# Getting the foundations right

Helen MacGillivray
*Queensland University of Technology, Brisbane, QLD, Australia*
*E-mail: h.macgillivray@qut.edu.au*

**Abstract.** There has been increasing interest in recent years in training in official statistics with reference to the 2030 Agenda, big data, diversification of data types and sources, and data science. Backgrounds for work in official statistics are becoming more varied than ever. The official statistics community has also become progressively more aware of the importance of statistical literacy in education and trust in official statistics. Hence foundation and introductory are of as much interest to official statistics as more specialised training. At the same time, greater access to data and vast technological capabilities have seen much emphasis and discussion of the statistical and data sciences and education therein, including development of educational resources in contexts such as civic data and statistics. Data science provides opportunities to renew the decades-long push for authentic learning that reflects the practice of 'greater statistics' and 'greater data science', and to examine progress to date in implementing and sustaining the extensive work and advocacy of many. This article discusses what is needed at the foundation and introductory levels to realize this advocacy, with commentary relevant to official statistics.

Keywords: Teaching statistics, teaching data science, statistical literacy, data investigations, introductory

## 1. Introduction

At the United Nations Statistical Commission (UNSC) in March, 2019, the side event of the International Statistical Institute (ISI) was titled *Preparation and Skilling for Official Statistics Careers*, with panellists Gabriella Vukovich (President, Hungarian Central Statistical Office), Steve McFeely (then Head of Statistics & Information, UNCTAD), Alex Shimuafeni (Statistician General, Namibia), Helen MacGillivray, President ISI), and discussant Haishan Fu (Director, Development Data Group, The World Bank Group). The session considered what type of generic skills and workplace preparation at the undergraduate level are of assistance for official statistics workplaces, whether and how these fit in university undergraduate programs, how to increase undergraduate understanding of what is involved in official statistics careers, and awareness of the longer-term career possibilities. Despite the early hour and freezing ($-9°$C) conditions, there was standing room only and request for more sessions, even conferences, on this theme.

There has been significant development of training for official statistics at the postgraduate level. More than 20 European universities participate in the Euro-

pean Master in Official Statistics. The wide variety of training of National Statistical Office (NSO) staff by global and regional institutions to strengthen statistical capacities was a key motivator in establishing GIST, the Global Network of Institutions for Statistical Training, GIST – https://unstats.un.org/GIST. Participants in this type of training and training at the postgraduate level, are usually working, or have worked, mostly in NSO's or closely associated government departments or agencies. Hence we see considerable in-service and postgraduate training developments, and emerging linkages at school level to increase awareness of official statistics, but what of the undergraduate level? Many NSO's put considerable effort into attracting graduates from suitable quantitative degree programs, including statistics, finance, economics, and mathematics, and some have established significant ongoing research and/or staff links with selected university departments or faculties. But official statistics is just one of innumerable destinations for such graduates, and is very different to many, not just in content but also in type of workplace.

At school level, there has been increasing collaboration between NSO's and those involved in school education, most notably through competitions such as the ISI's ISLP (International Statistical Literacy

Project) Poster competition (http://iase-web.org/islp/Poster_Competition_2018-2019.php). There is also increasing emphasis on school statistics teaching that enables students to engage with current social issues, such as ProCivicStat (http://iase-web.org/islp/pcs/), funded with support from the European Commission.

On the establishment of GIST at the 2018 UNSC, as inaugural chair I discussed with leaders in the UN Statistical Division (UNSD) how to tackle the ambitious and far-reaching aims in the projected strategic action areas, and proposed task forces of GIST members, with leaders selected by each task force on a GIST Board to facilitate sharing and liaison across GIST members and the GIST Stakeholder Advisory Group of NSO's covering world regions. One of the strategic action areas is "Promote training that enables data producers to improve data literacy and usage within stakeholder communities" and there has been a task team on statistical literacy, now called "Statistical Literacy in the context of the 2030 Agenda" since GIST's inception. There was considerable initial discussion within this task team on whether the focus was intended to be on improving statistical literacy in the general community, within government, or to improve statistical communication by NSO staff or all three. Indeed, this task team has done substantial work on the first two areas in particular, with significant support and input from UNITAR, as discussed elsewhere in this special issue.

All the above are indicators of the importance of statistics and data education for the world of official statistics, and the rapid increase of this importance with the 2030 agenda, with new sources and types of data, new gatherers and sources of information (and misinformation), increasing awareness of and emphasis on data science, and increasing concerns on issues such as trust in official statistics. Specific training of official statisticians is just one part of the entire spectrum of statistical and data science education which is pertinent to the core business of official statistics. Previously perceived precincts within the entirety of statistical, data and stochastic sciences, including official statistics, have become increasingly indistinct and vanishing as more and more complex real contexts, data and technological capabilities arise across disciplines, business, government and society. Hence sound foundations in statistical and data science education across school and tertiary levels are of consequence to official statistics from citizen statistical and data literacy right through to training of official statisticians. This article discusses a range of issues both general and specific in such foundations, including progress and problems, the contrasts

between extensive advocacy and examples of best practice with details of practicalities and realities across education, and noting points of particular interest to official statistics, and where official statistics leadership can assist.

## 2. Statistical and data literacy

### 2.1. Data literacy descriptions are the same as those for statistical literacy

There have been many descriptions of statistical literacy over many years, with examples including:

> ... good "statistical citizens," [are] able to consume the information that they are inundated with on a daily basis, think critically about it, and make good decisions based on that information [1]. People's ability to interpret and critically evaluate statistical information and data-based arguments appearing in diverse media channels, and their ability to discuss their opinions regarding such statistical information [2].

Like the above, many have referred to information in media sources. Jane Watson [3] initially developed a view of statistical literacy that centred on media reports and focused on the data consumer. The extensive and prolific work of Jane and her many colleagues, particularly Rosemary Callingham, include research, guidance and authentic classroom-ready activities developing statistical literacy (see publications under https://www.utas.edu.au/profiles/staff/education/Jane-Watson). Some recent excellent examples can be found at https://www.abc.net.au/education/media-literacy/a-guide-to-statistical-literacy-in-the-classroom/12789152. Here, Watson and Callingham describe statistical literacy as follows.

> Statistical literacy is the critical thinking needed when encountering claims based on statistics or data. It requires understanding of:
> – the terminology and representation used – what does it mean statistically?
> – the context – what does the terminology and representation mean in the context where it is presented?
> – critical thinking – what is the precise claim being made and is it reasonable?

In this blog, they also make use of the classic five steps of statistical data investigations not as a recipe

but as a framework for "*critical questions for students to ask of statistical claims in news media.*" This is of particular importance in considering data investigations in education – as discussed further below – as well as highly relevant to current concerns of misinformation, and especially for civic and official statistics.

For many years, the Statistics Department at the University of Auckland have successfully run a course https://www.auckland.ac.nz/en/study/study-options/undergraduate-study-options/general-education/course-descriptions/stats-150g.html called *Lies, Damned Lies, and Statistics* on statistical literacy designed to prepare everyone, regardless of statistical background, to become critical consumers of statistical information. The achievement of such a course, including its sustainability, lies not only in the deep understanding of, and commitment to, statistical thinking by its developers and deliverers, but also in its authentic experiential learning and workable assessment as outlined by Budgett and Pfannkuch [4].

There have been ambitions and attempts over many years to achieve statistical literacy for all at tertiary level. One such university was University of Wollongong which initiated a number of statistical literacy courses, for example, for law [5]. Statistics in this university had, and has, close ties with the Australian Bureau of Statistics. In 2017, their program was encouraged for all, with the following description:

> *The University of Wollongong statistical literacy programme is divided into the following 3 modules.*
>
> > *Module 1. Producing data*
> > *Module 2. Describing, Clarifying and Presenting Data*
> > *Module 3. Interpreting data*
>
> *Completing these modules will help you to develop the skills you need to:*
>
> – *look behind the data with which you are presented at University and in your everyday experiences,*
> – *ask why these data are being presented in those forms,*
> – *ask what questions can be answered or what arguments are being made with these data.*

Today in 2021, the above is an online section of a set of tertiary literacies run by the university's Learning Development, https://www.uow.edu.au/student/support-services/tertiary-literacies/, but whose description is available only through online registration.

In a number of forums, including editorials for the journal Teaching Statistics, an International Journal for Statistics and Data Science Teaching, I incorporated references to the UN World Data Forums (WDF) of 2017 and 2018 (for example, [6,7]). Some presentations at the first UN WDF in 2017 not only included long-established descriptions for statistical literacy as 'data literacy' but also denial that it was statistical literacy. In addition, a description of the statistical data investigation process (see Section 3) was renamed as "*Opportunities for engagement" in data literacy.* Further investigation revealed that the contexts for these were oriented to the Sustainable Development Goals (SDG's), statistical empowerment, inclusivity and "The desire and ability to constructively engage in society through and about data" as the following demonstrate:

> *"Being able to relate to data is being able to relate to power and ask the important questions." Who owns the data? Who is impacted by the data? What does data inclusion mean?*

http://datapopalliance.org/item/entering-the-age-of-data-a-focus-on-data-inclusion/.

In a workshop on data literacy http://datapopalliance.org/beyond-data-literacy-workshop/ Jeanne Bourgault, President and CEO of Internews, summarized that "*data literacy is really about the ability to create, share, consume, produce, and analyze for their own self*," and it is essential to build local capacity of data literacy so that individuals may make informed choices with data.

Such descriptions are highly relevant to official statistics and the 2030 Agenda, but it must be recognised and emphasized that these are part of statistical literacy; NSO's, GIST and UNSD play crucial roles in such emphasis.

In arguing that "*data literacy is statistical literacy*", Gould [8] also argues that previous descriptions of statistical literacy need augmenting to recognise the widening sources and nature of data as well as the increasingly sophisticated technology for data. Gould's [8] description is more extensive than previous, and possibly for many ventures into data science capabilities, but it does reflect the living and growing nature of the statistical and data sciences in being able to tackle more complex and more diversified data and problems for all aspects of data, including sourcing, querying, wrangling, dissecting, handling, exploring, processing, presenting, analysing, interpreting.

Some of those who have denied that data literacy is statistical literacy have viewed statistical literacy as capabilities such as being able to produce and use graphs

and statistical measures, but it is clear from the extensive literature, of which the above are just a few examples, that the common essential element of data and statistical literacy is to gain confidence as a critical consumer of data and statistical information. This naturally merges into one's statistical and data capabilities, with no clear-cut delineation, and, as always with statistics, is context-dependent. As Moore and Cobb [9] famously said "Data are numbers in a context". The extensive literature in statistics education also demonstrates that significant effort has been put in by many world-wide to embed statistical literacy in school education, in citizen/adult education and, more rarely, in learning support across disciplines at tertiary level.

### 2.2. CensusAtSchool and ISLP

Official statistics has been involved in such efforts in a number of ways. Following an idea to raise statistical awareness in primary schools in New Zealand by Forbes [10], the online *CensusAtSchool* project was developed by the then Royal Statistical Society's Centre for Statistical Education (RSSCSE) (disbanded in 2014) to coincide with the 2001 population census in the UK [11]. The RSSCSE and the *CensusAtSchool* project were supported by the UK's Office of National Statistics. The project, including its IT and activity resources, was subsequently also adopted by New Zealand, Australia, South Africa, Canada, Italy, Ireland and the USA. Each country's NSO participated, in some cases in a major way, to support this particularly successful activity that promoted the collection, analysis and communication of data for, and on behalf of, school students in primary and secondary schools.

Since 2010, the International Association for Official Statistics (IAOS) has been formally involved with the International Statistical Literacy Project (ISLP) http://iase-web.org/islp/. An ISI committee was established in 1994 to stimulate the spread of quantitative skills around the world in areas and populations especially in developing countries and among the young. The committee was known as the World Numeracy Program Advisory Committee (ISI-UNESCO) and was chaired by Prof. Luigi Biggeri of Italy. This program initiated projects, including the Multilingual Glossary of Statistical Terms, World Statistics Day, Statisticians of the Centuries, a committee on Professional Ethics, Assessment Challenges. In 2000, the ISI moved responsibility for this program to the International Association for Statistical Education (IASE) and in 2002 the name was changed to ISLP, and a website of resources

was created. In 2007, a pilot statistical literacy competition was held in Portugal, and in 2008–2009, an International Statistical Literacy Competition was conducted in written examination formats, culminating in a final examination in conjunction with ISIBALO with 10 countries taking part in Durban in 2009, but organised independently of IASE.

Following the 2009 ISI WSC, the IASE president re-organised the ISLP project with a formal and sustainable governance model consisting of an ISLP executive of director and deputy directors, an Advisory Board with IASE and IAOS representation, and country coordinators appointed by the executive, currently representing more than 80 countries. IAOS and IASE members provide significant leadership and support in all ISLP activities, including the ISLP International Poster Competition. Reija Helenius, Statistics Finland, has been ISLP Director since 2010, and Pedro Campos, University of Porto/Statistics Portugal a deputy director. Sharleen Forbes, Statistics New Zealand (SNZ), was a deputy director until her retirement, and was replaced by Steve MacFeely, UNCTAD and now WHO. It was decided that an international poster competition was more suitable to statistical literacy than examinations, and since 2010, the ISLP International Poster Competition has been held biennially, with international winning posters displayed and announced at ISI WSC's, and typically involving more than 12,000 students from more than 20 countries.

The *CensusAtSchool* project produced classroom-ready resources and was more than statistical literacy, and although the ISLP Poster Competition is essentially extra-curricular, it also applies and builds on student statistical and data learning in school curricula to at least some extent. Many NSOs offer access to, or specifically create, data sets for use in all levels of statistics education, aiming to provide resources that also in improving awareness and understanding of official statistics. However, as Forbes et al. [12] emphasize:

> *".. to ensure that these data products are accessible, interesting, valued and engaged with, requires that official statistics agencies and statistical literacy educators work together to inform the education community about these products and how to use them effectively in their everyday teaching."*

Forbes et al. [12] includes an excellent account of how SNZ worked collaboratively over many years with academics and teachers on developments ranging over statistical literacy resources, school curricula and assessment, to initiatives and eventually courses to raise

the statistical capabilities of employees across government. Elsewhere in this special issue, Sharleen Forbes and John Harraway discuss how this led to three free downloadable web apps in official statistics hosted on the ISLP website, and to supplying initial input for the United Nations Institute of Training and Research (UNITAR) e-learning course *Understanding data and statistics better – for more effective SDG decision making*, as discussed also in this issue by Elena Proden. It is important to stress that all this work is described, by both SNZ and UNITAR, as raising statistical capabilities. Although descriptions of statistical and data literacy depend to some extent on contexts, including educational level, to overuse the term, and to apply it beyond its remit into the domain of statistical capabilities, are counter-productive.

Thus we see how official statistics have collaborated, and are collaborating, with educators and academics to help support development of statistical literacy and natural extensions into improving statistical capabilities, but collaboration is key and direct impact such as SNZ in school and tertiary courses is rarer. Section 5 includes some comments and ideas on how to effect better influence and achieve greater educational impact.

## 3. Statistical data investigations and data science

### 3.1. The statistical and data sciences

At the second UN WDF, leading data science speakers from large organisations in communications, securities, information technologies, and official statistics, discussed data science as "*a label for work being done for years*", and as requiring a "*diverse collaborative team*" but "*all of whom need statistical foundation*" (https://undataforum.org/WorldDataForum/sessions/ta6-08-data-scientist-what-are-they/). Other interesting comments from the panel included that data scientists need:

- *to know what can and can't be done with data;*
- *curiosity in problems; identifying and posing problems so they can be tackled; investigative and problem-solving; bring together data, data problem-solving, technical;*
- *to know that data science may produce what, but statistics gives why/understanding.*

In recent years, the term 'big tent' of statistics has been increasingly used [13] to encompass the full spectrum of statistical science and statistical practice across everything to do with data, variation and uncertainty.

This builds on Chambers [14] 'greater statistics' which in turn builds on Tukey's [15] view of data analysis which is the forerunner of the confluence of the statistical and data sciences. Donoho [16] not only advocates 'greater statistics' but also 'greater data science', and initiates a vision for the latter that is far more than a 'mere scaling up to big data' and big technology, but an ongoing 'more intellectually productive and lasting' science. The term data science has been around for a while (see, for example, [17]) including the idea of calling 'statistics' 'data science' as mooted by Wu in 1986 [18].

The above comments tend to be oriented to professional work, research and higher education, but also apply at the school, foundation and introductory levels. In his exploration of what introductory data science courses should look like, Gould [19] argues "if we are to teach [secondary and introductory level] students to find meaning in data, then most of what students need sits firmly within the boundaries of statistics."

In the above and other articles and comments, we see that discussion and debate about what data science is and whether it is statistics are continuations of previous discussions on what statistics is. Statistics has always been part of developments in computing power, as contributor, motivator and user, and has fed into, and been fed by, increasing technological power to tackle more and more complex problems in wider contexts with broader as well as 'larger' data. Big data, data analytics, data deluges have been impacting on statistics as well as drawing data science 'out of the back room'.

Discussions and descriptions of what the statistical and data sciences are, are valuable in communicating, but to attempt to divide by internal boundaries is counter-productive. Certainly, the statistical and data sciences span a wide and increasingly broad range and diversity of topics and capabilities, but over the past 50 years, previous internal precincts in statistics have blurred into each other, often combining forces in tackling more complex real contexts and data. Topics are not boundaries. Indeed, it is time to finally remove any previous boundaries within statistics and to cogently advocate avoidance of such in what should now be known as the big tent of the statistical and data sciences. Official statistics is very much within this big tent, and any previous boundaries with other areas of statistics have become well-blurred in the era of big data, diversity of data types, provenance and ownership, and SDG's.

### 3.2. The statistical data investigation process

The statistical and data sciences are essentially investigative, problem-solving and driven by contexts

involving data, variation and hence uncertainty. Over many years, statisticians have advocated that students have authentic experience of the full statistical investigation process. In augmenting Chambers' [14] description of how statisticians practice '*greater statistics*', Cameron [20] describes this process as follows, commenting that such training is an appropriate foundation for most statisticians wherever they may be employed:

- *formulating a problem so that it can be tackled statistically*
- *preparing data (including planning, collecting, organising and validating)*
- *analysing data*
- *presenting information from data*
- *researching the interplay of observation, experiment and theory.*

Kenett and Thyregod [21] criticize university teaching that does not include sufficient focus on the first two and last two steps of their statistical consulting cycle described as follows:

- *problem elicitation*
- *data collection and/or aggregation*
- *data analysis using statistical methods*
- *formulation of findings, their consequences and derived conclusions*
- *presentation of findings and conclusions/recommendations.*

These and other descriptions of the full statistical investigation process have been used by leading international statistical educators in advocacy and initiatives across curricula, resources and guidance for statistics learning to reflect the practice of statistics. Those in statistical education will recognise the similarities to Wild and Pfannkuch's [22] popularization of [23] and the [24] stages of Question, Design, Collection, Analysis, Answer. The description of the data-handling cycle that featured in the UK National School Curriculum in the mid-seventies [25] became the PCPD (Plan, Collect, Process, Discuss) cycle.

Emphasis on embedding authentic experience of the statistical investigation process, was part of extensive work during the past 2–3 decades by statisticians and statistics educators worldwide in initiating a variety of changes:

- in teaching statistics at university, particularly introductory levels across disciplines, and at school level;
- in professional development workplaces and communities; and
- in statistics education research.

This work incorporated advocacy of:

- Data-driven concepts and statistical thinking;
- Real, 'large' contexts and data, treating simple concepts and procedures within complex;
- Statistics in its own right;
- Technological and data systems know-how;
- Student ownership and constructivism.

The American Statistical Association's Guidelines for Assessment and Instruction in Statistical Education (GAISE) at school and college levels, https://www.amstat.org/asa/education/Guidelines-for-Assessment-and-Instruction-in-Statistics-Education-Reports.aspx, have come about from, and are based on, this extensive work since the early 1990's.

All of the above, updated to reflect technological and data advances, have found their way into data science advocacy. The International Data Science in Schools Project (IDSSP) www.idssp.org, is a consortium of statisticians and computer scientists whose collaboration has produced curriculum frameworks intended to guide the writing, resourcing and teaching of data science curricula at senior school or introductory tertiary levels. PPDAC has been slightly edited in the IDSSP framework and renamed "the basic cycle of learning from data", without any referencing. Most of Unit 1 of this framework is familiar to leading statistical educators, and the approach and key elements of Unit 1 can be described as:

- Authentic and embedded learning of what has been called the statistical problem-solving process, the statistical/data investigation cycle, and has been renamed the data science learning cycle in Unit I.
- Authentic original contexts and data, with the simple embedded in the complex, and collecting/accessing data relevant to students' lives.
- Use of technology as in the *practice* of statistics.
- Multivariable contexts and data, with (again) the simple embedded in the complex and the emphasis on types of variables, and moving to more than two variables/multivariable data as quickly as possible.
- Visualisation and exploration.
- Student-centred learning.
- Hands-on learning of data acquisition and handling in the data-handling pipeline.
- The above applying in all introductory data science learning across disciplines.

The second last bullet point above reflects embedding more of the data science approach relevant to today. Unit 2 moves on to more statistical and data science sophistication.

### 3.3. Data investigation at foundation and introductory levels across disciplines

A challenge in foundation and introductory data science curricula is to ensure computing does not overshadow, or worse dominate, the learning from data. Because both mathematics and computer science serve statistics and data science, it is essential at foundation and introductory levels to ensure the focus is on development of statistical and data science learning. Burr et al. [26] provide invaluable discussion on how to gradually build computational capabilities by stealth. Gould [19] uses the experiences of developing and delivering a secondary data science course that pre-dates the IDSSP, the Mobilize *Introduction to Data Science* (IDS) course, to demonstrate how data scientific thinking has "*a strong core of statistical thinking, carefully selected components of computational thinking, and just a dash of mathematical thinking.*"

In describing the IDS Data Cycle (Ask questions, Consider data, Analyse data, Interpret data) as a template for the statistical investigative process, Gould [19] emphasizes the change from the older Collect data of codifications such as PPDAC and PCPD, to Consider data. This change was also made in the 2020 revision of GAISE preK-12 [27]. Such emphasis is of paramount importance in today's statistics and data science, and especially for official statistics. Almost the whole of Rubin's [28] excellent discussion of the Data Clubs project for young people ages 12–15 is on the value and interest for students in considering how data are generated as part of a larger paradigm shift from previous eras to today's statistics and data science.

Asking "who, how, when, where, and why?", as well as "can we measure what we want to measure?" have long been fundamental in the practice of statistics, but take on new significance in education now that technology enables students to access and explore large and complex data sources as well as collect diverse types of data themselves as described by Gould [19] and Rubin [28]. This is also of core focus in ProCivicStat (http://iase-web.org/islp/pcs/) which builds on and greatly extends prior work, such as that of the Smart Centre at Durham University, in providing educational resources using civic and government data. Embedding such work in mainstream curricula will be invaluable for official statistics, including the building of trust in official statistics that has received so much recent focus and concern.

My own journey with regard to student data investigations illustrates the journey as technological power increased and gradually became more accessible to students. In the 1980's, including use of statistical software in large introductory classes across disciplines, including health sciences, technology and engineering students, involved considerable logistic challenges. In the first half of the 1990's, asking "who, how, when, where, and why?", as well as "can we measure what we want to measure?" was core in my teaching across disciplines for the real datasets I provided to students.

From 1995 to 2011, I developed, implemented and embedded, with the increasing support of colleagues, student-based statistical data investigations as semester-long projects in parallel with the statistical knowledge and skills development in large (up to 600) introductory courses across disciplines and for statistics and mathematics majors. This strategy could most likely be classed as inquiry (or enquiry) oriented learning, but is essentially reflecting the practice of 'greater statistics'. It had some interesting similarities and contrasts with IOL in science [29]. The impact and value of this strategy reached beyond expectations not just in student learning and attitudes [30], but also into our curricula, teaching materials and resources, assessment, training of tutors and preparation of future statisticians for the workplace [31]. The emphasis started on students collecting their own data on issues of their choice with staff assistance in planning, but as time went on, there were increasing examples of students wishing to investigate data from other sources. The main problems they found were difficulties in accessing raw data and insufficient information on the "who, how, when, where, and why?". This is what has changed now with the vast array of technological capabilities in sourcing, accessing, scraping, wrangling and handling data in the official, scientific and public arena. So courses such as mine have needed to, and must, continue to evolve to reflect the practice of statistics and data science.

## 4. What are needed

### 4.1. Embedding authentic learning of data investigations

After reading the above, a reaction might be that, provided statistics and data science are recognised as a 'big tent' and work together, foundation and introductory learning looks sound, but unfortunately there is much that is not. The 'emergence' of data science is opportunity for the statistical and data sciences to greatly promote the understanding and advocacy outlined above,

but there are also significant lessons from the lack of penetration or sustainability of such advances and advocacy. At both school and university levels there needs to be quality information on the realities.

As just one recent example, [32], an international leader in statistics education, based in the US, spoke of doing some research about how countries are addressing statistics in the school curriculum, and asked for information in the Australian context. She mentioned her significant concerns, including that:

> "*In the US, while we have good documents on what should be happening, the reality is far from what those documents suggest – for a variety of reasons, most of the emphasis is on . . . content that is easy to assess.*"

At the university level, internships, work and clinical placements have long been embedded in professional programs, and across all programs there is now much good work focussing on work-integrated learning (WIL) and capstone projects in the final undergraduate years, replacing the vacation work experience or final year industry projects which received intermittent attention in past eras. In other disciplines, the roles of statistics and data science in such WIL and capstone courses depend critically on the foundations in those disciplines, and it is in the foundation and introductory levels that both general culture and curricula details need attention. This is important for official statistics not only because graduates of different disciplines go on to careers in official statistics but also, and possibly of greater significance, official statistics works with all of government as well as increasingly with business and industry, as well as needing informed citizens with trust in official statistics.

There needs to be greater value placed on statistical teaching expertise at the introductory level, by statistics, data science and all disciplines, accompanied by genuine sustained collaboration with other disciplines. Although teaching materials, resources and curricula details need nuancing for different disciplines, the essential of statistics and data science foundations are core to all, and no matter where statisticians are located in a university, establishing an active genuine community of practice in teaching statistics and data science enormously benefits efficiencies and effectiveness as well as student learning and staff morale and advancement. All universities should facilitate and support such a community.

The principles and practicalities of experiential learning of data investigations as discussed in Section 3 above, should be embedded in introductory statistics and data science, alongside well-scaffolded development of the relevant knowledge and skills. Beliefs that 'students won't do it right' and 'it's not serious enough' when students are encouraged to choose contexts and issues of interest to them are both misplaced and counterproductive. Fears of assessment workload are also misplaced. Data investigations are best done in groups because, as in the workplace, such investigations benefit from a group approach. Because data investigations both teach and assess the higher order statistical and data capabilities, thinking and usage, other forms of assessment can focus on knowledge and procedures in more easily-marked formats. In addition, staff involvement in advising students on their investigations during computer laboratory work throughout the course, builds a natural familiarity with the various investigations. It is also of interest to observe that multiple choice questions tend to be highly dependent on local culture/conditions and are course-specific, but criteria and standards for data investigations tend to be more universal, with exemplars which can be used across institutions and programs.

### 4.2. Curricula needs and cautions

Scrutiny of many introductory statistics textbooks illustrates that attempts to de-mathematicalize earlier introductory books (which were meant to be in mathematical statistics) without sufficient re-thinking of the purpose of the statistical 'story', lead to over-focus on new ways of teaching earlier content and not necessarily appropriate sequencing at the expense of developing data investigation skills, such as:

- Data: What? When? How? Limitations?
- Issues: What are we interested in? What can we investigate? What do we need?
  * Sources? Quality? Sufficient information? Access? Collect? Design?
- Turn research questions into statistical questions
- Identify variables and cases/subjects
- Do we need a pilot study/experiment?
- What do we need in data handling and preparation: organising, wrangling, checking, transferring, combining, coding, . . . , preparatory exploring

When students choose what to investigate, explore and source, they are motivated to find tools, they have ownership of data, context and questions; student ownership is the best motivator for learning.

Indeed there is need for developing statistical concepts and tools for exploration, visualisation, and analysis, but the following are needed:

- Framing of issues, identification of variables and understanding their types;
- Advantages and disadvantages of different visualisations, presentations;
- Understand what numerical codes can and cannot do, in order to prevent long-term incorrect use of types of data (unfortunately far too common in certain disciplines), for example,
  * Cannot turn nominal variables into numerical variables;
  * Cannot turn ordinal variables into continuous variables;
- Understand what aggregation is, its advantages and its limitations, from histograms to maps.
- Understand assumptions and how to evaluate assumptions graphically after models are fitted – far too many researchers in other disciplines ignore assumptions and neither use nor understand how to use graphical diagnostics;
- Learning to bring together findings in reporting using qualifications in language and identifying further issues;
- Real data and real contexts but
  * Contexts must not dominate statistical learning;
  * Contexts must be familiar/readily accessible to students;
  * Beware teacher-centred, top-down or context-complex case studies;
- Most importantly, **move to many variables and real empowerment as soon as possible.**

The above include cautions and indications of how to avoid foundation problems which have become self-perpetuating, particularly in other disciplines if students' introductory course does not include sufficient of the above statistical expertise. Below are some more specific cautions and problems which have not yet been tackled:

- fixation with restrictions to one and two variables, no matter what types of procedures are preferred by the instructor;
- isolated, single purpose clean data and questions and instructor-prescribed answer – the simple can be developed within the larger context or data;
- multiple procedures and forcing into discipline norms
  * the classic is the overuse of t;
- rigid, discipline-embedded approaches, top-down case studies, and too much orientation for research training – in any discipline, including statistics.

Scrutiny of many textbooks and introductory courses also highlights a very big problem requiring considerable attention, namely the need to reclaim and reform the teaching of probability and probabilistic thinking as integral to the statistical and data sciences. Probability must embed and be embedded in data, language and visualisation. Extensive student experience of the language of probability builds both familiarity and foundational understanding for statistics and data. This is particularly true of conditioning language in which familiarity is essential for understanding risk both as citizens and professionally. Conditional probability should be introduced before the special case of independence, and developed through data and estimates of conditional probabilities as well as through language. There are many examples in real contexts of misunderstanding of conditional probabilities, using inappropriate data for their estimation, and incorrect multiplication of probabilities, sometimes with appalling consequences. The term 'multiplication rule' should be banished forever. Tables are of particular importance in official statistics, and data on two and more categorical variables are ideal settings for engaging and invaluable learning and using conditional probabilities, including Bayes foundations, as well as splitting data, confounding and hidden variables.

### 4.3. Some points on school contexts

Much of the above also applies at school, gradually and more simply developed and experienced appropriately for the level, but in a slowly evolving coherent statistical 'story', with authentic student learning experiences at every stage, so that students own their foundation in probabilistic, statistical and data thinking for citizenship and further learning. There are many characteristics of schooling different to those of universities which must be taken into account. Firstly an obvious observation: who should teach statistics and data science as they should be taught, are those who have learnt it this way, demonstrating the importance of sufficient and appropriate statistical and data science learning in pre-service and in-service teacher education.

Comments on school education are too often generalised from the senior school context, but primary schooling has no specialisations, and middle school only some. Although many excellent resources have been developed for school levels, including extra-curricular activities, the big challenges lie in the 'parcelling up' of authentic statistical approaches to be embedded within classroom learning, activities, exercises

and a diversity of assessments, both formative and, later, summative. No matter where statistics (and now data science) is placed in curricula, we are speaking of a discipline that combines principles and procedures with the nuances of uncertainty and variation. The danger of viewing data science as merely up-scaling of the technology of data-handling is that coding and programming will have the same distorting effects on statistics and learning from data as through the eyes of other discipline-specific views.

There is an urgent need for more involvement of statistical and statistical teaching expertise in all aspects of schooling: curricula, educational authorities, pre-service, in-service, textbooks, assessment. Involvement of such in curricula may have improved, but implementation and sustainability need the full spectrum of involvement, as curricula interpretations, even when not ignored, depend critically on the user's background and understanding, especially in a discipline such as statistics. Textbooks are a major challenge, as are approaches which emphasize *the* question and *the* answer. Unless texts across year levels can become generally accepted, the statistics and data science sections of textbooks must be written by those with expertise in teaching authentic statistical and data science thinking, preferably with collaboration across expertise in educational levels.

## 5. Conclusion: How official statistics can help

A question from a number of participants at the UNSC 2019 ISI side event was, should there be bachelor degrees in official statistics? University authorities generally tend to be against 'boutique' degrees, and especially against 'boutique' courses at the introductory level. Official statistics is also just one of enormous number of possibilities for well-educated statistically-trained graduates, and it is currently difficult enough, in the face of competition for well-trained statistical graduates and university preoccupations with high end research, to prepare graduates for practicing statistician careers. NSO staff are also recruited from diverse degrees to provide strength in breadth in teams, with their subsequent careers emerging from individual capabilities, workplace experience and training. And, as discussed in Section 1, NSO's have increasingly strong vested interests in citizen statistical and data literacy, and in sufficient statistical capabilities in their many clients and collaborators in and outside government.

As mentioned above, NSO's are increasingly contributing valuable direct and support assistance in extra-curricular school and tertiary learning, and aim to improve access to official statistics data for educational purposes. There are topics of particular relevance to official statistics which could be strengthened and augmented at the school and tertiary levels. Tables play a major role in official statistics, and recent experiences in work on UNITAR's MOOC *Understanding data and statistics better – for more effective SDG decision making*, demonstrated how much is involved in the designing and using tables involving more than two variables and different possible representations. The design, data sources, evaluation and interpretation of indices are key to much official statistics, more so than ever in light of the 2030 Agenda, and aspects of these could be gradually introduced in engaging ways throughout schooling and into tertiary courses.

Another area which receives much attention from the critiquing viewpoint in statistics is sampling, but there is a need for constructive hands-on learning experiences in designing and using well-designed sampling. Michael Bulmer's *Islands* [33] provide an ideal setting for student experimentation in sampling, with the *Islands in Schools Project* https://sites.google.com/site/islandsinschoolsprojectwebsite/home providing ideas, guides, classroom activities, activity marking rubrics, exemplar responses and example datasets. Associated with this, just one of the statistical leftovers which are past their use-by-date is the dual meaning of the word 'population'; it is past time to return 'population' to its true meaning. Another concept which could easily and usefully be included in education with simple and engaging contexts and data, is utility, to help combat misunderstandings in risk between probability and outcome, and help build understanding of much in citizen life, from insurance to disease to climate change.

One of the best contributions official statistics can make to improving foundation, introductory and more specialised education in statistics and data science, is to use influence at the senior level to help within the 'big tent' of statistics and data science, in the advocacy, promotion and emphasis on the extent and importance of the many needs outlined in this article. Collaboration and collaborative leadership across all areas of statistics and data science have never been of greater value and significance.

## References

[1]    Rumsey, D.J. Statistical literacy as a goal for introductory statistics courses. Journal of Statistics Education. 2002; 10(3): ww2.amstat.org/publications/jse/v10n3/rumsey2.html.

[2]  Gal, I. (ed.) Adult Numeracy Development: Theory, Research, Practice, Cresskill, NJ: Hampton Press 2002.

[3]  Watson, J.M. Assessing Statistical Thinking Using the Media, The Assessment Challenge in Statistics Education, IOS Press and The International Statistical Institute, Gal I and Garfield JB (ed), Amsterdam. 1997; 107-121. ISBN 90-5199-333-1

[4]  Budgett, S., Pfannkuch, M. Assessing students' statistical literacy, ISI/IASE Satellite, Guimaraes, Portugal. 2007. www.stat.auckland.ac.nz/∼iase/publications/sat07/Budgett_Pfannkuch.pdf.

[5]  Porter, A. Statistical literacy for law students: six hours to teach! ICOTS 5, IASE, Singapore. 1998. THE ROYAL STATISTICAL SOCIETY'S (iase-web.org).

[6]  MacGillivray, H. You know what I mean. Teaching Statistics. 2017; 39(2): 39-41.

[7]  MacGillivray, H. Data science, statistical investigations, team sport and assessment. Teaching Statistics. 2019; 41(1): 1-2.

[8]  Gould, R. Data literacy is statistical literacy. Statistics Education Research Journal. 2017; 16: 22-2. http://iase-web.org/Publications.php?p=SERJ.

[9]  Moore, D., Cobb, G. Mathematics, statistics, and teaching. American Mathematical Monthly. 1997; 104: 801-823.

[10]  Forbes, S. Raising statistical awareness. Teaching Statistics. 1996; 18(3): 66-69.

[11]  Connor, D., Davies, N., Holmes, P. CensusAtSchool 2000. Teaching Statistics. 2000; 22: 66-70.

[12]  Forbes, S., Harraway, J., Chipperfield, J., Siu-Ming, T. Raising the Capability of Producers and Users of Official Statistics. In MacGillivray, H.L., Martin, M., Phillips, B. (eds.) Topics from Australian Conferences on Teaching Statistics: OZCOTS 2008–2012, Springer Science+Business Media, LLC, New York, 2014; pp. 246-265.

[13]  Rodriguez, R. Building the big tent for statistics. Journal of the American Statistical Association. 2013; 108: 501, 1-6. doi: 10.1080/01621459.2013.771010.

[14]  Chambers, J. Greater or lesser statistics: a choice for future research. Statistics and Computing. 1993; 3: 182-184.

[15]  Tukey, J. The future of data analysis. The Annals of Mathematical Statistics. 1962; 33: 1-67.

[16]  Donoho, D. 50 years of data science. Journal of Computational and Graphical Statistics. 2017; 26(4): 745-766. doi: 10.1080/10618600.2017.1384734.

[17]  Cleveland, W. Data science: an action plan for expanding the technical areas of the field of statistics. International Statistical Review. 2001; 69: 21-26.

[18]  Wu, C. Future directions of statistical research in China: a historical perspective (PDF). Application of Statistics and Management. 1986; 1: 1-7.

[19]  Gould, R. Towards data scientific thinking. Teaching statistics. Special Issue on Teaching Data Science and Statistics: Foundation and Introductory. 2021; 43(SI1): S11-S22.

[20]  Cameron, M. Training statisticians for a research organisation. Proceedings of the International Statistical Institute 57th Session. 2009. Durban, South Africa. ISI.

[21]  Kenett, R., Thyregod, P. Aspects of statistical consulting not taught by academia. Statistica Neerlandica. 2006; 60: 396-411.

[22]  Wild, C., Pfannkuch, M. Statistical thinking in empirical enquiry (with discussion). International Statistical Review. 1999; 67(3): 223-265.

[23]  MacKay, R., Oldford, W. Stat 231 Course Notes, Fall, 1994, University of Waterloo.

[24]  Tukey, J. We need both exploratory and confirmatory. The American Statistician. 1980; 34: 23-25.

[25]  Holmes, P. Assessing project work by external examiners. In I. Gal & J. Garfield (Eds.) The assessment challenge in statistics education, 1997: pp. 153-164. Amsterdam: IOS Press.

[26]  Burr, W., Chevalier, F., Collins, C., Gibbs, A., Ng, R., Wild, C. Computational skills by stealth in introductory data science teaching teaching statistics. Special Issue on Teaching Data Science and Statistics: Foundation and Introductory. 2021; 43(SI1): S34-S51.

[27]  Bargagliotti, A., Franklin, C., Arnold, P., Gould, Johnson, S., Perez, L., Spangler, D. Pre-K-12 Guidelines for assessment and instruction in statistics education (GAISE) report II. Alexandria, VA: American Statistical Association and Reston, VA: National Council of Teachers of Mathematics. 2020.

[28]  Rubin, A. What to consider when we consider data, teaching statistics. Special Issue on Teaching Data Science and Statistics: Foundation and Introductory. 2021; 43(SI1): S23-S33.

[29]  Kirkup, L., Pizzica, J., Waite, K., Srinivasan, L. Realizing a framework for enhancing the laboratory experiences of non-physics majors: from pilot to large-scale implementation. European Journal of Physics. 2010; 31(5): 1061-1070.

[30]  Forster, M., MacGillivray, H. Student discovery projects in data analysis. In Reading, C. (ed.) The Proceedings IASE/ISI 8th International Conference on Teaching Statistics, Ljubljana: ISI, Voorburg, The Netherlands, 2010. http://icots.net/8/cd/pdfs/invited/ICOTS8_4G2_FORSTER.pdf.

[31]  Gibbons, K., MacGillivray, H. Education for a workplace statistician. In MacGillivray, H.L., Martin, M., Phillips, B. (eds.) Topics from Australian Conferences on Teaching Statistics: OZCOTS 2008–2012, pringer Science+Business Media, LLC, New York, 2014; pp. 267-294.

[32]  MacGillivray, H. Statistics and data science must speak together, teaching statistics. Special Issue on Teaching Data Science and Statistics: Foundation and Introductory. 2021; 43(SI1): S5-S10.

[33]  Bulmer, M., Haladyn, J. Life on an Island: A simulated population to support student projects in statistics. Technology Innovations in Statistics Education. 2011; 5. Retrieved from http://escholarship.org/uc/item/2q0740hv.