

# The use of big data analytics and artificial intelligence in central banking

Okiriza Wibisono, Hidayah Dhini Ari, Anggraini Widjanarti, Alvin Andhika Zulen and Bruno Tissot<sup>1</sup>

## Executive summary

Information and internet technology has fostered new web-based services that affect every facet of today's economic and financial activity. This creates enormous quantities of "**big data**" – defined as "*the massive volume of data that is generated by the increasing use of digital tools and information systems*" (FSB (2017)). Such data are produced in real time, in differing formats, and by a wide range of institutions and individuals. For their part, central banks face a surge in "financial big data sets", reflecting the combination of new, rapidly developing electronic footprints as well as large and growing financial, administrative and commercial records.

This phenomenon has the potential to strengthen analysis for decision-making, by providing more complete, immediate and granular information as a complement to "traditional" macroeconomic indicators. To this end, a number of techniques are being developed, often referred to as "**big data analytics**" and "**artificial intelligence**" (AI). These promise faster, more holistic and more connected insights, as compared with traditional statistical techniques and analyses. An increasing number of central banks have launched specific big data initiatives to explore these issues. They are also sharing their expertise in collecting, working with, and using big data, especially in the context of the BIS's Irving Fisher Committee on Central Bank Statistics (IFC); see IFC (2017).

Getting the most out of these new developments is no trivial task for policymakers. Central banks, like other public authorities, face numerous challenges, especially in handling these new data and using them for policy purposes. In particular, significant resources are often required to handle large and complex data sets, while the benefits of such investments are not always clear-cut. For instance, to what extent should sophisticated techniques be used to deal with this type of information? What is the added value over more traditional approaches, and how should the results be interpreted? How can the associated insights be integrated into current decision-making processes and be communicated to the public? And, lastly, what are the best strategies for central banks seeking to realise the full potential of

<sup>1</sup> Respectively Big Data Analyst ([okiriza\\_w@bi.go.id](mailto:okiriza_w@bi.go.id)); Head of Digital Data Statistics and Big Data Analytics Development Division ([dhini\\_ari@bi.go.id](mailto:dhini_ari@bi.go.id)); Big Data Analyst ([anggraini\\_widjanarti@bi.go.id](mailto:anggraini_widjanarti@bi.go.id)); Big Data Analyst ([alvin\\_az@bi.go.id](mailto:alvin_az@bi.go.id)); and Head of Statistics and Research Support, BIS, and Head of the IFC Secretariat ([Bruno.Tissot@bis.org](mailto:Bruno.Tissot@bis.org)).

The views expressed here are those of the authors and do not necessarily reflect those of Bank Indonesia, the Bank for International Settlements (BIS), or the Irving Fisher Committee on Central Bank Statistics (IFC).

new big data information and analytical tools, considering in particular resource constraints and other priorities?

Against this backdrop, Bank Indonesia organised with support from the BIS and the IFC a workshop on *"Big data for central bank policies"* and a high-level policy-oriented seminar on *"Building pathways for policymaking with big data"*. Convening in July 2018 at Bali, Indonesia, the events were attended by officials from central banks, international organisations and national statistical offices from more than 30 jurisdictions across the globe, as well as by representatives from other public agencies, the financial sector and academia. This proved a useful opportunity to take stock of the various big data pilots conducted by the central bank community and of the growing use of big data analytics and associated AI techniques to support public policy. The following points of interest were highlighted:

- Big data offers **new types of data source** that complement more traditional varieties of statistics. These sources include Google searches, real estate and consumer prices displayed on the internet, and indicators of economic agents' sentiment and expectations (eg social media).
- Thanks to IT innovation, **new techniques** can be used to collect data (eg web-scraping), process textual information (text-mining), match different data sources (eg fuzzy-matching), extract relevant information (eg machine learning) and communicate or display pertinent indicators (eg interactive dashboards).
- In particular, big data techniques such as decision trees may shed **interesting light** on the decision-making process of economic agents, eg how investors behave in financial markets. As another example, indicators of economic uncertainty extracted from news articles, could help explain movements of macroeconomic indicators. This illustrates big data's potential in providing insights not only into what happened, but also into what might happen and why.
- In turn, these new insights can usefully **support central bank policies** in a wide range of areas, such as market information (eg credit risk analysis), economic forecasting (eg nowcasting), financial stability assessments (eg network analysis) and external communication (eg measurement of agents' perceptions). Interestingly, the approach can be very granular, helping to target specific markets, institutions, instruments and locations (eg zip codes) and, in particular, to support macroprudential policies. Moreover, big data indicators are often more timely than "traditional" statistics – for instance labour indicators can be extracted from online job advertisements almost in real time. In addition, big data indicators are often more timely than "traditional" statistics – for instance, labour indicators can be extracted from online job advertisements almost in real time.
- As a note of caution, feedback from central bank pilot projects consistently highlights the complex privacy implications of dealing with big data, and the associated **reputational risks**. Moreover, while big data applications such as machine learning algorithms can excel in terms of predictive performance, they can lend themselves more to explaining what is happening rather than why. As such, they may be exposed to public criticism when insights gained in this way are used to justify policy decisions.

- Another **concern** is that, as big data samples are often far from representative (eg not everyone is on Facebook, and even fewer are on Twitter), they may not be as reliable as they seem. Lastly, there is a risk that collecting and processing big data will be hindered by privacy laws and/or change in market participants. Relevant authorities should coordinate their efforts so that they can utilise the advantages of big data analytics without compromising data privacy and confidentiality.

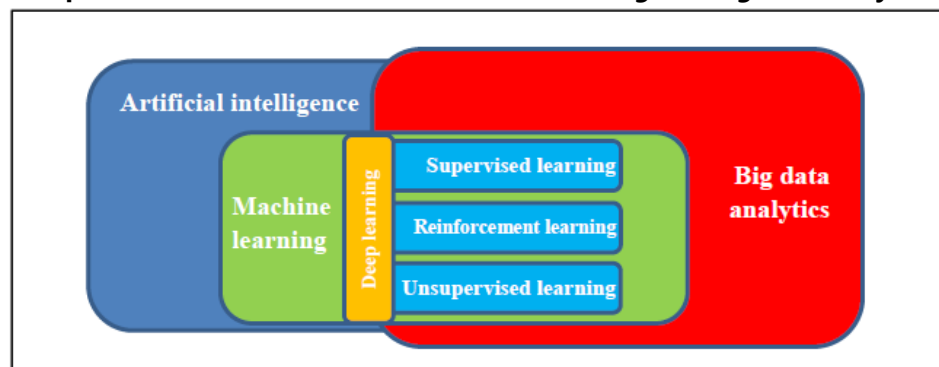
The related presentations, referred to in this overview and included in this *IFC Bulletin*, analysed the various aspects related to the use of big data and associated techniques by central banks. They cover three main aspects: (1) an assessment of the main big data sources and associated analytical techniques that are relevant for central banks; (2) the insights provided by big data for economic policy, with an overview of concrete central bank projects aiming at improving statistical information, macroeconomic analysis and forecasting, financial market monitoring and financial risk assessment; and (3) the use of big data in crafting central bank policies, including organisational aspects and related challenges.

## 1. The big data revolution: new data sources and analytical techniques

As emphasised in the opening remarks by Erwin Rijanto, Deputy Governor of Bank Indonesia, policymakers should not miss out on the opportunities provided by big data – described by some as the new oil of the 21st century (The Economist (2017)). Public institutions are not the main producers of big data sets, and some of this information may have little relevance for their daily work. Yet central banks are increasingly dealing with “financial big data” sources that impinge on a wide range of their activities, as noted by Claudia Buch, IFC Chair and Vice President of the Deutsche Bundesbank.

Data volumes have surged hand in hand with the development of specific techniques for their analysis, thanks to “*big data analytics*” – broadly referring to the general analysis of these data sets – and “*artificial intelligence*” (AI) – defined as “*the theory and development of computer systems able to perform tasks that traditionally have required human intelligence*” (FSB (2017)). Strictly speaking, these two concepts can differ somewhat (for instance, one can develop tools to analyse big data sets that are not based on AI techniques), as shown in Graph 1.

**Graph 1: A schematic view of AI, machine learning and big data analytics**



Source: FSB (2017).

In practice, and as underlined by Yati Kurniati (Bank Indonesia) in her welcoming remarks, big data analytics are not very different from traditional econometrical techniques, and indeed they borrow from many long-established methodologies and techniques developed for general statistics; for instance, principal component analysis, developed at the beginning of the last century. Yet one key characteristic is that they are applied to modern data sets that can be both very large and complex. Extracting relevant information from these sources is not straightforward, often requiring a distinct set of skills, depending on the type of information involved. As a result, big data analytics and AI techniques comprise a variety of statistical/modelling approaches, such as machine learning, text-mining techniques, network analysis, agent-based modelling<sup>2</sup> etc.

The seminar and workshop provided an opportunity to review, first, the main big data sources relevant for central banks, and, second, the principal techniques developed in recent years for analysing big data – focusing on the classification and clustering of information derived from large quantitative data sets, with machine learning, text-mining and network analysis all playing an important role.

## Big data information for central banks

**Three main sources of big data can be identified**, as Paul Robinson (Bank of England) noted in his introductory presentation on fundamental concepts and frameworks.<sup>3</sup> These categories are related to (i) social networks (human-sourced information such as blogs and searches); (ii) traditional business systems (process-mediated data, such as files produced by commercial transactions, e-commerce, credit card operations); and (iii) the internet of things (machine-generated data, such as information produced by pollution/traffic sensors, mobile phones, computer logs etc). These are very generic categories and, in practice, big data will comprise multiple types of heterogeneous data set derived from these three main sources.

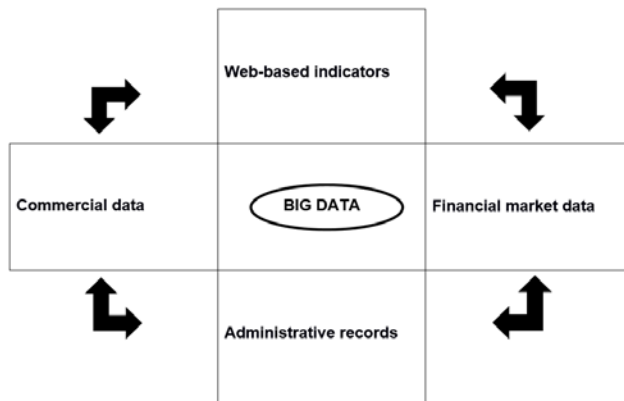
Focusing more specifically on central banks, the presentation by Bruno Tissot (BIS) identified **four types of data set** that would usually be described as financial big data (see Graph 2): internet-based indicators, commercial data sets, financial market indicators and administrative records.<sup>4</sup>

<sup>2</sup> See Haldane (2018), who argues that big data can facilitate policymakers' understanding of economic agents' reactions through the exploration of behaviours in a "virtual economy".

<sup>3</sup> Following the work conducted under the aegis of the United Nations (see Meeting of the Expert Group on International Statistical Classifications (2015)).

<sup>4</sup> For the use of administrative data sources for official statistics, see for instance Bean (2016) in the UK context.

**Graph 2: Four main types of financial big data set**



Compared with the private sector,<sup>5</sup> central banks' use of **web-based indicators** may be somewhat more limited, especially as regards unstructured data such as images. Even so, several projects are under way to make use of data collected on the internet to support monetary and financial policymaking (see Section 2). Moreover, an important aspect relates to the increased access to **digitalised information**, reflecting both the fact that more and more textual information is becoming available on the web (eg social media) and also that "traditional" printed documents can now be easily digitalised, searched and analysed in much the same way as web-based indicators.

In reality, however, the bulk of the financial big data sets relevant to central banks consists of the very granular information provided by large and growing records covering commercial transactions, financial market developments and administrative operations. This type of information has been spurred by the expansion of the **micro-level data sets** collected in the aftermath of the Great Financial Crisis (GFC) of 2007–09, especially in the context of the Data Gaps Initiative (DGI) endorsed by the G20 (FSB-IMF (2009)). For instance, significant efforts have been made globally to compile large and granular loan-by-loan and security-by-security databases as well as records of individual derivatives trades (IFC (2018)). As a result, central banks now have at their disposal very detailed information on the financial system, including at the level of specific institutions, transactions or instruments.

### Extracting knowledge from large quantitative data sets: classification and clustering

The expansion of big data sources has gone hand in hand with the development of new analytical tools to deal with them. The first, and particularly important, category of these big data techniques aims at extracting summary information from large quantitative data sets. This is an area that is relatively close to "traditional statistics", as it does not involve the treatment of unstructured information (eg text, images). In fact, many big data sets are well structured, and can be appropriately dealt with using statistical algorithms developed for numerical data sets. The main goal is to obtain summary indicators by condensing the large amount of data points available,

<sup>5</sup> Especially the major US technology companies ("GAFAs"): Google, Apple, Facebook and Amazon.

basically by finding similarities between them (through classification) and regrouping them (through clustering).

Many of these techniques involve so-called **machine learning**. This is a subset of AI techniques, that can be defined as *"a method of designing a sequence of actions to solve a problem that optimise automatically through experience and with limited or no human intervention"* (FSB (2017)). This approach is quite close to conventional econometrics, albeit with three distinct features. First, machine learning is typically focused on prediction rather than identifying a causal relationship. Second, the aim is to choose an algorithm that fits with the actual data observed, rather than with a theoretical model. Third, and linked to the previous point, the techniques are selected by looking at their goodness-to-fit, and less at the more traditional statistical tests used in econometrics.

The lecture by Sanjiv Das (Santa Clara University) recalled that there are several categories of machine learning, which can be split into two main groups. First, in **supervised machine learning**, *"an algorithm is fed a set of 'training' data that contain labels on the observations"* (FSB (2017)). The goal is to classify individual data points, by identifying, among several classes (ie categories of observations), the one to which a new observation belongs. This is inferred from the analysis of a sample of past observations, ie the training data set, for which their group (category) is known. The objective of the algorithm is to predict the category of a new observation, depending on its characteristics. For instance, to predict the approval of a new loan ("yes" or "no", depending on its features and in comparison with an observed historical data set of loans that have been approved or rejected); or whether a firm is likely to default in a few months. Various algorithms can be implemented for this purpose, including logistic regression techniques, linear discriminant analysis, Naïve Bayes classifier, support vector machines, k-nearest neighbours, decision trees, random forest etc.

The second group is **unsupervised machine learning**, for which *"the data provided to the algorithm do not contain labels"*. This means that categories have not been identified ex ante for a specific set of observations, so that the algorithm has to identify the clusters, regrouping observations for which it detects similar characteristics or "patterns". Two prominent examples are clustering and dimensionality reduction algorithms. In **clustering**, the aim is to detect the underlying groups that exist in the granular data set – for example, to identify groups of customers or firms that have similar characteristics – by putting the most similar observations in the same cluster in an agglomerative way (bottom-up approach).<sup>6</sup> **Dimensionality reduction** relates to the rearrangement of the original information in a smaller number of pockets, in a divisive (top-down) way; the objective is that the number of independent variables becomes (significantly) smaller, without too much compromise in terms of information loss.

There are, of course, additional types of algorithm. One is **reinforcement learning**, which complements unsupervised learning with additional information feedback, for instance through human intervention. Another is **deep learning** (or artificial neural networks), based on data representations inspired by the function of neurons in the brain. Recent evaluations suggest that deep learning can perform better than traditional classification algorithms when dealing with unstructured data

<sup>6</sup> More precisely, cluster analysis can be defined as *"a statistical technique whereby data or objects are classified into groups (clusters) that are similar to one another but different from data or objects in other clusters"* (FSB (2017)).

such as texts and images – one reason being that applying traditional quantitative algorithms is problematic, as it requires unstructured information to be converted into a numerical format. In contrast, deep learning techniques can be used to deal directly with the original raw data.

In view of this diversity, the choice of a specific algorithm will depend on the assumptions made regarding the features of the data set of interest – for instance, a Naïve Bayes classifier would be appropriate when the variables are assumed to be independent and follow a Gaussian distribution. In practice, data scientists will have to identify which algorithm works best for the problem at hand, often requiring a rigorous and repetitive process of trial and error; this is often as much art as it is science.

In choosing the right “model”, it is important to define an evaluation metric. The aim is to measure how well a specific algorithm fits, and to compare the performance of alternative algorithms. The most straightforward metric for classification is **prediction accuracy**, which is simply the percentage of observations for which the algorithm predicts the class variable correctly (this will usually be done by comparing the result of the algorithm with what a human evaluator would conclude on a specific data sample). But an accuracy metric may not be suitable for all exercises, particularly in the case of an unbalanced distribution of classes. For example, when looking at whether a transaction is legitimate or fraudulent, a very simplistic model could be adopted that assumes that all transactions are legitimate: its accuracy will look very high, because *a priori* most transactions are not fraudulent; but the usefulness of such a simple model would be quite limited. Hence, other metrics have to be found for evaluating algorithms when the distribution of classes is highly unbalanced.<sup>7</sup> Another possible approach is to address the class imbalance issue at the observation level, for instance, by duplicating (over-sampling) elements from the minority class or, conversely, by discarding those (under-sampling) from the dominant class.

## Text mining

Another rapidly developing area of big data analytics is text-mining, ie analysis of semantic information – through the automated analysis of large quantities of natural language text and the detection of lexical or linguistic patterns with the aim of extracting useful insights. While most empirical work in economics deals with numerical indicators, such as prices or sales data, a large and increasing amount of textual information is also generated by economic and financial activities – including internet-based activities (eg social media posts), but also the wider range of textual information provided by, say, company financial reports, media articles, public authorities’ deliberations etc. Analysing this unstructured information has become of key interest to policymakers, not least in view of the important role played by “soft” indicators such as confidence and expectations during the GFC. As illustrated in the lecture delivered by Stephen Hansen (University of Oxford), text-mining techniques can usefully be applied to dealing with these data in a structured, quantitative way.

<sup>7</sup> Such other metrics include, for instance, precision, recall and F1-score. For binary (two-class) classification, precision is defined as the percentage of times an algorithm makes a correct prediction for the positive class; recall is defined as the percent of positive class that the algorithm discovers from a given data set; and the F1-score is the harmonic average of precision and recall.

Text analysis typically starts with some standard **pre-processing steps**, such as tokenisation (splitting text into words), stopword removal (discarding very frequent/non-topical words eg “a”, “the”, “to”), stemming or lemmatising (converting words into their root forms, for instance “prediction” and “predicted” into “predict”), and merging words within a common message (eg “Bank” and “Indonesia” grouped into “Bank Indonesia”). Once this is done, the initial document can be transformed into a document-term matrix, which indicates for each specific text a term’s degree of appearance (or non-appearance). This vectoral text representation is made of numerical values that can then be analysed by quantitative algorithms; for example, to measure the degree of similarity between documents by comparing the related matrixes (Graph 3).

One popular algorithm for working on textual information is the **Latent Dirichlet Allocation** (LDA).<sup>8</sup> This assumes that documents are distributed by topics, which in turn are distributed by keywords. For example, one document may combine, for a respective 20% and 80%, a “monetary” and an “employment” topic, based on the number of words reflecting this topic distribution (ie 20% of them related to words such as “inflation” or “interest rate”, and the remaining 80% related to words such as “jobs” and “labour”). Based on these calculations, one can build an indicator measuring how frequently a specific topic appears over time, for instance, to gauge the frequency of the messages related to “recession” – providing useful insights when monitoring the state of the economy.

Besides quantitative algorithms, simpler **dictionary-based methods** can be also employed for analysing text data. A set of keywords can be selected that are relevant to the topic of interest – for example, a keyword related to “business confidence”. Then an index can be constructed based on how frequently these selected keywords appear in a given document, allowing the subject indicator to be assessed (eg the evolution of business sentiment).<sup>9</sup> A prominent example is the Economic Policy Uncertainty (EPU), which quantifies the degree of uncertainty based on the appearance of a set of economic-, policy-, and uncertainty-related keywords in news articles; by the end of 2018, more than 20 country-specific EPU indexes had been compiled.<sup>10</sup>

<sup>8</sup> See Blei et al (2003).

<sup>9</sup> See for example Tetlock (2007) and Loughran and McDonald (2011).

<sup>10</sup> See Baker et al (2016) and [www.policyuncertainty.com/](http://www.policyuncertainty.com/).





complex network by regrouping nodes in clusters and filtering noise, through the use of specific machine learning algorithms (see above).

This sort of analysis appears particularly well suited to representing **interconnectedness** within a system, for instance, by mapping the global value chain across countries and sectors or the types of exposure incurred by financial institutions.<sup>12</sup> One example is recent work to assess the role of CCPs in the financial system by looking at the connections between them as well as with other financial institutions such as banking groups, in particular, by considering subsidiary-parent relationships (CPMI-IOSCO (2018)). This can help to reveal how a disruption originating in one single CCP would affect that CCP's clearing members and, in turn, other CCPs.

## 2. Opportunities for central banks

Big data can play an important role in improving the quality of economic analysis and research, as increasingly recognised by policymakers. This was the starting point for the presentation by Gabriel Quirós-Romero (IMF). The IMF is researching big data as a new way of measuring economic indicators, such as prices, labour market conditions, the housing market, business sentiment etc (Hammer et al (2017)).

Many central banks are now working on how to make use of the characteristics of financial big data sets in pursuing their mandates (Cœuré (2017)). As recalled in the introduction by Paul Robinson (Bank of England), big data has many advantages in terms of details, flexibility, timeliness and efficiency, as summarised in the list of their so-called "Vs" – eg volumes, variety, velocity, veracity and value; see Laney (2001) and Nymand-Andersen (2015). Central banks are interested in developing various pilot projects to better understand the new data sets and techniques, assess their value added in comparison with traditional approaches, and develop concrete "use cases" (IFC (2015)).

This experience has highlighted the opportunities that big data analytics can provide in key areas of interest to central banks, namely (i) the production of statistical information; (ii) macroeconomic analysis and forecasting; (iii) financial market monitoring; and (iv) financial risk assessment.

### More and enhanced statistical information

Big data can be a useful means of improving the official statistical apparatus. First, it can be an **innovative source of support for the current production of official statistics**, offering access to a wider set of data, in particular to those that are available in an "*organic*" way. Unlike statistical surveys and censuses, these data are usually not collected ("*designed*") for a specific statistical purpose, being the by-product of other activities (Groves (2011)). Their range is quite large, covering transaction data (eg prices recorded online), aspirational data (eg social media posts, product reviews displayed on the internet), but also various commercial, financial and administrative indicators. In addition, they present various advantages for statistical

<sup>12</sup> For a recent example of the monitoring of network effects for global systemic institutions in the context of the DGI, see FSB (2011) and Bese Goksu and Tissot (2018).

compilers, such as their rapid availability and the relative ease of collection and processing with modern computing techniques (see Graph 4) – always noting, however, that actual access to such sources, private or public, can be restricted by commercial and/or confidentiality considerations.

**Graph 4: Relative advantages of designed versus organic data**

	Designed data	Organic data
<b>Structure</b>	Geographic and socio-economic	Behaviour
<b>Representative</b>	Yes	No
<b>Sample selection</b>	Response rates deteriorating	Extreme
<b>Intrusive</b>	Extremely intrusive	Non-intrusive
<b>Cost</b>	Large	Small
<b>Curation</b>	Well studied	Unclear
<b>Privacy</b>	Well protected	Large violations of privacy

Source: Rigobon (2018).

Organic data can be used to **enhance existing statistical exercises**, especially in improving coverage when this is incomplete. In some advanced economies, the direct web-scraping<sup>13</sup> of online retailers’ prices data can, for instance, be used to better measure some specific components of inflation, such as fresh food prices.<sup>14</sup> At the extreme, these data can replace traditional indicators in countries where the official statistical system is underdeveloped. As noted by Roberto Rigobon (MIT Sloan School of Management), one example is the Billion Prices Project,<sup>15</sup> which allows inflation indices to be constructed for countries that lack an official and/or comprehensive index. Similarly, a number of central banks in emerging market economies have compiled quick price estimates for selected goods and properties, by directly scraping the information displayed on the web, instead of setting up specific surveys that can be quite time- and resource-intensive.

Second, big data can support a **timelier publication of official data**, by bridging the time lags before these statistics become available. In particular, the information generated instantaneously by the wide range of web and electronic devices – eg search queries – provides high-frequency indicators that can help current economic developments to be tracked more promptly (ie through the compilation of advance estimates). Indeed, another objective of the Billion Prices Project is to provide advance information on inflation in a large number of countries, including advanced economies, and with greater frequency – eg daily instead of monthly, as with a consumer price index (CPI). Turning to the real economy side, the real-time evolution of some “hard” indicators, such as GDP, can now be estimated in advance (“nowcast”) by using web-based information combined with machine learning algorithms, as presented by Tugrul Vehbi (Reserve Bank of New Zealand) in the case of New Zealand.

<sup>13</sup> Web-scraping can be defined as the automated capturing of online information.

<sup>14</sup> Hill (2018) reports that 15% of the US CPI is now collected online.

<sup>15</sup> See [www.thebillionpricesproject.com](http://www.thebillionpricesproject.com), and Cavallo and Rigobon (2016).

The high velocity of big data sources helps to provide more timely information, which can be particularly important during a crisis.

A third benefit is to provide **new types of statistics that complement “traditional” statistical data sets**, as emphasised in the presentation by Naruki Mori (Bank of Japan). Two important developments should be noted here. One is the increased availability of digitalised textual information, which allows the measurement of “soft” indicators such as economic agents’ sentiment and expectations – derived, for instance, from social media posts. Traditional statistical surveys can also provide this kind of information, but they typically focus on specific items eg firms’ production expectations and consumer confidence. In contrast, internet-based sources can cover a much wider range of topics. In addition, they are less intrusive than face-to-face surveys, and may therefore better reflect true behaviours and thoughts. A second important element has been the increased use of large granular data sets to improve the compilation of macroeconomic aggregates, allowing for a better understanding of their dispersion (IFC (2016a)) – this type of distribution information is generally missing in the System of National Accounts (SNA; European Commission et al (2009)) framework.<sup>16</sup>

## Macroeconomic forecasting with big data

Many central banks are already using big data sets for macroeconomic forecasting. Indeed, nowcasting applications as described above can be seen as a specific type of forecasting exercise. For instance, the presentation by Per Nymand-Andersen (ECB) showed how Google Trends data can be used to compile short-term projections of estimates of car sales in the euro area, with a lead time of several weeks over actual publication dates. Moreover, and as argued in the presentation by Alberto Urtaşun (Bank of Spain), big data allows a wider range of indicators to be used for forecasting headline indicators – for instance Google Trends,<sup>17</sup> uncertainty measures such as the EPU index (see above), or credit card operations as well as more traditional indicators. The devil is in the details, though, and statisticians need to try several approaches. For instance, some indicators may work well in **nowcasting** GDP (ie its growth rate over the current quarter) but less so in forecasting its future evolution (say, GDP growth one year ahead). Another point is that the internet is not the sole source of indicators that can be used in this context; in fact, some web-based indicators may work less well in nowcasting/forecasting exercises than do traditional business confidence surveys.<sup>18</sup>

In view of these caveats, and considering the vast amount of data potentially available, it may be useful to follow a **structured process** when conducting such exercises. The presentation of Paphatsorn Sawaengsuksant (Bank of Thailand)

<sup>16</sup> Indeed the SNA highlights the importance of considering the skewed distribution of income and wealth across households but recognises that getting this information is “*not straightforward and not a standard part of the SNA*” (2008 SNA, #24.69). It also emphasises that “*there would be considerable analytical advantages in having microdatabases that are fully compatible with the corresponding macroeconomic accounts*” (2008 SNA, #1.59). An important recommendation of the second phase of the DGI aims at addressing these issues (FSB-IMF (2015)).

<sup>17</sup> See <https://trends.google.com/trends/>. Google Trends provide indexes of the number of Google searches of given keywords. The indexes can be further segregated based on countries and provinces.

<sup>18</sup> For the use of nowcasting in forecasting “bridge models” using traditional statistics and confidence surveys, see Carnot et al (2011).

recommended a systematic approach when selecting the indicators of interest such as internet search queries. For instance, key words in Google Trends data could be selected if they satisfied several criteria, depending on their degree of generality, their popularity (ie number of searches recorded), their robustness (ie sensitivity to small semantic changes), their predictive value (ie correlation with macro indicators), and whether the relationship being tested makes sense from an economic perspective.

## Financial market forecasting and monitoring

As in the macroeconomic arena, big data analytics have also proved useful in monitoring and forecasting financial market developments, a key area for central banks. A number of projects in this area facilitate the processing of huge volume of **quantitative information** in large financial data sets. For instance, the presentation by Tom Fong (Hong Kong Monetary Authority) showed that returns in a number of emerging sovereign bond markets can be predicted using various technical trading rules and machine learning techniques, to assess their robustness as well as the relative contributions of specific foreign (eg US monetary policy) and domestic factors.

Other types of project are looking at **less structured data**. For instance, the presentation by Okiriza Wibisono (Bank Indonesia) described how a text-mining algorithm could be used to measure public expectations for the direction of interest rates in Indonesia.<sup>19</sup> Specifically, a classification algorithm is trained to predict whether a given piece of text indicates an expectation for the future tightening, loosening or stability of the central bank policy rate. All the newspaper articles discussing potential developments in the policy rate from two weeks prior to monthly policy meetings are collected, and an index of policy rate expectation is produced. This index has facilitated the analysis of the formation of policy rate expectations, usefully complementing other sources (eg Bloomberg surveys of market participants). Similarly, the presentation by Stephen Hansen (University of Oxford) looked at the information content of news articles grouped into several categories using an “emotion dictionary sample”, to predict equity returns. Other types of textual information, such as social media posts and official public statements, could also be usefully considered.

Experience reported by several central banks shows that new big data sources can also help to **elucidate** developments in financial markets, and shed light on their potential future direction. As regards the Bank of Japan, the use of high-frequency “tick data” has facilitated the assessment of market liquidity in the government bond market, and hence the risk of potential abrupt price changes. Similarly, the Bank of England has set up specific projects to identify forex market dynamics and liquidity at times of large market movements – eg when the Swiss National Bank decided to remove the EUR/CHF floor in January 2015 (Cielinska et al (2017)).

## Financial risk assessment

Big data sources and techniques can also facilitate financial risk assessment and surveillance exercises that sit at the core of central banks’ mandates – for both those in charge of micro financial supervision and those focusing mainly on financial

<sup>19</sup> See Zulen and Wibisono (2018).

stability issues and macro financial supervision (Tissot (2019)). In particular, the development of big data analytics has opened promising avenues for using the vast amounts of information entailed in granular financial data sets to assess financial risks.

To start with, they allow **new types of indicator** to be derived, as highlighted by the work presented by Vasilis Siakoulis (Bank of Greece) on the analysis of the financial strength of individual firms. Based on the granular information collected in the central bank's supervisory database (covering around 200,000 borrowers over a decade), a deep learning technique with a specific classification algorithm<sup>20</sup> was used to forecast the default for each loan outstanding. To facilitate policy monitoring work, this approach was complemented by a dimensionality reduction algorithm to reduce the number of variables to be considered.

Moreover, big data analytics can help to **enhance existing financial sector assessment processes**, by extending conventional methodologies and providing additional insights – in terms of eg financial sentiment analysis, early warning systems, stress-test exercises and network analysis. For instance, Sanjiv Das (Santa Clara University) and Kimmo Soramäki (Financial Network Analytics) presented a number of cases, including the use of network analytics for systemic risk measurement; the application of text analysis techniques to corporate e-mails and news for risk assessment; the measuring of interconnectedness to identify risk concentrations in CCPs and contagion effects; the identification of liquidity and solvency problems in payment systems; and the simulation of a financial institution's operational failure. These various experiences underlined the importance of having a sound theoretical framework to interpret the signals provided by disparate sources as well as to detect unusual, odd patterns in the data. They also highlighted the role played by model simplicity and transparency in the success of these initiatives, the benefit of a multidisciplinary approach, and the high IT and staff costs involved.

### 3. The use of big data in crafting central bank policy: organisation and challenges

Central bank experience suggests that the opportunities provided by big data sources and related analytical techniques can be significant, supporting a wide range of areas of policy interest. But how should central banks organise themselves to make the most of these opportunities? And what are the key challenges?

#### Organisational issues

Central banks' tasks cover a **wide range of topics** that can greatly benefit from big data. For instance, and as noted by Per Nymand-Andersen (ECB), central bankers need near-real-time and higher-frequency snapshots of the macro economy's state, its potential evolution (central scenario), and the risks associated with this outlook (eg early warning indicators and assessment of turning points). At the same time, their interest in financial stability issues calls for the ability to zoom in and get insights at

<sup>20</sup> eXtreme Gradient Boosting (XGBoost), which is commonly used in decision tree-based algorithms; see Chen and Guestrin (2016) and <https://xgboost.readthedocs.io/en/latest/>.

the micro level – see the ongoing initiatives among European central banks to develop very granular data sets on security-by-security issuance and holdings as well as on loan-by-loan transactions (the AnaCredit project; Schubert (2016)).

This puts a premium on **information systems** that can support this diversity of approaches. The presentation by Renaud Lacroix (Bank of France) argued that this requires a multidisciplinary and granular data platform, to supply flexible and innovative services to a wide range of internal users. The data lake platform project being developed at the French central bank will provide key data management services underlying multiple activities, covering data collection, supply (access), quality management, storage, sharing, analytics and dissemination. As presented by Robert Kirchner, the Deutsche Bundesbank has set up an integrated microdata-based information and analysis system (IMIDIAS) to facilitate the handling of granular data used to support its activities. It has also worked on fostering internal as well as external research on this information to gain new insights and facilitate policy analysis. Moreover, the Bundesbank actively supports the International Network for Exchanging Experience on Statistical Handling of Granular Data (INEXDA) (Bender et al (2018)).

A key takeaway was that the development of an adequate information system is only one element of a more **comprehensive strategy** to make the most of big data at central banks. As presented by David Roi Hardoon (Monetary Authority of Singapore), the “MAS Digital Supervision” initiative relies heavily on the use of machine learning, text-mining, natural language processing and visualisation techniques; and it is also backed by an extensive staff training programme on data analytics. Another example provided by Iman van Lelyveld relates to data science at the Netherlands Bank. Several use cases have been developed there – eg in the areas of credit risk, contagion risk, CCP risk, and stress testing in specific market segments. A key outcome has been the recognition of the important role played not only by the techniques used but also by the staff, organisation and culture.

## Challenges and limitations

In practice, important challenges remain, especially in handling and using big data sources and techniques.

**Handling big data** can be resource-intensive, especially in collecting and accessing the information, which can require new, expensive IT equipment as well as state-of-the-art data security. Staff costs should not be underestimated too, as suggested by the experience reported for the Bank of Japan. First, large micro data sets on financial transactions often have to be corrected for false attributes, missing points, outliers etc (Cagala (2017)); this cleaning work may often require the bulk of the time of the statisticians working with these data. Second, a much wider set of profiles – eg statisticians, IT specialists, data scientists and also lawyers – are needed to work in big data multidisciplinary teams; ensuring a balanced skillset and working culture may be challenging. Third, there is a risk of a “war for talent” when attracting the right candidates, especially vis-à-vis private sector firms that are heavily investing in big data; public compensation and career systems may be less than ideally calibrated for this competition. In addition, and as seen above, a key organisational consideration is how to integrate the data collected into a coherent and comprehensive information model. The challenge for central bank statisticians is thus to make the best use of available data that were not originally designed for specifically



statistical purposes and can be overwhelming (with the risk of too much “fat data”, and too little valuable information). In most cases, this requires significant preparatory work and sound data governance principles, covering data quality management processes (eg deletion of redundant information), the setup of adequate documentation (eg metadata), and the allocation of clear responsibilities (eg “who does what for what purpose”) and controls.

**Using big data** is also challenging for public authorities. One key limitation relates to the underlying quality of the information as noted above. This challenge can be reinforced by the large variety of big data formats, especially when the information collected is not well structured. Moreover, big data analytics rely frequently on correlation analysis, which can reflect coincidence as well as causality patterns.<sup>21</sup> Furthermore, the veracity of the information collected may prove insufficient. Big data sets may often cover entire populations, so by construction there is little sampling error to correct for, unlike with traditional statistical surveys. But a common public misperception is that, because big data sets are extremely large, they are automatically representative of the true population of interest. Yet this is not guaranteed, and in fact the composition bias can be quite significant, in particular as compared with much smaller traditional probabilistic samples (Meng (2014)). For example, when measuring prices online, one must realise that not all transactions are conducted on the internet. The measurement bias can be problematic if online prices are significantly different from the prices observed in physical stores, or if the products consumers buy online are different to those they buy offline.

These challenges are reinforced by two distinctive features of central banks – the first being their independence and the importance they accord to preserving public trust. Since the **quality** of big data sets may not be at the standard required for official statistics, “misusing” them as the basis of policy actions could raise ethical, reputational as well as efficiency issues. Similarly, if the **confidentiality** of the data analysed is not carefully protected, this could undermine public confidence, in turn calling into question the authorities’ competence in collecting, processing and disseminating information derived from big data as well as in taking policy decisions inferred from such data. This implies that central banks would generally seek to provide reassurance that data are used only for appropriate reasons, that only a limited number of staff can access them, and that they are stored securely. The ongoing push to access more detailed data (often down to individual transaction level) reinforces the need for careful consideration of the need to safeguard the privacy of the individuals and firms involved.

A second feature is that central banks are policymaking institutions whose actions are influencing the financial system and thereby the information collected on it. Hence, there is a feedback loop between the financial big data collected, its use for designing policy measures, and the actions taken by market participants in response. As a result, any move to measure a phenomenon can lead to a change in the underlying reality, underscoring the relevance of the famous Lucas critique for policymakers (Lucas (1976)).

<sup>21</sup> See the words of Stephen Jay Gould as quoted in Per Nymand-Andersen’s presentation: “*the invalid assumption that correlation implies cause is probably among the two or three most serious and common errors of human reasoning*”.



## Looking forward

The workshop and seminar provided a unique opportunity to take stock of the implementation of big data projects in the central bank community. These show that new big data-related sources of information and analytical techniques can provide various benefits for policymakers. Yet big data is still seen as complementing rather than replacing present statistical frameworks. It raises a number of difficult challenges, not least in terms of accuracy, transparency, confidentiality and ethical considerations. These limitations apply to big data sources as well as to the techniques that are being developed for their analysis. In particular, one major drawback of big data analytics is their black-box character, a difficulty reinforced by the frequent use of fancy names even for simple things (“buzzwords”). This can be a challenge for policymakers who need to communicate the rationale behind their analysis and decisions as transparently as possible. Moreover, **important uncertainties** remain on a number of aspects related to information technology and infrastructures, such as the potential use of cloud-based services and the development of new processes (eg cryptography, anonymisation techniques) to facilitate the use of micro-level data without compromising confidentiality.

One important point when discussing these issues is that **central banks do not work in isolation**. They need to explain to the general public how the new data can be used for crafting better policies, say, by providing new insights into the functioning of the financial system, clarifying its changing structure, improving policy design, and evaluating the result of policy actions (Bholat (2015)). But they also need to transparently recognise the associated risks, and to clearly state the safeguards provided in terms of confidentiality protection, access rights and data governance. Ideally, if big data is to be used for policymaking, the same quality of standards and frameworks that relate to traditional official statistics should be applied, such as transparency of sources, methodology, reliability and consistency over time. This is will be key to facilitating a greater use of this new information as well as its effective sharing between public bodies.<sup>22</sup>

Looking forward, it is still unclear whether and how far big data developments will trigger a change in the **business models of central banks**, given that they are relatively new to exploiting this new type of information and techniques. As noted in the presentation by Bruno Tissot (BIS), central banks have historically focused more on analysing data and less on compiling them. They are now increasingly engaged in statistical activities, reflecting the data collections initiated after the GFC as well as the growing importance of financial channels in economic activities – see, for instance, the substantial involvement of central banks in the compilation of financial accounts, a key element of the SNA framework (van de Ven and Fano (2017)). As both data users and data producers, they are therefore in an ideal position to ensure that big data can be transformed into useful information in support of policy.

<sup>22</sup> On the general data-sharing issues faced by central banks, see IFC (2016b).

## References

- Baker, S, N Bloom and S Davis (2016): "Measuring economic policy uncertainty", *Quarterly Journal of Economics*, vol 131, no 4, pp 1593–636.
- Bean, C (2016): *Independent review of UK economic statistics*, March.
- Bender, S, C Hirsch, R Kirchner, O Bover, M Ortega, G D'Alessio, L Teles Dias, P Guimarães, R Lacroix, M Lyon and E Witt (2016): "INEXDA – the Granular Data Network", *IFC Working Papers*, no 18, October.
- Bese Goksu, E and B Tissot (2018): "Monitoring systemic institutions for the analysis of micro-macro linkages and network effects", *Journal of Mathematics and Statistical Science*, vol 4, no 4, April.
- Bholat, D (2015): "Big data and central banks", Bank of England, *Quarterly Bulletin*, March.
- Blei, D, A Ng and M Jordan (2003): "Latent Dirichlet Allocation", *Journal of Machine Learning Research*, pp 993–1022.
- Cagala, T (2017): "Improving data quality and closing data gaps with machine learning", *IFC Bulletin*, no 46, December.
- Carnot, N, V Koen and B Tissot (2011): *Economic Forecasting and Policy*, second edition, Palgrave Macmillan.
- Cavallo, A and R Rigobon (2016): "The Billion Prices Project: Using online prices for measurement and research", *Journal of Economic Perspectives*, Spring 2016, vol 30, no 2, pp 151–78.
- Chen, T and C Guestrin (2016): *Xgboost: A scalable tree boosting system*, Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp 785–94.
- Cielinska, O, A Joseph, U Shreyas, J Tanner and M Vasios (2017): "Gauging market dynamics using trade repository data: the case of the Swiss franc de-pegging", Bank of England, *Financial Stability Papers*, no 41, January.
- Coeuré, B (2017): "Policy analysis with big data", speech at the conference on "Economic and financial regulation in the era of big data", Bank of France, Paris, November.
- Committee on Payments and Market Infrastructures (CPMI) and Board of the International Organization of Securities Commissions (IOSCO) (2018): *Framework for supervisory stress testing of central counterparties (CCPs)*, April.
- European Commission, International Monetary Fund, Organisation for Economic Cooperation and Development, United Nations and World Bank (2009): *System of National Accounts 2008*.
- Financial Stability Board (2011), "Understanding financial linkages: a common data template for global systemically important banks", *FSB Consultation Papers*.
- (2017): *Artificial intelligence and machine learning in financial services - Market developments and financial stability implications*, November.
- Financial Stability Board and International Monetary Fund (2009): *The financial crisis and information gaps*.

——— (2015): *The financial crisis and information gaps – Sixth Implementation Progress Report of the G20 Data Gaps Initiative*.

Groves, R (2011): "Designed data and organic data", in the Director's Blog of the US Census Bureau, [www.census.gov/newsroom/blogs/director/2011/05/designed-data-and-organic-data.html](http://www.census.gov/newsroom/blogs/director/2011/05/designed-data-and-organic-data.html).

Haldane, A (2018): "Will big data keep its promise?", speech at the Bank of England Data Analytics for Finance and Macro Research Centre, King's Business School, 19 April.

Hammer, C, D Kostroch, G Quirós and staff of the IMF Statistics Department (STA) Internal Group (2017): "Big data: potential, challenges, and statistical implications", *IMF Staff Discussion Notes*, no 17/06, September.

Hansen, S (2018): "Introduction to text mining", lecture delivered at the workshop on "Big data for central bank policies", Bank Indonesia, Bali, 23–25 July.

Hill, S (2018): "The big data revolution in economic statistics: waiting for Godot ... and government funding", *Goldman Sachs US Economics Analyst*, 6 May.

Irving Fisher Committee on Central Bank Statistics (IFC) (2015): *Central banks' use of and interest in 'big data'*, IFC Report, October.

——— (2016a): "Combining micro and macro statistical data for financial stability analysis", *IFC Bulletin*, no 41, May.

——— (2016b): *The sharing of micro data – a central bank perspective*, IFC Report, December.

——— (2017): "Big data", *IFC Bulletin*, no 44, September.

——— (2018): *Central banks and trade repositories derivatives data*, IFC Report, October.

Laney, D (2001): "3D data management: controlling data volume, velocity, and variety", META Group (now Gartner).

Loughran, T and B McDonald (2011): "When is a liability not a liability? Textual analysis, dictionaries, and 10-Ks", *Journal of Finance*, vol 66, no 1, pp 35–65.

Lucas, R (1976): "Econometric policy evaluation: A critique", *Carnegie-Rochester Conference Series on Public Policy*, vol 1, no 1, pp 19–46.

Meeting of the Expert Group on International Statistical Classifications (2015): *Classification of Types of Big Data*, United Nations Department of Economic and Social Affairs, ESA/STAT/AC.289/26, May.

Meng, X (2014): "A trio of inference problems that could win you a Nobel Prize in statistics (if you help fund it)", in X Lin, C Genest, D Banks, G Molenberghs, D Scott and J-L Wang (eds), *Past, present, and future of statistical science*, Chapman and Hall, pp 537–62.

Nyman-Andersen, P (2015): "Big data – the hunt for timely insights and decision certainty: Central banking reflections on the use of big data for policy purposes", *IFC Working Papers*, no 14.

Petropoulos, A, V Siakoulis, E Stravroulakis and A Klamargias (2018): "A robust machine learning approach for credit risk analysis of large loan level data sets using deep learning and extreme gradient boosting", paper presented at the workshop on "Big data for central bank policies", Bank Indonesia, Bali, 23–25 July.

Rigobon, R (2018): "Promise: measuring from inflation to discrimination", presentation given at the workshop on "*Big data for central bank policies*", Bank Indonesia, Bali, 23–25 July.

Schubert, A (2016): "AnaCredit: banking with (pretty) big data", *Central Banking Focus Report*.

Tetlock, P (2007): "Giving content to investor sentiment: the role of media in the stock market", *Journal of Finance*, vol 62, no 3, pp 1139–68.

The Economist (2017): "The world's most valuable resource is no longer oil, but data", 6 May edition.

Tissot, B (2019): "Making the most of big data for financial stability purposes", in S Strydom and M Strydom (eds), *Big Data Governance and Perspectives in Knowledge Management*, IGI Global, pp 1–24.

van de Ven, P and D Fano (2017): *Understanding Financial Accounts*, OECD Publishing, Paris.

Zulen, A and O Wibisono (2018): "Measuring stakeholders' expectations for the central bank's policy rate", paper presented at the workshop on "*Big data for central bank policies*", Bank Indonesia, Bali, 23–25 July.