf{C}_{j,i}^{(m)}$, contains features associated with countries and months, which we will discuss further in Section 3.2:

³Google trends provides *related topics* by associating some standardized groups to the search terms. The topics are more reliable than search terms because they include not only the exact search terms, but also their misspellings and acronyms, and they are harmonized over all languages Google News Initiative (2023). As an example, the search term ‘R&D expenditure’ relates to the topics ‘Research and development’, ‘Innovation’, ‘Technology’, ‘Patents’, etc.

⁴Considering the scalars $s_{k_s,t,j,i}$ with j corresponding to the $(12 - j)^{th}$ month of the year t for country i , is equivalent to $s_{k_s,t-j/12,i}$. It would be interpreted as follows: December $s_{k_s,t-0/12,i}$, November $s_{k_s,t-1/12,i}$, ..., June $s_{k_s,t-0.5,i}$, ..., and January $s_{k_s,t-11/12,i}$.

$$\mathbf{C}_{j,i}^{(m)} = \begin{bmatrix} c_{1,i} & c_{2,i} & m_1^{(m)} & \dots & m_{12}^{(m)} \end{bmatrix}^\top.$$

As a baseline, expressing the model in equation (1) considering a classical regression setup, in which the input and output data satisfy a linear relation, yields:

$$y_{t,i} = \mathbf{S}_{t,j,i}^{(m)\top} \boldsymbol{\omega}_S^{(m)} + \mathbf{Y}_{t-\tau,i}^\top \boldsymbol{\omega}_{AR}^{(m)} + \mathbf{Z}_{t-\tau,i}^\top \boldsymbol{\omega}_Z^{(m)} + \mathbf{C}_{j,i}^{(m)\top} \boldsymbol{\omega}_C^{(m)} + \varepsilon_{t,i}^{(m)},$$

and the vector of model parameters is defined as:

$$\boldsymbol{\theta}^{(m)} = \begin{bmatrix} \boldsymbol{\omega}_{AR}^{(m)\top} & \boldsymbol{\omega}_X^{(m)\top} & \boldsymbol{\omega}_Z^{(m)\top} & \boldsymbol{\omega}_C^{(m)\top} \end{bmatrix}^\top.$$

We depart from this linear regression model by developing a neural network architecture from the multilayer perceptron (MLP) class (Goodfellow et al., 2016). In this approach, the function $f(\cdot)$ is represented by a neural network that will learn weights and biases through the course of training for a suitable representation of the data. This problem falls under the class of supervised learning. We observe some data (train set) $S_{\text{train}} = \{\mathbf{x}_n, y_n\}_{n=1}^N \in \mathcal{X} \times \mathcal{Y}$ and given a new \mathbf{x} , we aim to predict its label y , corresponding to R&D expenditures.

Our approach is informed by the expressive power of neural networks—having superior feature learning capabilities compared to traditional regression models. It aligns with Radhakrishnan et al. (2022)’s insights on the importance of feature learning models compared to non-feature learning ones. Also, according to the double-descent risk curve introduced by Belkin et al. (2019), in which they incorporate both classical and modern regimes, rich models such as neural networks that are over-parameterized and considered to be over-fitted, still exhibit high accuracy on out-of-sample data. The tendency of neural networks for better generalization and reduced risk of overfitting compared to under-parameterized classes of model (low-capacity function classes) further motivated our choice.

We adopt a series of configurations for modeling $f(\cdot)$ by varying the input space—and noting that all configurations contain features associated with country and month, $\mathbf{C}_{j,i}^{(m)}$. We start by considering two ‘conventional’ input spaces. The configuration with the minimal dimensional input space includes only autoregressive terms, referred to as *LagRD* model. Then, restricting the set of predictors to macroeconomic variables and autoregressive terms, leads to a configuration labeled as the *Macros* model. Next, we consider an input space that contains exclusively Google trends data and we explore two configurations. Model *AGT* includes the annual Google trends data from previous years, so that $\mathbf{S}_{t,i} = [\bar{\mathbf{S}}_{t-\tau,i}]$. Model *MGT* integrates the yearly data with the monthly Google trends data from the most recent months, that is $\mathbf{S}_{t,i} = [\bar{\mathbf{S}}_{t-\tau,i} \quad \bar{\mathbf{S}}_{t,i}^{\text{YTD}}]$. We then expand the input space by including autoregressive terms of R&D expenditures as additional predictors, hereafter referred to as the *AGTwRD* and *MGTwRD* models, respectively. Finally, we consider the broadest possible input space by combining macroeconomic variables, historical target values, and Google trends data at both yearly and monthly intervals, hereafter called *AllVar*. This model is the general configuration introduced in equation (1). By varying the input space, we will be able to assess the improvements in predictive accuracy that the various pieces of data offer. Table 1 summarizes the seven configurations that we will explore.

Table 1: Summary of Model Configurations

Configuration	Prediction function given an input space
<i>LagRD</i>	$f(\mathbf{x}) := f^{(m)} \left(\mathbf{Y}_{t-\tau,i}, \mathbf{C}_{j,i}^{(m)} \right)$
<i>Macros</i>	$f(\mathbf{x}) := f^{(m)} \left(\mathbf{Y}_{t-\tau,i}, \mathbf{Z}_{t-\tau,i}, \mathbf{C}_{j,i}^{(m)} \right)$
<i>AGT</i>	$f(\mathbf{x}) := f^{(m)} \left(\mathbf{S}_{t,j,i}(\bar{\mathbf{S}}_{t-\tau,i}), \mathbf{C}_{j,i}^{(m)} \right)$
<i>MGT</i>	$f(\mathbf{x}) := f^{(m)} \left(\mathbf{S}_{t,j,i}^{(m)}(\bar{\mathbf{S}}_{t-\tau,i}, \bar{\mathbf{S}}_{t,j,i}^{\text{YTD}^{(m)}}), \mathbf{C}_{j,i}^{(m)} \right)$
<i>AGTwRD</i>	$f(\mathbf{x}) := f^{(m)} \left(\mathbf{Y}_{t-\tau,i}, \mathbf{S}_{t,j,i}(\bar{\mathbf{S}}_{t-\tau,i}), \mathbf{C}_{j,i}^{(m)} \right)$
<i>MGTwRD</i>	$f(\mathbf{x}) := f^{(m)} \left(\mathbf{Y}_{t-\tau,i}, \mathbf{S}_{t,j,i}^{(m)}(\bar{\mathbf{S}}_{t-\tau,i}, \bar{\mathbf{S}}_{t,j,i}^{\text{YTD}^{(m)}}), \mathbf{C}_{j,i}^{(m)} \right)$
<i>AllVar</i>	$f(\mathbf{x}) := f^{(m)} \left(\mathbf{Y}_{t-\tau,i}, \mathbf{S}_{t,j,i}^{(m)}(\bar{\mathbf{S}}_{t-\tau,i}, \bar{\mathbf{S}}_{t,j,i}^{\text{YTD}^{(m)}}), \mathbf{Z}_{t-\tau,i}, \mathbf{C}_{j,i}^{(m)} \right)$

The next section elaborates on the MLP set-up to model the function $f(\cdot)$.

3.2 MLP set-up

As mentioned in Section 3.1, unlike non-feature learning algorithms that require explicit programming for feature extraction, neural networks autonomously learn and utilize features inherent in the data (Radhakrishnan et al., 2022). Also, neural networks have the ability to interpolate data and generalize effectively, while being over-parameterized and considered as high-capacity function classes (Belkin, 2021). In this subsection, we present the neural networks set-up from the MLP class developed for our context.

Mathematically, an MLP architecture can be expressed as a series of function compositions that map an input vector to an output prediction. We can define $f : \mathbb{R}^d \rightarrow \mathbb{R}$ as a fully connected network with L hidden layers for $L > 1$, weight matrices $\{\mathbf{W}^{(l)}\}_{l=1}^{L+1}$, bias vectors $\{b^{(l)}\}_{l=1}^{L+1}$, and activation function ϕ , of the form:

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{h}(\mathbf{x})^\top \mathbf{W}^{(L+1)} + b^{(L+1)}; \\ \mathbf{x}^{(\ell)} &= h^{(\ell)}(\mathbf{x}^{(\ell-1)}) := \phi(\mathbf{W}^{(\ell)\top} \mathbf{x}^{(\ell-1)} + b^{(\ell)}) \quad \text{for } \ell \in \{2, \dots, L\}; \\ \mathbf{h}(\mathbf{x}) &:= \mathbf{x}^{(L)} = h^{(L)}(\mathbf{x}^{(L-1)}) = \phi(\mathbf{W}^{(L)\top} \mathbf{x}^{(L-1)} + b^{(L)}). \end{aligned}$$

with $\mathbf{x}^{(0)} = \mathbf{x}$ for $\ell = 1$, and $\mathbf{h} : \mathbb{R}^d \rightarrow \mathbb{R}^K$. The function $\mathbf{h}(\cdot)$ represents the feature extractor part of the architecture, with activation function ϕ that performs an element-wise non-linear transformation on its input.

The model learns the weight matrices and bias vectors with the objective of minimizing the training loss (also known as empirical risk minimization). The implemented architecture is MLP with a rectified linear unit (ReLU) activation function and $L = 3$, which is a feed-forward type comprising three fully connected hidden layers consisting of 200, 20, and 20 neurons at each layer, respectively.⁵ These dense layers are followed by an output layer with a single neuron, such that we reference the network as ‘NN4.’ We assessed the model’s performance using both the pyramid strategy, as recommended by Masters (1993), and the flatter strategy, *i.e.* same size for all layers. LeCun et al. (2002) points out that using the same number of neurons for all dense layers works at least as well as a pyramid-like setup or an upside-down pyramid. Also, they add that in most cases, the performance of an architecture with an overcomplete first hidden layer is better than an undercomplete one (LeCun et al., 2002).⁶ Lastly, they mention that the set-up choice ultimately depends on the data. Hence, settling on the present set-up is informed by having tested several set-ups for our context. The selected architecture, in fact, has an overcomplete first hidden layer, which aligns with LeCun et al. (2002)’s findings.

In order to diminish the sensitivity of the model’s output to the random initial parametrization, we construct an ensemble of ten such neural networks. In this setup, each neural network is initialized with random parameters, and we take the average of all predictions. This approach stabilizes the model’s performance by combining predictions from multiple neural networks (Woloszko, 2020).

Furthermore, as mentioned in Section 3.1, the vector of predictors $\mathbf{C}_{j,i}^{(m)}$ contains features associated with months and countries. As for months, we implement ‘one-hot encoding,’ treating each month separately as an independent feature. In regard to countries, first we use ‘LabelEncoder’ to encode the country names to transform them into numerical values, and then we map them to vectors—as an ‘embedding layer.’ More precisely, given that we consider an embedding dimension of size two for the countries, the embedding layer maps each unique integer (representing a country) to a two-dimensional continuous vector (*i.e.*, a tensor of rank one). Having countries in vector format allows the neural network to understand latent relationships between countries in terms of the prediction task, which outperforms one-hot encoding methods (Guo and Berkahn, 2016). As training progresses, countries with similar behaviors or characteristics might get vectors that are closer together, indicating their similarity in terms of the target variable the model is predicting (Guo and Berkahn, 2016). The embedding vector of each country is learned over the training process and constitutes the weight matrix associated with the embedding layer.

Moreover, in the context of model specificity, Woloszko (2020) elaborates on the trade-offs between country-specific and cross-country modeling. While the former provides nuanced, country-centric insights, it requires a richer variable set. In contrast, the latter aggregates data across countries, thereby increasing the sample size and improving estimation robustness. This trade-off exemplifies the Bias-Variance trade-off. Although pooling data across countries might introduce certain biases, it significantly diminishes the estimator’s variance. Hence, we consider a cross-country set-up in which we pool all countries together.

⁵The ReLU activation function $\phi(x)$ is defined as $\phi(x) = (x)_+ = \max\{0, x\}$.

⁶An ‘overcomplete’ first hidden layer stands for cases in which the number of neurons in the first hidden layer is larger than the input vector.

In addition, we incorporate batch normalization into the set-up, which is applied post-ReLU activation for each hidden layer. It acts as a stabilizer for neuron activation, ensuring smoother learning, and might fasten the convergence process (Ioffe and Szegedy, 2015; Santurkar et al., 2018).

The model utilizes the AdamW optimizer, which complements the popular Adam optimizer with weight decay. This choice is informed by empirical evidence by Loshchilov and Hutter (2017) suggesting superior generalization capabilities.

Finally, we also implement scenarios with ‘early stopping’ according to the validation set. To check and mitigate the risks of overfitting, we interrupt the training regimen for each MLP in the ensemble if the validation loss stagnates or worsens over a particular number of iterations, defined by the ‘patience’ parameter. Using ‘early stopping’ in scenarios ensures that the results are not driven by overfitting. In support of open scientific research, all codes and datasets used in this study are available in a dedicated GitHub repository.⁷

4 Data

This section provides an overview of the data underpinning the empirical application. The target variable is R&D expenditure, formally known as Gross Domestic Expenditures on Research and Development (GERD). It is published on a yearly basis by the OECD. As explained in Section 3, we generate nowcasting values for yearly R&D using mixed-frequency data, consisting of macro variables and Google trends data. In this study, we focus particularly on training and calibrating the model for eight selected countries out of the top twenty innovative countries, according to WIPO’s Global Innovation Index (GII) (WIPO, 2023). The selected countries in our sample set are as follows: Switzerland, the United States, the United Kingdom, Germany, South Korea, China, Japan, and Canada.⁸

4.1 Gross domestic expenditures on research and development (GERD)

GERD measures the total expenditures (both current and capital) on R&D activities conducted by all organizations within a nation’s territory, including companies, research institutes, universities, and government laboratories (OECD, 2023). It includes R&D funded from external sources but excludes domestic funds designated for R&D activities outside the domestic economy. These data are primarily gathered through surveys conducted by national authorities and then reported to the OECD (OECD, 2015). National authorities follow the international norms set forth by the Frascati Manual (OECD, 2015). The R&D series are expressed in billions of 2015 USD PPPs to ensure comparability over time and across countries.

The GERD series have data gaps for certain countries, such as Switzerland. We deal with data gaps by estimating the missing values using linear interpolation, thereby guaranteeing the integrity and continuity of the time series data. We also introduce a binary indicator variable $y'_{t,i}$ that takes value 1 when the original R&D figure is missing and 0 otherwise. However, note that we use the interpolated figures only for the lagged predictor variables and not for the target variable we aim to predict. We exclude target variables with missing values from the learning phase to ensure the reliability of our predictive model.

4.2 Google trends

Google trends, developed by Google, measures search intensities for any given queried keyword(s) by geographical area and over any selected time period starting in 2004 (Google News Initiative, 2023). The tool provides normalized data on the relative search volume of queries. Each data point in Google trends is normalized by the total search volume of the geography and time range considered, on a scale between 0 and 100. The search volume index (SVI) for a given search term or topic k_s is denoted as $s_{k_s,t,j,i}^{(m)}$ and defined as follows:

$$s_{k_s,t,j,i}^{(m)} = \frac{SV_{k_s,t,j,i}^{(m)}}{TSV_{t,j,i}^{(m)}} * c_{i,k_s} \in [0, 100]$$

in which $SV_{k_s,t,j,i}^{(m)}$ and $TSV_{t,j,i}^{(m)}$ stand for volume of searches for k_s in country i at a given time (year t and month j), and total number of searches in country i at the selected time range, *i.e.* from 2004 to the given time (year t and month j).

⁷The dedicated GitHub repository is available at https://github.com/AtiinA1/Nowcasting_RD_Expenditures.git.

⁸Selected countries in our sample are arbitrary selection from different continents, based on the top twenty innovative countries according to WIPO’s GI.

j), respectively. The SVI is normalized and also indexed on a scale of 0 to 100 by the constant c_{i,k_s} for a given time series $s_{k_s,t,j,i}^{(m)}$.

We use Google trends data to predict R&D expenditures dynamics with a fine level of granularity. To keep the relevant search terms that might correlate with R&D activities, we start by reconstructing the ecosystem of stakeholders associated with R&D activities. This ecosystem comprises entities that either directly perform R&D or play a supporting role in the broader R&D landscape. We identify the various stakeholders and the specific search terms by studying the ‘systems of innovation’ literature (Edquist, 2013; Lundvall, 2010). Table 2 presents the ecosystem developed for this study, containing the finalized selection of stakeholders and the corresponding search terms that we identified.⁹

Table 2: Stakeholders and their respective search terms for R&D expenditure.

Stakeholder	Search Terms
Businesses	R&D Expenditure, Product Development
Consulting Firms	Innovation Management, Innovation Strategy
Government Agencies	Government Grants, Research Funding
Innovation Hubs	Startup Incubation, Technology Park
Patent Attorneys	Patent Attorney, Patent Registration
R&D Employees	R&D Jobs
Research Institutions	Collaboration with Industry, Research Grant
Tax Authorities	R&D Tax Credit
Venture Capitalists (VCs)	VC Investment, Startup Funding

After identifying specific search terms, our method involves extracting the related topics from Google trends. While the primary use of Google trends is often to search for individual search terms, it is not limited to that. In order to have the big picture, the data are also grouped into topics and we can retrieve the topic(s) associated with a given search term. The use of topics allows us to address the challenges posed by language-specific search terms and to mitigate any potential ambiguity inherent in relying only on individual search term searches (Woloszko, 2020). Consequently, the use of topics ensures that our analysis remains comparable across countries with diverse linguistic backgrounds, as these identifiers represent standardized, language-neutral markers of user interest. Nevertheless, the selection of relevant topics in Google trends, given the absence of a fixed topic list, demands careful exploration. We build upon the specific search terms (Table 2) as anchors to extract associated topic identifiers from Google trends. In the end, we capture 57 different Google trends topics, such that, $k_s = 57$ in the setting of this problem.

One feature of Google trends is that search indices are derived from a sample of total search volume, for computational tractability reasons. Indices that represent low-volume searches may exhibit significant sampling variance (Woloszko, 2020). Accordingly, in order to minimize sampling variance associated with Google’s methods, we retrieve five samples per query and average search volume index across the samples for each topic to obtain a more reliable time series (Woloszko, 2020; Medeiros and Pires, 2021; Woloszko, 2023).

4.3 Macroeconomic variables

The set of predictors also contains six ‘traditional’ macroeconomic variables capturing elements of countries’ wealth, economic conditions, and competitiveness ($k_z = 6$). These variables, which correlate with R&D expenditures, include:

- Gross domestic product per capita: Reflects the wealth and economic strength of an economy.
- Unemployment rate: Indicates the health of the labor market and broader economic conditions.
- Population: Used as a scale variable, represents the size of an economy.
- Inflation rate: Represents the rate of price growth and indicates the state of the economy.
- Export and import volumes: Chosen as indicators of a nation’s competitiveness and its engagement with the global economy.

We sourced these variables from the International Monetary Fund (IMF)’s World Economic Outlook Database, April 2023 release (International Monetary Fund, 2023). We filled missing values with the mean of the respective variable for that country, calculated over all available years.

⁹The initial ecosystem developed for this study encompassed a broader range of search terms, detailed in the appendix, Table A.1.

5 Results

5.1 Out-of-sample predictions

As explained earlier, we present the performance of the MLP setup over different configurations of the input space. By varying the input space, we aim to understand the value of the Google trends data compared to traditional economic variables.

Figures 2a and 2b respectively present box plots of the Root Mean Square Errors (RMSE) and Mean Absolute Percentage Errors (MAPE) of the MLP predictions for each of the seven input spaces considered. The figures also present RMSE and MAPE for corresponding OLS regression models as benchmarks. Two remarks are in order.

First, the neural network predictions generally outperform the linear regression results, having lower MAPE, with the notable exception of the parsimonious *LagRD* model, which exhibits similar performances. This exception can be explained by the low dimensionality of the input space (*i.e.*, smaller input space cardinality), resulting in a low signal-to-noise ratio and underfitting of the model. Focusing for a moment on the linear regression model, the scale of metrics values for the *AGT* and *MGT* configurations have the highest median values and the broadest interquartile range. All other configurations exhibit very similar median values and interquartile ranges. The common point between all these configurations is the presence of AR terms (lagged target values), which have the highest coefficients and dominantly drive the prediction outcome. Turning now to the neural network models, it is apparent that they perform better than the linear regression models. They perform particularly well with a large number of predictors, reflecting their ability to capture complex, non-linear relationships between the input features and the target variable—relationships that linear models are inherently unable to represent.

Second, the model exploiting the annual Google trends data as input (*AGT*) using a neural network exhibits the best predictive accuracy on average. The RMSE and MAPE are the lowest of all configurations—even lower than the models that also integrate lagged R&D data. The intuition behind this performance is that there are too few data points available on a yearly basis to train the model properly and at the same time too much degree-of-freedom associated with those data points available due to the cardinality of the input space, compared to other configurations involving the Google trends data. This intuition also explains the higher variance in the MAPE for *AGT* compared to *AGTwRD*. Using as analogy a classification problem with data points in the input space, it means that we have too many hyperplanes to classify them. While we may achieve perfect classification on our training data, the performance can vary significantly on diverse test sets, indicating high variance. Therefore, relying solely on the performance of *AGT* can be misleading, as it signals overfitting. It is more reliable and informative to consider *AGT* with respect to all configurations.

Overall, the observations of Figure 2 highlight the potential of neural networks and high-frequency data to capture complex dynamics and generate insightful forecasts. This can be interpreted from the outperformance of the configurations with search volume data included, compared to those without any high-frequency data in neural network set-up.

Figure 2: Comparison of RMSE and MAPE Values for Different Configurations.

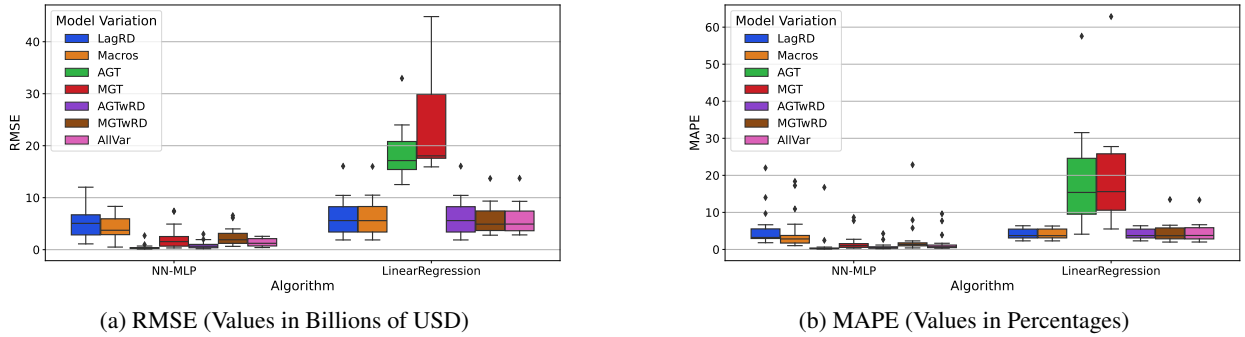
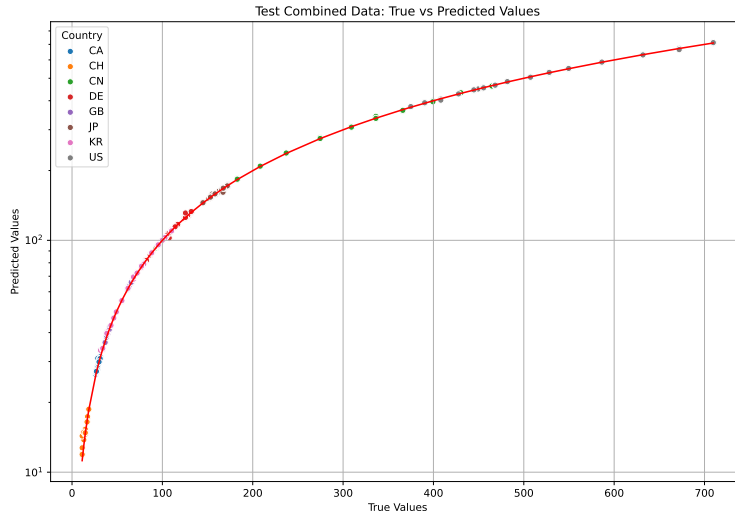


Figure 3 depicts the out-of-sample performance of the *AllVar* neural-network prediction model. We focus on this model because it exploits all the available information. The figure suggests that the model's predictions are fairly consistent across countries. Appendix C reports the out-of-sample performance for the individual countries. Zooming into countries suggests the presence of some outlier predictions, which could arise from country-specific factors, and the lack of representativeness of the Google trends data for those countries, either in the early days of Google trends or in general.

Regarding Switzerland, the R&D expenditures data are reported biennially since 2015 and prior to that, only once every four years. This irregular frequency and sparse data availability can hinder the model’s capability to accurately predict for Switzerland. Increasing the frequency and consistency of data points could potentially improve the model’s performance. In the case of China, the influence of government censorship, commonly referred to as the Great Firewall, has significantly restricted access to Google. This limitation affects the representativeness and reliability of Google trends data, as a substantial portion of internet searches can be routed through local search engines, *e.g.* Baidu (StatCounter, 2024a). For countries like Japan and South Korea, Google competes with other search engines that are more integrated with the respective languages and cultural contexts, *e.g.* Yahoo! in Japan (StatCounter, 2024b) and Naver in South Korea (StatCounter, 2024c). This competition affects the volume and characteristics of data available from Google trends in these countries and how comprehensive the view of search behaviors are. Lastly, a common factor affecting all models’ performance for all countries is the varying rate of Google’s adoption over time, which directly influences the values obtained from Google trends. If these factors can be effectively addressed, it may enhance the model’s accuracy even more, leading to more reliable and robust outcomes across different countries.

Figure 3: True vs. Predicted R&D Expenditures in USD (bn) for Different Countries in *AllVar* Configuration.

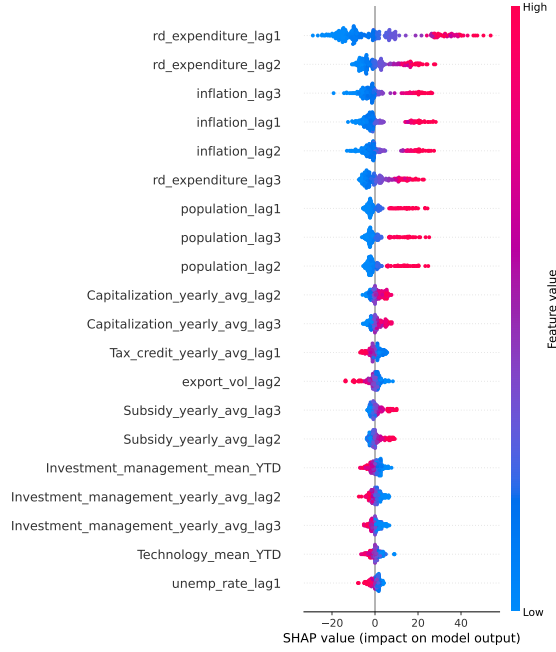


5.2 Global interpretability of model: SHAP

In response to the tension between accuracy and interpretability inherent in complex ML models, Lundberg and Lee (2017) offers a unified framework for interpreting feature importance, so-called SHAP (SHapley Additive exPlanations).

SHAP values have become a standard tool to obtain both local and global interpretability (Lundberg and Lee, 2017; Woloszko, 2020). They enable local interpretability by showing how a model makes a decision for an individual prediction. They also enable global interpretability by explaining the model’s behavior across the entire sample. In this regard, we implement the recent model-agnostic approximation method, *Kernel SHAP*, introduced by Lundberg and Lee (2017). This method applies to any ML model, including deep learning models like ours. We use k -means to sample our data based on k clusters, for which we considered five representative points per country and in total $k = 40$. The chosen method ensures that the sample is representative of the larger dataset.

The SHAP summary plot in Figure 4 illustrates the top 20 most important features and their contributions to the prediction. It is evident that the contributions are not confined to a few dominant features; rather, a wide array of features, particularly Google trends features, collectively contribute to high-quality predictions. This result is a manifestation of the ‘Illusion of Sparsity,’ which posits that while a small number of variables may be identified as primary predictors in a prediction model, a broad range of features could still collectively exhibit significant influence (Giannone et al., 2021). It highlights the fact that we are dealing with a *dense*-modelling problem, in which the target value is recovered by many features with small contributions (Giannone et al., 2021; Woloszko, 2023). Also, it validates the choice of neural networks, as part of *dense* class of estimators, to address non-linearities between input features and the output (Woloszko, 2023). In addition, Figure 4 reveal that the model effectively learns the importance of historical data on R&D expenditures, as captured in autoregressive terms, and also integrates signals from other macroeconomic variables.

Figure 4: SHAP Summary Plot for *AllVar* Configuration.

6 Interpolating R&D expenditures at a higher frequency

We have established so far that Google trends data are well suited to the nowcasting of yearly R&D expenditures. Another advantage of the data lies in the fact that they are available at a high frequency, making it possible to nowcast monthly R&D expenditures. We take a step in this direction by interpolating R&D expenditures, *i.e.* applying temporal disaggregation, which involves using high-frequency indicators to construct high-frequency renditions of low-frequency information (Mosley et al., 2022).¹⁰

We interpolate the series using three methods. We start with the classical regression-based temporal disaggregation method introduced by Chow and Lin (1971), followed by a recent extension for high-dimensional settings proposed by Mosley et al. (2022). We then propose a full neural-network-based approach.

In order to construct the unobserved monthly time series $y_{t,j,i}^{(m)}$, Chow and Lin (1971) and Mosley et al. (2022) assume the following regression model at the monthly frequency (Mosley et al., 2022, Equation 1):

$$y_{t,j,i} = \mathbf{x}_{t,j,i}^{(m)\top} \boldsymbol{\theta} + \mathbf{u}_{t,j,i}^{(m)}. \quad (2)$$

The residual vector $\mathbf{u}_{t,j,i}^{(m)}$ is mean-zero, has covariance matrix \mathbf{V}_m , and follows a first-order autoregressive process $\mathbf{u}_{t,j,i}^{(m)} = \rho \mathbf{u}_{t,j-1,i}^{(m)} + \varepsilon_{t,j,i}^{(m)}$ with $\varepsilon_{t,j,i}^{(m)} \sim \mathcal{N}(0, \sigma^2)$ and $|\rho| < 1$. Due to $y_{t,j,i}$ being unobserved, the counterpart of equation (2) is written as (Mosley et al., 2022, Equation 3)

$$y_{t,i} = \mathbf{x}_{t,i}^\top \boldsymbol{\theta} + \mathbf{u}_{t,i}. \quad (3)$$

We obtain it by multiplying equation (2) by an aggregation matrix $\mathbf{A} = \mathbf{I}_n \otimes \mathbf{1}_{12}$. In our context, the twelve monthly figures in a year must sum to their corresponding yearly figure, without loss of generality (Mosley et al., 2022). Similarly, we transform the covariance matrix \mathbf{V}_m to \mathbf{V}_a using the aggregation matrix and its transpose.

Given the yearly figures being observed, we can solve equation (3) using standard techniques, for which Chow and Lin (1971) establishes a generalized least squares (GLS) estimator for $\boldsymbol{\theta}$. This estimator can be derived from the

¹⁰The task of nowcasting monthly R&D expenditures differs from the task of interpolating yearly series into monthly figures in an important way. In the latter task, we apply temporal disaggregation on the yearly amount of R&D expenditures, by factoring in the records of full year monthly Google trends data.

Chow-Lin cost function term in equation (4) and it depends on the unknown parameter ρ , which needs to be estimated (Chow and Lin, 1971, Equation 15). In an extension suited to high-dimensional data, Mosley et al. (2022) proposes the estimator as the form in equation (4) (Mosley et al., 2022, Equation 7). They incorporate a penalty function with the classic Chow-Lin cost function to regularize and encode the sparsity assumption in high-dimensional settings. As for the regularizer term, they consider LASSO (l_1) penalty $P_\lambda(\boldsymbol{\theta}) = \lambda \|\boldsymbol{\theta}\|_1 := \lambda_\rho \sum_{j=1}^d |\boldsymbol{\theta}_j|$, and refer to this method as l_1 -spTD. The index λ is a non-negative regularization parameter that controls the degree of shrinkage.

$$\hat{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}_\rho = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \left\{ \underbrace{\left\| \mathbf{V}_a^{-1/2} (y_{t,i} - \mathbf{x}_{t,i}^\top \boldsymbol{\theta}) \right\|_2^2}_{\text{Chow-Lin cost function}} + \underbrace{P_\lambda(\boldsymbol{\theta})}_{\text{Regularizer}} \right\}. \quad (4)$$

We implement the classical regression-based temporal disaggregation method proposed by Chow and Lin (1971), in particular the ‘chow-lin-maxlog’ method using the ‘tempdisagg’ R package by Sax and Steiner (2013). As for the sparse temporal disaggregation method by Mosley et al. (2022), we use the ‘DisaggregateTS’ R package by Mosley and Nobari (2022).

While these regression-based estimates are well-accepted temporal disaggregation methods, they may fail to capture the inherent characteristics of the data. Accordingly, we propose a new method that leverages our neural-network model.

In broad terms, we corrupt the input and record its corresponding output, allowing us to derive a neural network-driven elasticity value for each feature of the AGT configuration. Considering the elasticities for all input features, in addition to the proportion of each input feature at the monthly level with respect to its aggregated value at the yearly level, we distribute the yearly figures of R&D expenditures to monthly figures, labeled $\hat{y}_{t,j,i}$.¹¹

More formally, we calculate the *empirical predicted expected elasticity* $\mathbb{E}[\hat{\eta}_{s_{k_s},i}]$, for a particular input feature, *i.e.* topic s_{k_s} , and a given country i , from the prediction model on the observations (train set) as in equation (5). $\mathbb{E}_{S_{\text{train}}}$ denotes the expectation over the samples in the training set, while $\mathbb{E}_{\delta \sim \Delta}$ represents the expectation over the perturbations δ , which are sampled from the distribution Δ . In our setting, we apply perturbations to the input features by sampling from a normal distribution with a mean of 0.01 and a standard deviation of 0.005.

$$\bar{\eta}_{s_{k_s},i} = \mathbb{E}_{S_{\text{train}}} [\hat{\eta}_{s_{k_s},i}] = \mathbb{E}_{S_{\text{train}}} [\mathbb{E}_{\delta \sim \Delta} [\eta_{s_{k_s},t-\tau,i}]] = \mathbb{E}_{S_{\text{train}}} \left[\mathbb{E}_{\delta \sim \Delta} \left[\frac{\frac{y_{t,i}^\delta - y_{t,i}}{y_{t,i}}}{\frac{\bar{s}_{k_s,t-\tau,i}^\delta - \bar{s}_{k_s,t-\tau,i}}{\bar{s}_{k_s,t-\tau,i}}} \right] \right]. \quad (5)$$

Letting $p_{s_{k_s},t,j,i}^{(m)}$ denote the proportion of Google trends value for a topic associated with search terms k_s on a monthly frequency, with respect to its aggregated value on a yearly frequency of year t , for month j , and country i , we have:

$$p_{s_{k_s},t,j,i}^{(m)} = \frac{s_{k_s,t,j,i}^{(m)}}{\sum_{j=0}^{11} s_{k_s,t,j,i}^{(m)}}.$$

Given the estimator for the expected elasticity value $\hat{\eta}_{s_{k_s},i}$, the adjusted proportion of the Google trends value of year t for month j and country i can be defined as:

$$\tilde{p}_{s_{k_s},t,j,i}^{(m)} = \bar{\eta}_{s_{k_s},i} \times p_{s_{k_s},t,j,i}^{(m)}.$$

As a result, estimated monthly R&D expenditures of year t for month j and country i is given by:¹²

$$\hat{y}_{t,j,i}^{(m)} = y_{t,i} \times \frac{\sum_{k_s} \tilde{p}_{s_{k_s},t,j,i}^{(m)}}{\sum_j \sum_{k_s} \tilde{p}_{s_{k_s},t,j,i}^{(m)}} \quad \text{s.t.} \quad y_{t,i} = \sum_{j=0}^{11} \hat{y}_{t,j,i}^{(m)}.$$

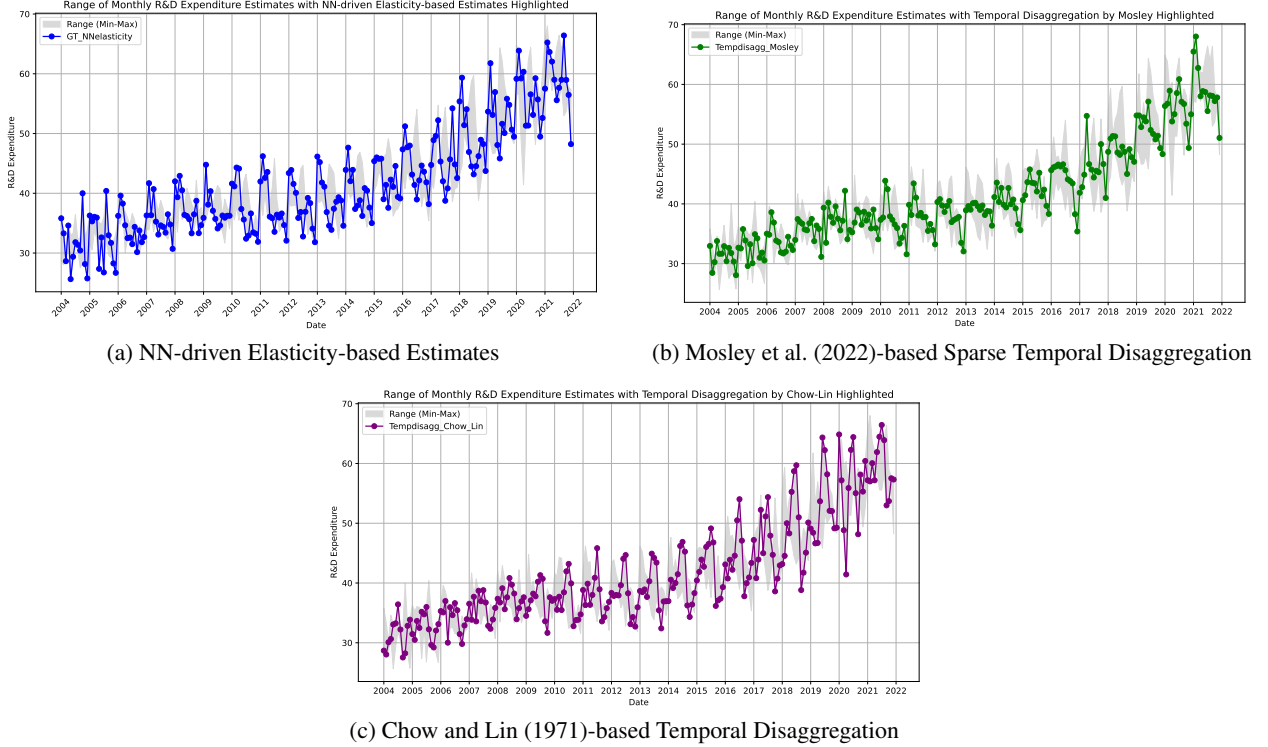
Figure 5 depicts different estimated monthly R&D expenditures over time for the United States. The neural-network-driven elasticity-based estimates in Figure 5a strongly correlate with the sparse temporal disaggregation

¹¹Appendix D proposes an alternative method based on the corrupted input approach. While this method can be very intuitive from computer science perspective, it does not perform as well as the proposed temporal disaggregation method.

¹²In case of unavailability of the original data point $y_{t,i}$, we can estimate it using the neural network-based nowcasting model $\hat{y}_{t,i} = f(\mathbf{x})$.

estimates (Mosley et al., 2022) shown in Figure 5b, with correlation coefficient $r(N - 2) = r(214) = 0.93$, $p < 0.001$ in level, and $r(214) = 0.66$, $p < 0.001$ on the growth rates. While these two estimators exploit full information by allowing for high-dimensional indicator matrices, the classical temporal disaggregation method of Chow and Lin (1971) requires a dimensionality-reduction step to limit the set of predictors to a handful of indicators. In this regard, we leverage the insights generated by the neural networks model, in particular, the top six important features in terms of their contribution to the prediction based on the shapley values in the AGT configuration. The correlation between the resulting temporal disaggregation shown in Figure 5c and the NN-driven elasticity-based estimates is $r(214) = 0.73$, $p < 0.001$ in level, and $r(214) = -0.17$, $p = 0.02$ on the growth rates which is not statistically significantly different from 0.¹³

Figure 5: Comparison of Estimated Monthly R&D Expenditures for the United States.

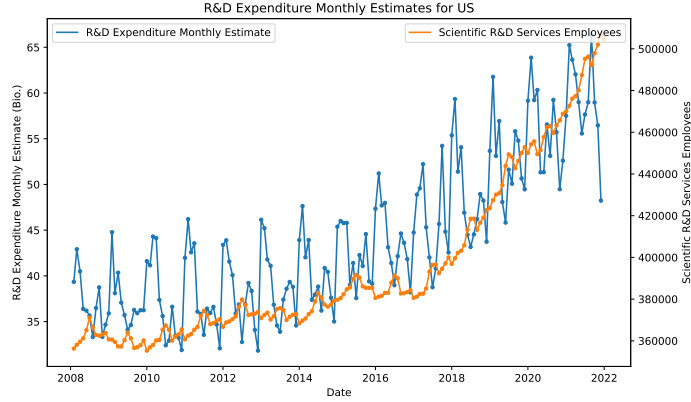


In a second validation exercise, we correlate the monthly series to another series that we expect to relate strongly to R&D expenditures. We use monthly figures for employment in scientific R&D services in the United States from Data USA (Data USA, 2023), as shown in Figure 6. A priori, it is unclear how monthly R&D expenditures should correlate with monthly R&D employment. The time lag between R&D expenditures and employment is ambiguous because R&D expenditures include wages as well as consumables and equipment costs. An increase in employment at time $t = 0$ could result in a subsequent rise in R&D expenditures at time $t + x$, as new employees begin their research, thereby increasing future non-employee expenditures for equipment and consumables. Conversely, the purchase of expensive equipment might occur before the hiring of engineers and scientists needed to operate it. Consequently, the time lags between these variables are not straightforward. Furthermore, the correlation could even be negative or null depending on the degree of substitution between human capital and physical capital.

We compute the correlation between estimated R&D expenditures (in growth rates) and scientific R&D services employees (in growth rates) for different lags and different disaggregation methods, as shown in Table 3. We report the correlation coefficients for lags statistically significantly different from zero, specifically with a p-value less than 0.01. A positive lag means a backward shift in the R&D expenditures growth rate, *i.e.*, capturing how the R&D expenditures growth rate from a previous period (month) correlates to the current R&D services employees growth rate. Conversely, a negative lag value indicates a forward shift, comparing the future R&D expenditures growth rate with the current employees' growth rates.

¹³The correlation between the temporal disaggregation shown in Figure 5c and Figure 5b is $r(214) = 0.87$, $p < 0.001$ in level, and $r(214) = 0.25$, $p = 0.0011$ on the growth rates.

Figure 6: Estimated Monthly R&D Expenditures vs. Employees in Scientific R&D Services.



We observe significant correlations at various lags across all three methods, indicating a relationship between R&D expenditure and employee growth. The estimates obtained using Mosley et al. (2022) seem to be captured by the neural-network-based estimates, with roughly similar correlation coefficients at lags -5, -4, and 4. The estimates based on Chow and Lin (1971) exhibit a strong positive correlation coefficient at lag 0 and negative at lags -3 and 3. The correlation coefficient at lag 5 is similar to the neural-network-based estimates.

It is worth emphasizing that all methods are based on the predicted features of our neural network model. Chow and Lin (1971) exploits only a handful of the features and exhibit correlation coefficients that are markedly different from the other two methods, which exploit the full set of features. We have no way of establishing which method delivers the most ‘valid’ results. Lacking ground truth, we cannot rule out or favor a method. The fact that all methods significantly correlate with the R&D series indicates that they all capture meaningful aspects of the R&D dynamics.¹⁴

Table 3: Correlation between Monthly R&D Expenditure and R&D Services Employees Growths

Monthly R&D Expenditure Estimates Growth	Lag	Correlation	P-Value
Neural Network-driven elasticity-based estimates	0	-0.42	<0.0001
	-5	-0.39	<0.0001
	5	0.29	0.0002
	-2	0.27	0.0004
	4	0.21	0.0080
	-4	-0.20	0.0097
Mosley et al. (2022)-based estimates	-5	-0.32	<0.0001
	-4	-0.25	0.0010
	4	0.22	0.0049
Chow and Lin (1971)-based estimates	0	0.46	<0.0001
	-2	-0.37	<0.0001
	1	0.31	<0.0001
	3	-0.25	0.0010
	-3	-0.22	0.0050
	5	0.20	0.0090

7 Conclusion and future work

In this study, we propose a nowcasting model to predict annual gross domestic expenditures on R&D. We develop an MLP set-up and leverage a large set of covariates, including macroeconomic variables and Internet search volume data. Our choice of a neural network prediction model is based on these models’ superior ability in feature learning and

¹⁴Data on patent filings are available at high frequency, but do not relate to current R&D expenditures since the capture the output of the R&D process. The gap between R&D expenditures and patent filings is a least one year, as documented by de Rassenfosse and Jaffe (2018). Furthermore, patent statistics are published 18 months after filing, making it a poor candidate for nowcasting.

their tendency for better generalization and, accordingly, a reduced risk of overfitting. The empirical evidence from our research highlights the existence of non-linear mapping between the covariates and our target. Indeed, the neural network models not only outperform traditional linear regression models but also demonstrate their ability to capture the complex interplay of features. This finding also validates that traditional linear models often fail to perform well when considering high-dimensional settings.

Besides predicting annual R&D expenditures, we explore the feasibility of offering higher-frequency R&D series. We leverage the neural-network-based annual nowcasting model and integrate additional steps to produce a distribution of monthly R&D expenditures. More specifically, by perturbing the input and analyzing its corresponding output, we derive neural-network-driven elasticity values for each feature, that combined with the monthly proportion of each feature relative to its annual aggregate, enable the distribution of yearly R&D expenditures to monthly figures. The results from this extension indicate the model's potential to provide more frequent and detailed insights on R&D. It offers a step in the direction of a monthly nowcasting model. However, lack of ground truth data prevented us from recommending one temporal disaggregation method over the others. More studies are needed to evaluate the robustness, accuracy, and responsiveness of the model, particularly regarding its ability to adapt and respond to economic shocks.

Finally, we hope that policymakers will adopt our model to produce R&D statistics and to inform their decisions. A valuable feature of the framework is that one can adapt it at different levels, *e.g.* regional level, to offer more granular insights. Furthermore, the method can be adapted to capture a different set of topics that go beyond technological innovation reflected by R&D expenditures. This method's potential to track, say, social innovation (Mulgan et al., 2007) or open innovation (Chesbrough, 2003) is a particularly exciting extension.

References

- Aghion, P. and Howitt, P. W. (2008). *The economics of growth*. MIT press.
- Ashtiani, M. N. and Raahemi, B. (2023). News-based intelligent prediction of financial markets using text mining and machine learning: A systematic literature review. *Expert Systems with Applications*, 217:119509.
- Banbura, M., Giannone, D., and Reichlin, L. (2010). Nowcasting. *ECB working paper*.
- Belkin, M. (2021). Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation. *Acta Numerica*, 30:203–248.
- Belkin, M., Hsu, D., Ma, S., and Mandal, S. (2019). Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854.
- Bloom, N. (2007). Uncertainty and the dynamics of R&D. *American Economic Review*, 97(2):250–255.
- Bollen, J., Mao, H., and Zeng, X. (2011). Twitter mood predicts the stock market. *Journal of Computational Science*, 2(1):1–8.
- Borup, D., Rapach, D. E., and Schütte, E. C. M. (2023). Mixed-frequency machine learning: Nowcasting and backcasting weekly initial claims with daily internet search volume data. *International Journal of Forecasting*, 39(3):1122–1144.
- Borup, D. and Schütte, E. C. M. (2022). In search of a job: Forecasting employment growth using Google trends. *Journal of Business & Economic Statistics*, 40(1):186–200.
- Cavallo, A. and Rigobon, R. (2016). The billion prices project: Using online prices for measurement and research. *Journal of Economic Perspectives*, 30(2):151–178.
- Cheng, B., He, R., Yang, H., and Yang, J. (2005). Quantitative method and model for forecasting R&D expenditures in China. *Research Evaluation*, 14(1):51–56.
- Chesbrough, H. W. (2003). *Open innovation: The new imperative for creating and profiting from technology*. Harvard Business Press.
- Choi, H. and Varian, H. (2009). Predicting initial claims for unemployment benefits. *Google Inc*, 1(2009):1–5.
- Choi, H. and Varian, H. (2012). Predicting the present with google trends. *Economic Record*, 88:2–9.
- Chow, G. C. and Lin, A.-I. (1971). Best linear unbiased interpolation, distribution, and extrapolation of time series by related series. *The Review of Economics and Statistics*, pages 372–375.
- Clements, M. P. and Galvão, A. B. (2008). Macroeconomic forecasting with mixed-frequency data: Forecasting output growth in the united states. *Journal of Business & Economic Statistics*, 26(4):546–554.
- Dai, Z., Hu, Y., and Zhao, G. (2017). The suitability of different nighttime light data for GDP estimation at different spatial scales and regional levels. *Sustainability*, 9(2).

- Data USA (2023). Scientific research & development services. <https://datausa.io/profile/naics/scientific-research-development-services#:~:text=As%20of%20February%202023%2C%20there,when%20compared%20to%20February%202022>. Accessed on 10 July 2024.
- De Caigny, A., Coussement, K., De Bock, K. W., and Lessmann, S. (2020). Incorporating textual information in customer churn prediction models based on a convolutional neural network. *International Journal of Forecasting*, 36(4):1563–1578.
- de Rassenfosse, G. and Jaffe, A. B. (2018). Econometric evidence on the depreciation of innovations. *European Economic Review*, 101:625–642.
- Diebold, F. X., Göbel, M., Goulet Coulombe, P., Rudebusch, G. D., and Zhang, B. (2021). Optimal combination of arctic sea ice extent measures: A dynamic factor modeling approach. *International Journal of Forecasting*, 37(4):1509–1519.
- Edler, J. and Fagerberg, J. (2017). Innovation policy: what, why, and how. *Oxford Review of Economic Policy*, 33(1):2–23.
- Edquist, C. (2013). *Systems of Innovation: Technologies, Institutions and Organizations*. Routledge.
- European Commission (2020). Aiming for more: R&D investment scenarios for the next decade. Technical report, European Commission, Directorate-General for Research and Innovation. ISBN 978-92-76-19834-5, DOI: 10.2777/85096.
- Evans, M. (2005). Where are we now? real-time estimates of the macro economy.
- Ferrara, L. and Simoni, A. (2019). When are Google data useful to nowcast GDP? an approach via pre-selection and shrinkage. Working Paper 717, Banque de France.
- Foroni, C., Marcellino, M., and Schumacher, C. (2015). Unrestricted mixed data sampling (midas): Midas regressions with unrestricted lag polynomials. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 178(1):57–82.
- Ghysels, E., Santa-Clara, P., and Valkanov, R. (2004). The midas touch: Mixed data sampling regression models. *UCLA: Finance*.
- Giannone, D., Lenza, M., and Primiceri, G. E. (2021). Economic predictions with big data: The illusion of sparsity. *Econometrica*, 89(5):2409–2437.
- Giannone, D., Reichlin, L., and Small, D. (2008). Nowcasting: The real-time informational content of macroeconomic data. *Journal of Monetary Economics*, 55(4):665–676.
- González, X. and Pazó, C. (2008). Do public subsidies stimulate private R&D spending? *Research Policy*, 37(3):371–389.
- Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep learning*. MIT press.
- Google News Initiative (2023). Understanding Google trends data. <https://newsinitiative.withgoogle.com/resources/trainings/fundamentals/google-trends-understanding-the-data/>. Accessed on 27 September 2023.
- Götz, T. B. and Knetsch, T. A. (2019). Google data in bridge equation models for german GDP. *International Journal of Forecasting*, 35(1):45–66.
- Guellec, D. and van Pottelsberghe, B. (2003). The impact of public R&D expenditure on business R&D. *Economics of Innovation and New Technology*, 12(3):225–243.
- Guo, C. and Berkhahn, F. (2016). Entity embeddings of categorical variables. *arXiv preprint arXiv:1604.06737*.
- Henderson, J. V., Storeygard, A., and Weil, D. N. (2012). Measuring economic growth from outer space. *American Economic Review*, 102(2):994–1028.
- International Monetary Fund (2023). World economic outlook database. <https://www.imf.org/en/Publications/WEQ/weo-database/2023/April/download-entire-database>. Accessed on 20 June 2023.
- Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448–456. PMLR.
- LeCun, Y., Bottou, L., Orr, G. B., and Müller, K.-R. (2002). Efficient backprop. In *Neural networks: Tricks of the trade*, pages 9–50. Springer.
- Loshchilov, I. and Hutter, F. (2017). Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.
- Lundberg, S. M. and Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30.

- Lundvall, B.-A. (2010). *National systems of innovation: towards a theory of innovation and interactive learning*, volume 2. Anthem Press.
- Masters, T. (1993). *Practical neural network recipes in C++*. Morgan Kaufmann.
- Medeiros, M. C. and Pires, H. F. (2021). The proper use of Google trends in forecasting models. *arXiv preprint arXiv:2104.03065*.
- Mosley, L., Eckley, I. A., and Gibberd, A. (2022). Sparse temporal disaggregation. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 185(4):2203–2233.
- Mosley, L. and Nobari, K. S. (2022). *High-Dimensional Temporal Disaggregation*. CRAN, CRAN Repository. R package version 2.0.0.
- Mulgan, G., Tucker, S., Ali, R., Sanders, B., et al. (2007). *Social innovation: what it is, why it matters and how it can be accelerated*. Young Foundation.
- OECD (2009). *OECD Science, Technology and Industry Scoreboard 2009: R&D in the economic crisis*. OECD Publishing.
- OECD (2015). *Measuring R&D: Methodologies and procedures*. OECD Publishing.
- OECD (2023). Gross domestic spending on R&D (indicator). Accessed on 27 October 2023.
- Preis, T., Moat, H. S., and Stanley, H. E. (2013). Quantifying trading behavior in financial markets using google trends. *Nature Scientific Reports*, 3(1):1–6.
- Radhakrishnan, A., Beaglehole, D., Pandit, P., and Belkin, M. (2022). Feature learning in neural networks and kernel machines that recursively learn features. *arXiv preprint arXiv:2212.13881*.
- Romer, P. M. (1990). Endogenous technological change. *Journal of Political Economy*, 98(5, Part 2):S71–S102.
- Santurkar, S., Tsipras, D., Ilyas, A., and Madry, A. (2018). How does batch normalization help optimization? *Advances in Neural Information Processing Systems*, 31.
- Sax, C. and Steiner, P. (2013). Temporal disaggregation of time series. *The R Journal*, 5(2):80–87.
- Shu, X. and Ye, Y. (2023). Knowledge discovery: Methods from data mining and machine learning. *Social Science Research*, 110:102817.
- Sokolov-Mladenović, S., Milovančević, M., Mladenović, I., and Alizamir, M. (2016). Economic growth forecasting by artificial neural network with extreme learning machine based on trade, import and export parameters. *Computers in Human Behavior*, 65:43–45.
- StatCounter (2024a). Search engine market share in china. <https://gs.statcounter.com/search-engine-market-share/all/china/#monthly-200901-202406>. Accessed on 8 July 2024.
- StatCounter (2024b). Search engine market share in japan. <https://gs.statcounter.com/search-engine-market-share/all/japan/#monthly-200901-202406>. Accessed on 8 July 2024.
- StatCounter (2024c). Search engine market share in south korea. <https://gs.statcounter.com/search-engine-market-share/all/south-korea/#monthly-200901-202406>. Accessed on 8 July 2024.
- Tajaddini, R. and Gholipour, H. F. (2021). Economic policy uncertainty, R&D expenditures and innovation outputs. *Journal of Economic Studies*, 48(2):413–427.
- Tümer, A. E. and Akkuş, A. (2018). Forecasting gross domestic product per capita using artificial neural networks with non-economical parameters. *Physica A: Statistical Mechanics and its Applications*, 512:468–473.
- WIPO (2023). *Global Innovation Index 2023: Innovation in the face of uncertainty*. World Intellectual Property Organization (WIPO).
- Woloszko, N. (2020). Tracking activity in real time with google trends. *OECD*.
- Woloszko, N. (2023). Nowcasting with panels and alternative data: The oecd weekly tracker. *International Journal of Forecasting*.
- Wu, L. and Brynjolfsson, E. (2015). The future of prediction: How Google searches foreshadow housing prices and sales. In *Economic Analysis of the Digital Economy*, pages 89–118. University of Chicago Press.
- Zhao, Y. and Yang, G. (2023). Deep learning-based integrated framework for stock price movement prediction. *Applied Soft Computing*, 133:109921.

A Google search terms

The initial ecosystem developed for this study encompassed a broader range of search terms, detailed in table A.1.

Table A.1: Initial search terms to build a network of stakeholders.

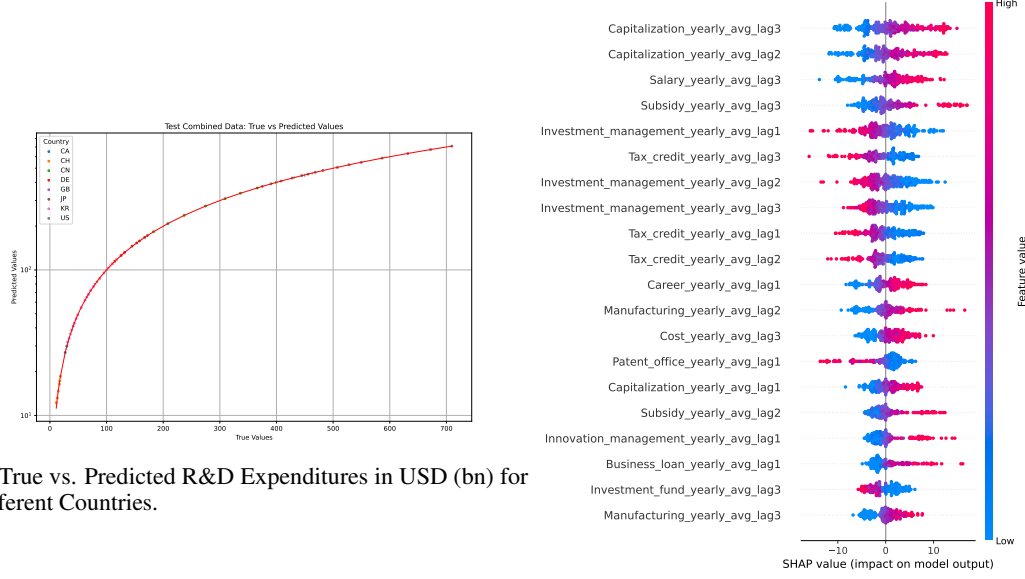
Stakeholder	Keywords/Topics/Categories
Firms/Companies	R&D Expenditure, Product Development, Technology Innovation, Patent Application, Tech Research, Pharma Research, New Drug Application, Research Grants, Intellectual Property.
Venture Capitalists (VCs)	Startup Funding, Technology Startups, Pharma Startups, VC Investment in R&D, Innovation Investment, Return on Investment, Exit Strategy, Seed Funding, Angel Investment.
Banks and Financial Institutions	Business Loans, R&D Loans, Investment Banking, Corporate Finance, Financial Risk, Credit Assessment, Interest Rate, Loan Application, Credit Score.
Universities and Research Institutions	Academic Research, Collaboration with Industry, Research Funding, University Patents, Postgraduate Studies, Doctoral Research, Research Publication, Research Grant.
Government agencies	R&D Policy, Research Funding, Government Grants, Innovation Policy, Public-Private Partnership, Tax Incentives for R&D, Technology Transfer, Patent Law, Economic Development.
R&D Employees	Research Methods, Data Analysis, Patent Filing, Lab Equipment, Scientific Journal, Professional Development, Research Ethics, Project Management, Collaboration Tools.
Tax Authorities	R&D Tax Credit, Tax Incentives, Tax Deduction, Tax Filing, Corporate Tax, Tax Law, Tax Consultancy.
Consulting Firms	Business Strategy, Market Analysis, Risk Assessment, Business Growth, Innovation Strategy, Portfolio Management, Project Planning, Financial Modeling.
Innovation Hubs/Incubators	Startup Incubation, Innovation Hub, Technology Park, Business Accelerator, Entrepreneurship, Mentorship, Networking, Business Pitch, Startup Ecosystem.
Patent Attorneys	Patent registration, Intellectual property rights, Patent law, Technology patents, Patent disputes.
Tax Consultants/Accountants	R&D tax credits, Corporate tax, Business expenses, Tax deductions, Tax advice for R&D, Accounting for R&D.

B Out-of-sample performance and global interpretability of different configurations

In this section, we show the out-of-sample performance of the prediction model, for all the configuration individually.

B.1 AGT Configuration

Figure B.1: Analysis for AGT Configuration.

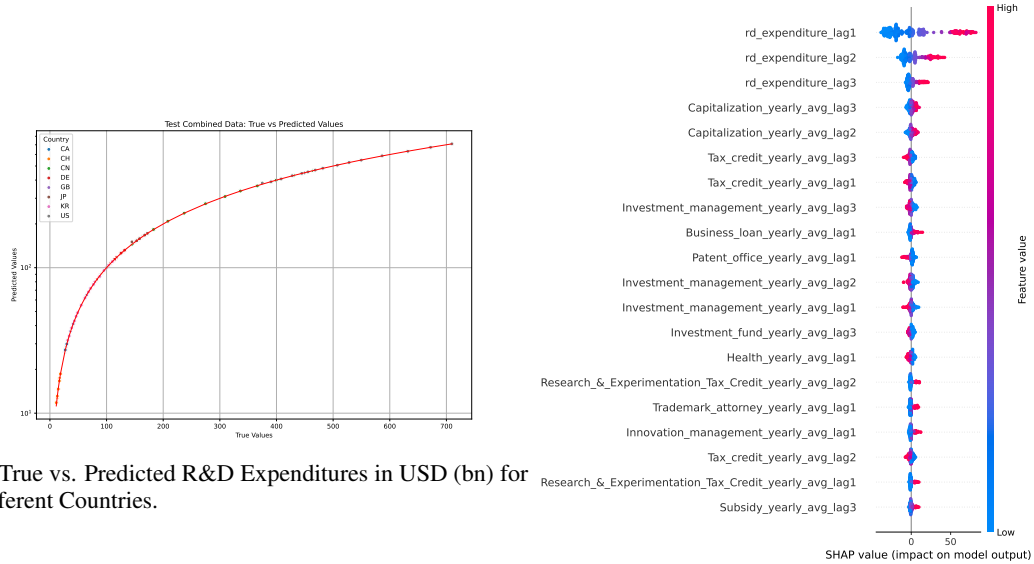


(a) True vs. Predicted R&D Expenditures in USD (bn) for Different Countries.

(b) SHAP summary plot.

B.2 AGTwRD Configuration

Figure B.2: Analysis for AGTwRD Configuration.

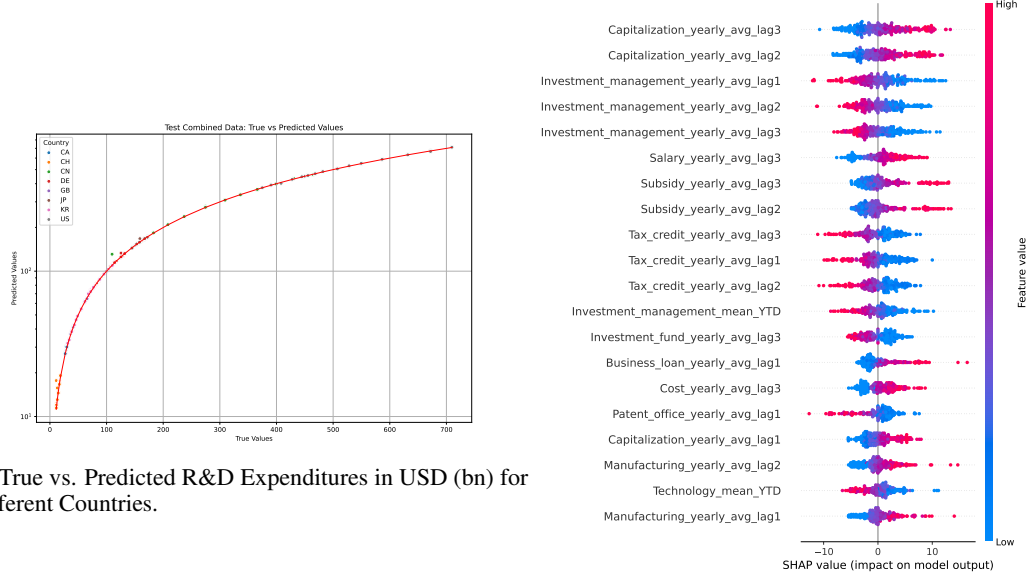


(a) True vs. Predicted R&D Expenditures in USD (bn) for Different Countries.

(b) SHAP summary plot.

B.3 MGT Configuration

Figure B.3: Analysis for *MGT* Configuration.

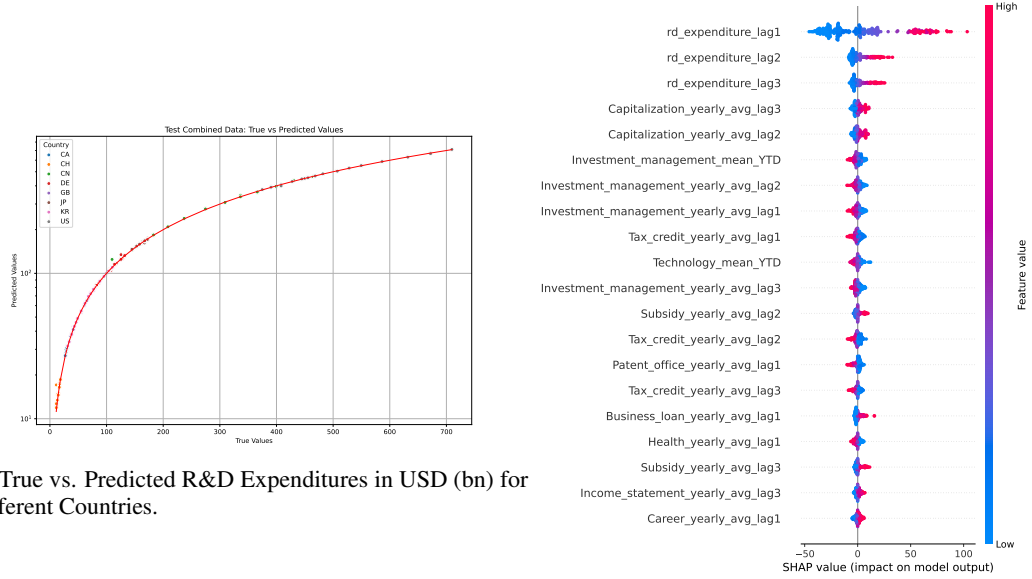


(a) True vs. Predicted R&D Expenditures in USD (bn) for Different Countries.

(b) SHAP summary plot.

B.4 MGTwRD Configuration

Figure B.4: Analysis for *MGTwRD* Configuration.

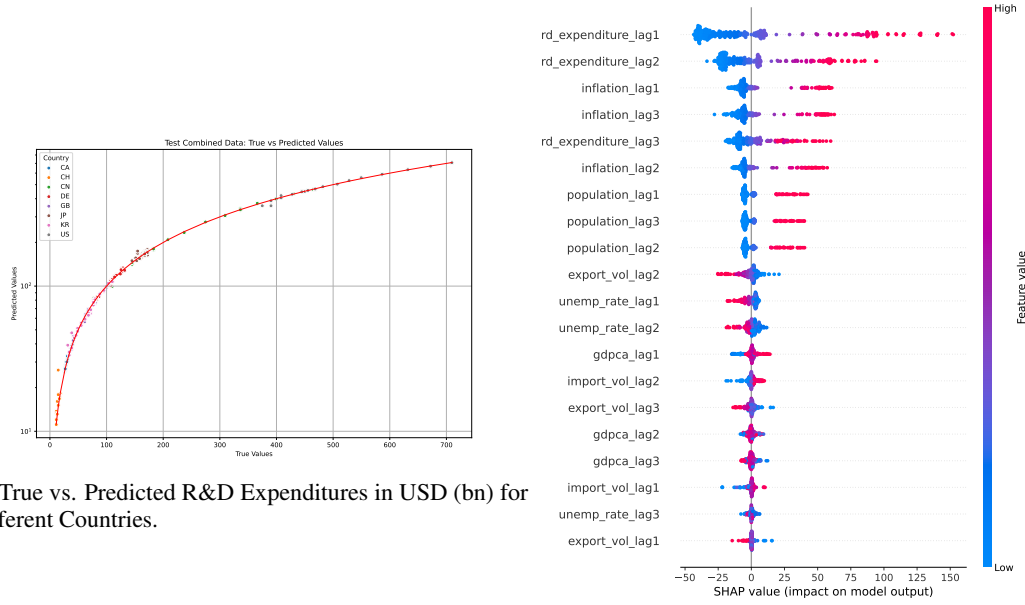


(a) True vs. Predicted R&D Expenditures in USD (bn) for Different Countries.

(b) SHAP summary plot.

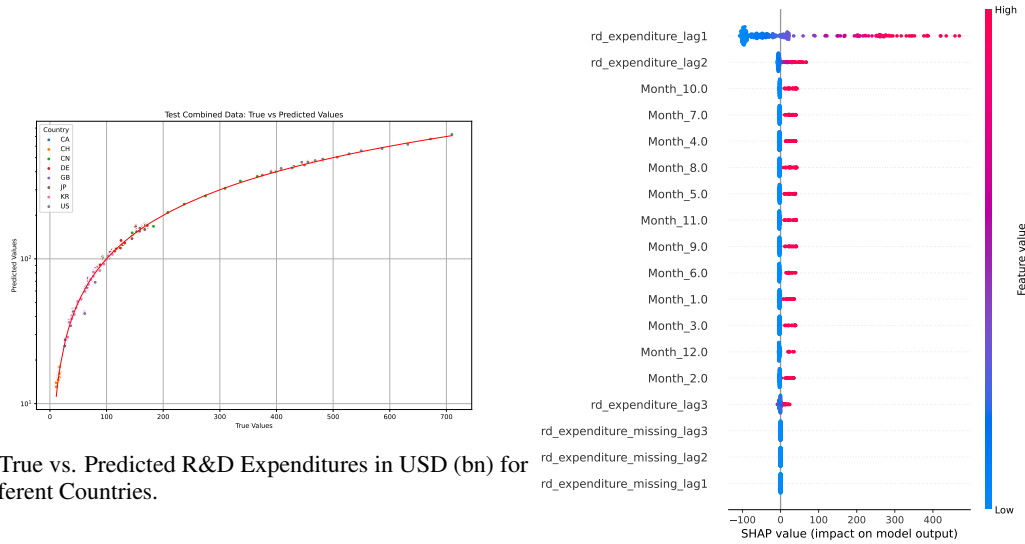
B.5 Macros Configuration

Figure B.5: Analysis for *Macros* Configuration.



B.6 LagRD Configuration

Figure B.6: Analysis for *LagRD* Configuration.



C Country-level out-of-sample performances

C.1 *AllVar* Configuration

In the following figures C.1, C.2, C.3 and C.4, we illustrate the out-of-sample performance of the prediction model, for each country, in *AllVar* configuration, as the most comprehensive configuration in terms of input vector.

Figure C.1: Yearly R&D expenditures for selected countries in North America and Asia: True Values vs. Estimates.

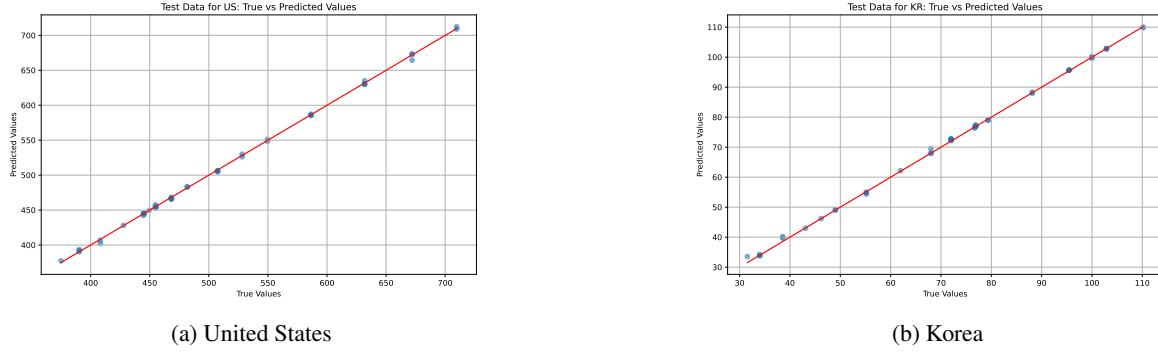


Figure C.2: Yearly R&D expenditures for selected countries in Europe: True Values vs. Estimates.

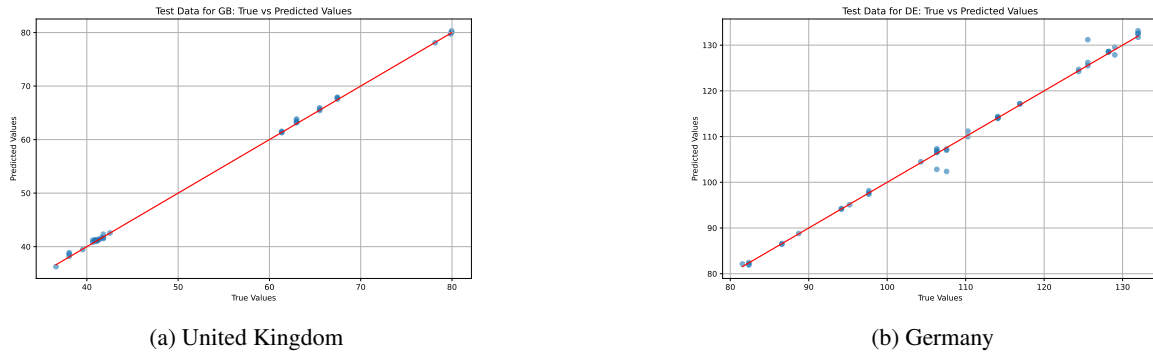


Figure C.3: Yearly R&D expenditures for selected countries in North America and East Asia: True Values vs. Estimates.

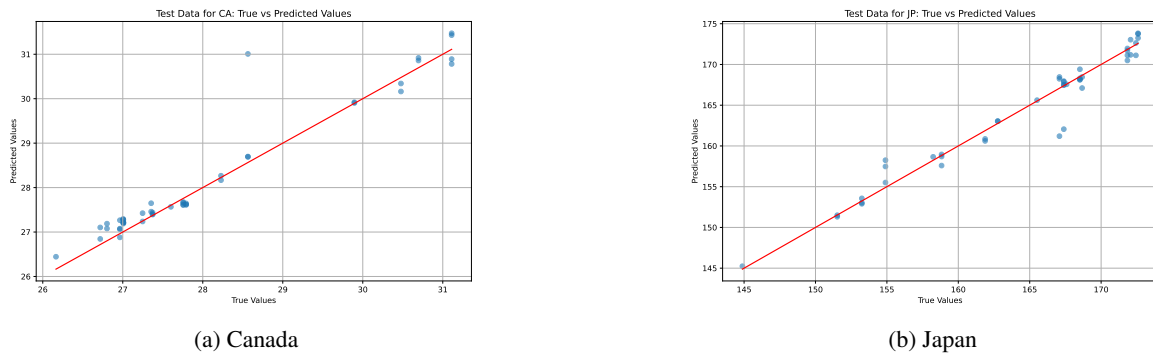
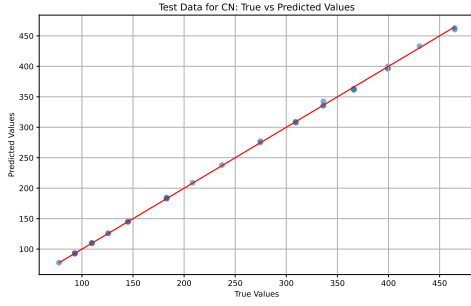
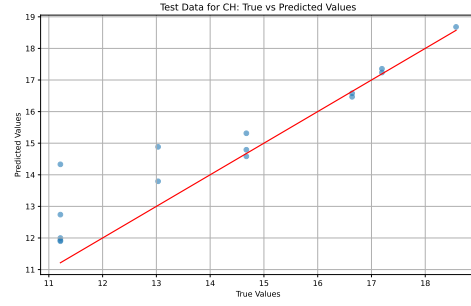


Figure C.4: Yearly R&D expenditures for selected countries in East Asia and Europe: True Values vs. Estimates.



(a) China

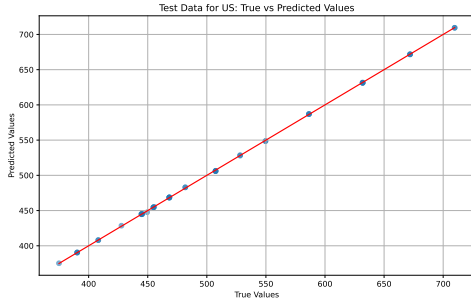


(b) Switzerland

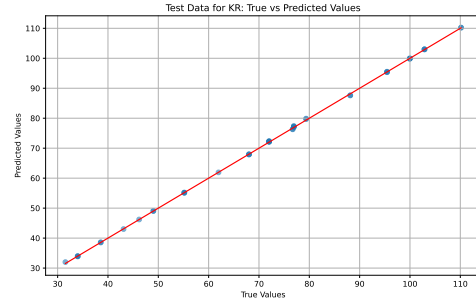
C.2 AGT Configuration

Similarly, in the next batch of figures C.5, C.6, C.7 and C.8, we present the out-of-sample performance of the prediction model, for each country, in *AGT* configuration, as the outperformer.

Figure C.5: Yearly R&D expenditures for selected countries in North America and Asia: True Values vs. Estimates.

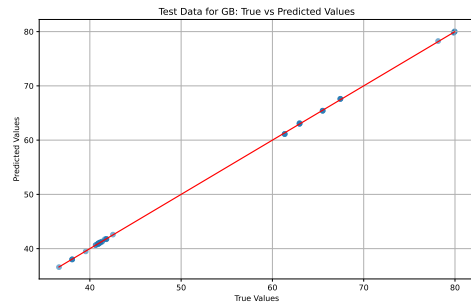


(a) United States

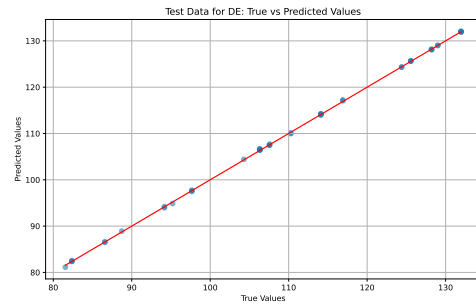


(b) Korea

Figure C.6: Yearly R&D expenditures for selected countries in Europe: True Values vs. Estimates.

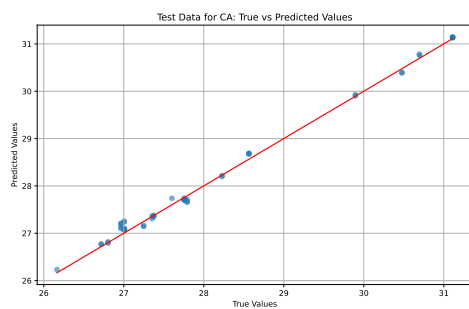


(a) United Kingdom

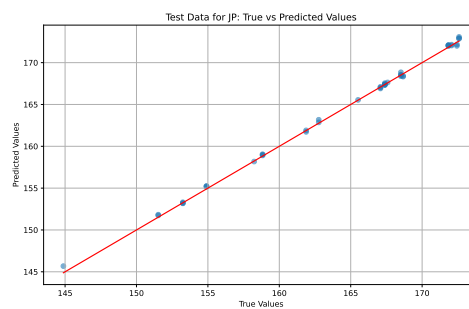


(b) Germany

Figure C.7: Yearly R&D expenditures for selected countries in North America and East Asia: True Values vs. Estimates.

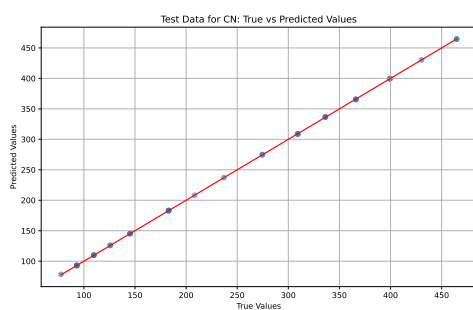


(a) Canada

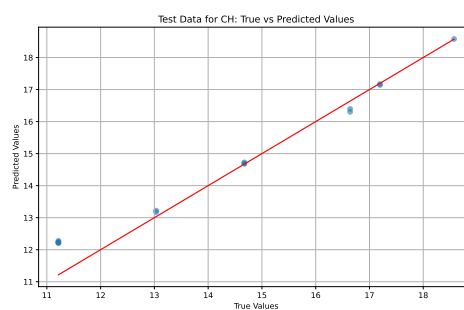


(b) Japan

Figure C.8: Yearly R&D expenditures for selected countries in East Asia and Europe: True Values vs. Estimates.



(a) China



(b) Switzerland

D Interpolating R&D expenditures at a higher frequency via Perturbed/Corrupted input

As an alternative approach for interpolating R&D expenditures, we elaborate on broadening our model by perturbing the input. We leverage on a trained model for yearly prediction, as a forward pass to have distribution of estimated values with respect to various input levels, and then develop a backward pass to breakdown yearly values and assign the resulted values to each month as its estimated contribution to the full-year expenditure. The model on the forward pass, is similar to the nowcasting model developed for yearly frequency in Section 3, in particular the *AGT* configuration:

$$y_{t,i} = f^{(m)} \left(\mathcal{S}_{t-\tau,i}^{(m)}, \mathbf{C}_{j,i}^{(m)}; \boldsymbol{\theta}^{(m)} \right) + \varepsilon_{t,j,i}^{(m)}$$

Since the *AGT* model stands out as the top performer, according to the results in Section 5, we build up on that with one modification on the aggregation approach for the annual GT values, as follows:

$$\begin{aligned} \mathcal{S}_{t-\tau,i} &= [\tilde{\mathbf{s}}_{t,i} \dots \tilde{\mathbf{s}}_{t-\tau,i}] \\ \tilde{\mathbf{s}}_{t-\tau,i} &= [\tilde{s}_{1,t-\tau,i} \dots \tilde{s}_{k_s,t-\tau,i}] \\ \tilde{s}_{k_s,t-\tau,i} &= \sum_{j=0}^{11} s_{k_s,t-\tau,j,i}^{(m)} \end{aligned}$$

After training the model in this setting, we run the model over different input levels, which can be considered as corrupted versions of the original inputs, to get a distribution of outputs that are diverged from the original predictions, accordingly. In other words, we create different sub-sets for $\tilde{\mathbf{s}}_{t,i}$ in format of $\tilde{\mathbf{s}}_{t-j/12,i}^{(m)}$ in which we sum only GT values up to the values of a particular month (the lower bound varies accordingly). j indicates a particular month in which the nowcasting is being done, and it takes values of $j = 0, \dots, 11$ accordingly; and (m) superscript reflects the monthly dynamics and sampling frequency, in accordance to the yearly frequency model (Section 3). Following the forward pass across all months and for each country, we proceed to the backward pass phase. This phase begins with the last month of the year (December, where $j = 0$) and involves comparing its estimated target value with that of the preceding month. This comparison is unique in that the preceding month's input vector contains one less month of data. This backward pass continues sequentially for each month, effectively covering the entire year and also all available years. The difference observed for each month, referred to as the discrepancy, is interpreted as the estimated contribution of that month to the overall model output.

As a preliminary result to the higher-frequency approach, figure D.1 depicts estimated monthly R&D expenditure over time for the US.

Figure D.1: Cumulative Estimated Monthly R&D Expenditures over time for the US.

