

"Statistical thinking will one day be as necessary for efficient citizenship as the ability to read and write." H.G. Wells

Why Statistical Thinking is Transforming Programming Language Research

And why the dream of usable probabilistic programming in the spreadsheet is getting closer...

Andrew D. Gordon, Microsoft Research and University of Edinburgh

Based on joint work at Microsoft Research with Judith Borghouts, Dany Fabian, Matt McCutchen, Simon Peyton Jones, Claudio Russo, Advait Sarkar, Sruti Srinivasa Ragavan, Neil Toronto, Jack Williams, Nick Wilson

> Colloquium on Probabilistic Programming Collège de France, June 29 and 30, 2022

The Dream of Probabilistic Programming in the Spreadsheet



To empower more people to get more value from data via statistical thinking and representing uncertainty

Clara decides whether to buy her sofa



	A	В	C		D		E	F	(
1						€	2,000.00	Income	
2									
3	Rent			€1	,100.00	_			
4	Commute	_		€	85.00				
5	Utilities	€ 50.00	€ 150.00	€	107.05	=ra	nge(B5,C5)		
6	Sofa			€	700.00				
7									
8	8					€	1,992.05	Expenses	
9						€	7.95	Balance	
10	8						1	Above zero	o?
44									

1	A	В	С		D		E	F	Tod
1						€	2,000.00	Income	AJ
2									
3	Rent			€ 1	1,100.00				
4	Commute	1		€	85.00				
5	Utilities	€ 50.00	€ 150.00	€	75.46	=ra	nge(B5,C5)		
6	Sofa			€	700.00				
7									
8						€	1,960.46	Expenses	
9		N	<mark>ew:</mark> cell ho	lds a	а	€	39.54	Balance	
10		prob	probability distr				1	Above zero?	
11									
12					1000	0:	3 fields	=Sample(AboveZero,D12)	
13			lara's form	ula +			78%	=E12.expectation	
14		hor cho	$\frac{1}{2}$ a $\frac{7}{2}$ a $\frac{7}{2}$						
15		st	aying solve	ent					



End-user programming

"Programming to achieve the result of a program primarily **for personal, rather than public** use. [Nardi 1993]

End-user programmers might be secretaries, accountants, children, teachers, interaction designers, scientists, or anyone who finds themselves writing programs to support their work or hobbies." [Ko et al, 2011]

Felienne Hermans Spreadsheets are Code (SANER 2016)



- Ambitious spreadsheets have thousands of formulas
- Sometimes developed over years

ALGORITHMS BY COMPLEXITY								
LEFTPAD	QUICKSORT	GIT	SELF- DRIWNG CAR	GOOGLE. SEARCH BACKEND	SPRAULING EXCEL SPREADSHEET BUILT UP OVER 20 YEARS BY A CHURCH GROUP IN NEBRASKA TO COORDINATE THEIR SCHEDULING			

Sheet defined functions (SDFs)

Peyton Jones, S.L., Blackwell, A.F., Burnett, M.M.: A user-centred approach to functions in Excel. ICFP 165–176 (2003)



Aim of Microsoft's Project Yellow

Remove the "glass ceiling" that limits the scope and reach of what an end-user programmer can do with Excel:

- Make Excel functions reflect the abstractions of our end users, by allowing end-users to define new functions using ordinary formulas.
- Make Excel's data values reflect the datatypes of our end users' domains, by adding arrays, vectors, records, and hence represent domain-specific data types such as probability distributions.





https://xkcd.com/2453/

New function range(a,b) via LAMBDA





The dream of usable probabilistic programming in the spreadsheet

- was there since the 1980s because of RAND(),
- but getting easier because of transforming the spreadsheet programming language with general constructs like compound data, LAMBDA, SDF

But why is statistical thinking transforming programming language research?

Amy J. Ko, "A human view of programming languages" SPLASH 2016



• "...equal access to even the most basic elements of computation requires an epistemological pluralism, accepting the validity of multiple ways of knowing and thinking."

Sherry Turkle and Seymour Papert, "Epistemological pluralism: styles and voices within the computer culture", Signs: Journal of Women in Culture and Society 1990:16(1)

Definition →

- PL is math \rightarrow
- PL is interface \rightarrow
- PL is design →
- PL is notation \rightarrow
 - PL is media →
 - PL is power →
- → certainty
 → efficiency
 → utility
 → sharing
 → expression

Value

agency

Discoveries (my impression)

"What values are behind our research? When I look at my impression of where we spend time, most of our discoveries are about certainty, **but there are a lot of values here that we don't explore all that deeply**."

Amy J. Ko

A Checklist Manifesto for **Empirical Evaluation: A Preemptive Strike Against** a Replication Crisis in **Computer Science**

SIGPLAN 201

by Emery D. Berger, Stephen M. Blackburn, Matthias Hauswirth, Michael W. Hicks

https://www.sigplan.org/Resources/EmpiricalEvaluation/

Clearly Stated Claims

Example Violational

Suitable Comparison

Example Violations

Principled Benchmark Choice

Example Violations

Adequate Data Analysis

Example Violations

:525





"Because user studies are currently relatively infrequent in the papers we examined, we have not included them among the category examples."

https://www.sigplan.org/Resources/EmpiricalEvaluation/

SIGPLAN 207-

A Glorious Dream... but no user study





An example of quantitative statistical thinking in programming languages

Elastic SDFs

1	А	в	с	D	E	F
1	Average	functi	on			
2						
3	Input:	4		Count:	5	=COUNT(B3:B7)
4		5		Sum:	42	=SUM(B3:B7)
5		23		Average:	8.4	=E4/E3
6		4				
7		6				

```
function AVERAGE( B3:B7 ) returns E5 {
E3 = COUNT( B3:B7 )
E4 = SUM( B3:B7 )
E5 = E4/E3
```

Problem:

- =AVERAGE(X5:X7) too small!
- =AVERAGE(G2:G200) too big!

User's eye view

- \cdot Write a function with fixed-size inputs, using familiar copy/paste
- · Magic happens
- \cdot The function works on input of arbitrary size

Main point:

we think that automatically inferred elasticity will dramatically broaden the audience that can use SDFs effectively.

Participants

- · 20 participants
- Students

(Computer Science, Statistics, Management, Mathematics)

 10 participants industry experience (teacher, administrator, economist, consultant)

Procedure

- 1. Video tutorial of ~10 minutes
- 2. Elastic SDFs
 - 1 practice task
 - 3 study tasks
 - Cognitive load questionnaire
- 3. Array SDFs
- 4. Semi-structured interview



Findings from user study (N=20, 7 female)

People perceived **significantly lower cognitive workload** for elastic SDFs than with SDFs based on map/reduce.

I think elastic functions are easier to work with, also with the "mental model" that you have of Excel, because you can more just follow your normal Excel **workflow**. – P9

It'd be nice to have this kind of **middle ground**, of not having to write the same things over and over again, but not having to persuade someone to make a macro either. – P7



Reactions to LAMBDA in Online Forums

An example of **qualitative** statistical thinking in programming languages



G

н

F



E

Leila Gharani 📀

1.68M subscribers

Advanced Formula Environment





RECURSIVE

BRAND

FUNCTIONS

A **thematic analysis** of nearly 2,700 comments posted on the Reddit, Hacker News, YouTube, and Microsoft Tech Community online forums, regarding LAMBDA in spreadsheets.

Findings:

- computational abstractions are viewed both as helpful and harmful,
- users encounter learning and understanding barriers to applying them,
- there are deficiencies and opportunities in tooling such as in formula editing, versioning, reuse and sharing.

LAMBDA prompts new debate around:

- whether spreadsheets are code,
- whether writing formulas can be considered programming
- whether spreadsheet users identify themselves as programmers.

Advait Sarkar, Sruti Srinivasa Ragavan, Jack Williams, Andrew D. Gordon, "End-user encounters with lambda abstraction in spreadsheets: Apollo's bow or Achilles' heel?", VL/HCC'22, to appear.

V. Braun and V. Clarke, "Using thematic analysis in psychology," Qualitative research in psychology, 3(2):77–101, 2006

Embrace the human in programming research

Statistical thinking is transforming programming language research

- Quantitative methods assess effectiveness of new features with users
- Qualitative methods find patterns across population of users

Consider assessing your research with **human-centric methods** If you review for PL venues, **embrace the diversity of research disciplines**

Empower people with statistical thinking

The dream of usable probabilistic programming in the spreadsheet

- was there since the 1980s because of RAND(),
- but getting easier because of transforming the spreadsheet programming language with general constructs like compound data, LAMBDA, SDF

Come join us

https://aka.ms/CalcIntel