

# REEEVR - Automated Conversion of Excel Cost-Effectiveness Models to R Models

Presenting Author: Michael J. O'Donnell

Authors:

Michael J. O'Donnell <sup>1</sup>

Qian Xin <sup>1</sup>

Gabriel Rogers <sup>2</sup>

Abdul-Lateef Haji-Ali <sup>3</sup>

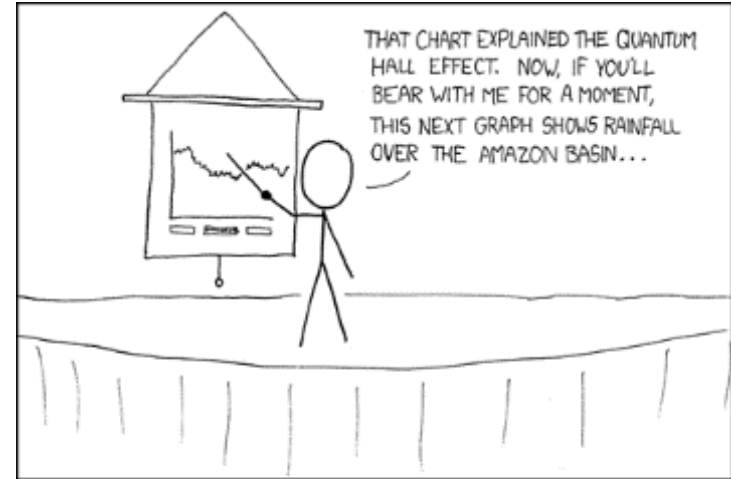
Howard Thom <sup>1</sup>

<sup>1</sup>University of Bristol, UK; <sup>2</sup>University of Manchester, UK; <sup>3</sup>Heriot-Watt University, UK



# The presentation

- Previous work
- Why a conversion process
- The Technical Bit
- Examples
- Demonstration



IF YOU KEEP SAYING "BEAR WITH ME FOR A MOMENT", PEOPLE TAKE A WHILE TO FIGURE OUT THAT YOU'RE JUST SHOWING THEM RANDOM SLIDES.

<https://xkcd.com/365/>

# Previous improvements to HTA

## Excel focused:

- SAVI
  - Generates CE/CEAC automatically
  - EVPI
  - EVPPI
  - Requires PSA already performed

## R focused

- BCEA
  - Encourages use of Bayesian modelling
  - Generates CE/CEAC
  - EVPPI
- MLMC
  - R Package based on C++ package.
  - Utilises MLMC to sample the parameter space for VOI calculations

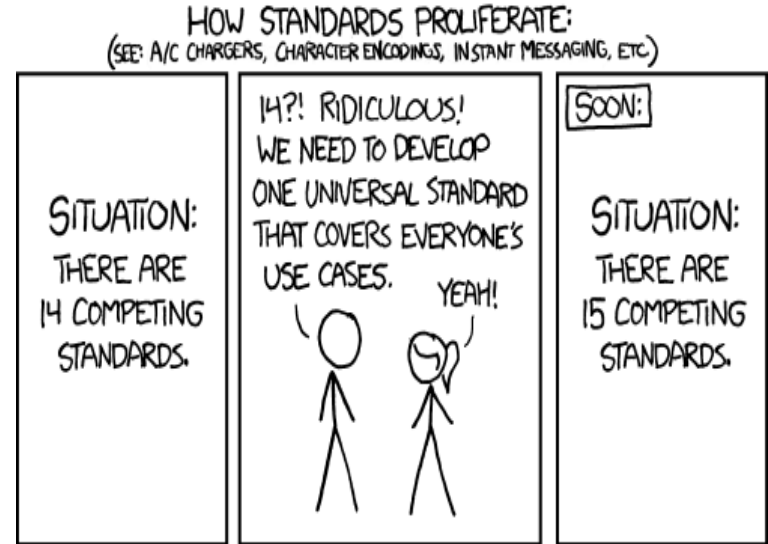
# Advantages of pure R

- Can utilise various VOI methods
  - Types: EVPI, EVPPI
  - Premade R libraries: MLMC, BCEA
- Can take advantage of high performance computing techniques
  - In hardware: multithreading, SIMD, multi-nodal (MPI)
  - In Software: vectorisation, pre-compilation
- Reduces effect of HTA PSA runtime restrictions (e.g. CADTH)
- It's free



# Disadvantages of pure R

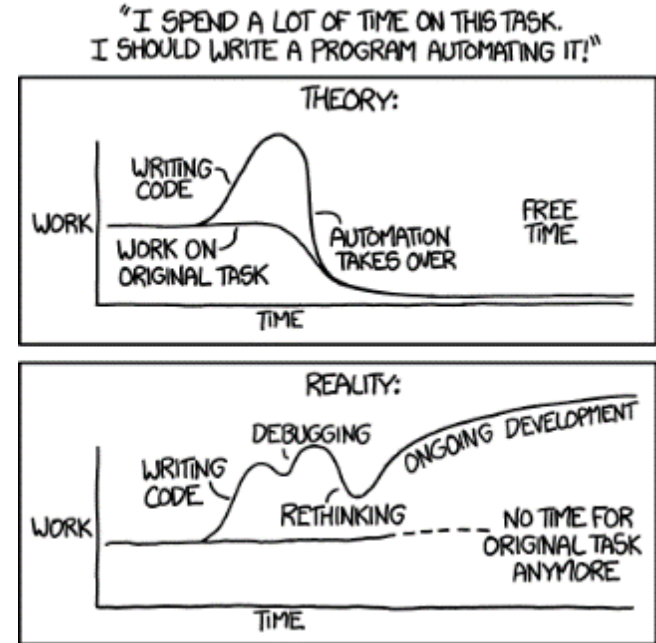
- Has a higher technical bar to entry
  - Need to learn a programming language
  - Generally requires stronger mathematical knowledge to take full advantage
- Is not the current industry standard
  - Many in the industry unfamiliar
  - May not have time to learn
  - Is not accepted by most HTA bodies



<https://xkcd.com/927/>

# Advantages of Automated Conversion

- Allows those unable to use R to still produce models suitable for use with modern techniques
- Allows for old models to be converted easily with minimal loss of productivity
- Prevents the loss of reproducibility if everyone moves away from Excel

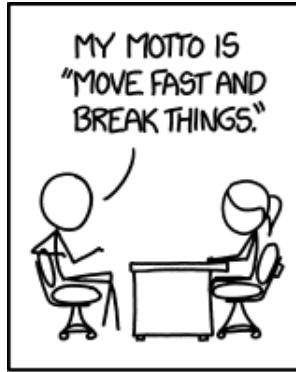


<https://xkcd.com/1319/>

# Advantages of Automated Conversion

- Faster PSA calculations – Allows more complex models, or larger PSA sample sizes.
- Retains usability/familiarity of Excel
- Allows those with less time to learn the language to benefit from R (e.g. clinicians, HTA bodies, etc.)

# The Conversion Process



JOBS I'VE BEEN  
FIRED FROM

FEDEX DRIVER  
CRANE OPERATOR  
SURGEON  
AIR TRAFFIC CONTROLLER  
PHARMACIST  
MUSEUM CURATOR  
WAITER  
DOG WALKER  
OIL TANKER CAPTAIN  
VIOLINIST  
MARS ROVER DRIVER  
MASSAGE THERAPIST



# AI, Chat-GPT and machine learning



<https://xkcd.com/1838/>

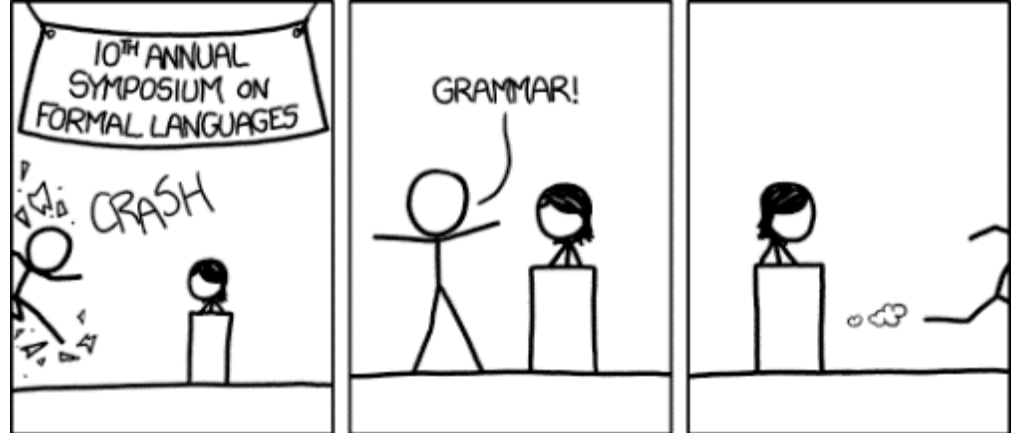
bristol.ac.uk

- Chat-GPT is a large language model, not suitable for complex programming
- Chat-GPT 4 will act confident while giving the wrong answer, so all inputs need checking via an expert
- Other AI style assistants more useful as trained specifically on code, however they face legal challenges: GitHub Copilot
- Machine learning excellent for result analysis and interpretation, but poor at code conversion

# How to convert Excel to R

- Read Excel with openpyxl
- Tokenize and create an Abstract Syntax Tree (AST)
- Traverse the AST and apply a grammar
- Order/sort the code
- Cull dead code
- Add peripherals

<https://xkcd.com/1090/>



# Excel Grammar

## 2.2.2 Formulas

The following ABNF grammar is used by formulas in other parts of this document.

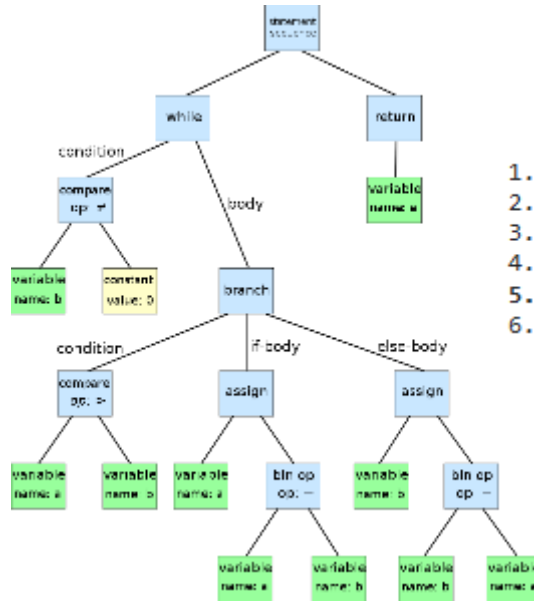
```
formula = expression

expression= ref-expression / "whitespace nospace-expression" whitespace
ref-expression= "whitespace ref-nospace-expression" whitespace
nospace-expression = "(" expression ")" / constant / prefix-operator expression /
expression infix-operator expression / expression postfix-operator / function-call
ref-nospace-expression = "(" ref-expression ")" / ref-constant / ref-expression ref-infix-
operator ref-expression / cell-reference / ref-function-call / name-reference / structure-
reference
constant = error-constant / logical-constant / numerical-constant / string-constant / array-
constant
ref-constant = "#REF!"
error-constant = ref-constant / "#DIV/0!" / "#N/A" / "#NAME?" / "#NULL!" / "#NUM!" /
"#VALUE!" / "#GETTING DATA"
logical-constant = "FALSE" / "TRUE"
numerical-constant = [neg-sign] significand-part [exponent-part]
significand-part = whole-number-part [fractional-part] / fractional-part
whole-number-part = digit-sequence
fractional-part = full-stop digit-sequence
exponent-part = exponent-character [ sign ] digit-sequence
full-stop = "."
sign = "+" / neg-sign
neg-sign = "-"
exponent-character = "E"
digit-sequence = 1*decimal-digit
decimal-digit= %x30-39
nonzero-decimal-digit = %x31-39
string-constant = double-quote [string-chars] double-quote
string-chars = string-char *string-char
string-char = escaped-double-quote / character ; MUST NOT be a double-quote
escaped-double-quote = 2double-quote
double-quote = %x22
;character = as defined by the production Char in the [W3C-XML] section 2.2
array-constant = "{" constant-list-rows "}"
constant-list-rows = constant-list-row *(semicolon constant-list-row)
semicolon = ";"
constant-list-row = constant *(comma constant)
;An array-constant MUST NOT contain an array-constant or columns of unequal length or rows of
unequal length.
operator = ":" / comma / space / "^" / "*" / "/" / "+" / "-" / "&" / "=" / "<>" / "<" / ">"
/ ">" / ">=" / "<=" / "<"
infix-operator = ref-infix-operator / value-infix-operator
value-infix-operator = "&" / "*" / "/" / "+" / "-" / "&" / "=" / "<>" / "<" / ">" / ">"
/ ">"
ref-infix-operator = range-operator / union-operator / intersection-operator
union-operator = comma
intersection-operator = space
range-operator = ":"
postfix-operator = "&"
prefix-operator = "+" / "-"
cell-reference = external-cell-reference / local-cell-reference
local-cell-reference = A1-reference
external-cell-reference = bang-reference / sheet-range-reference / single-sheet-reference
book-prefix = workbook-index "!"
```

- 25 Pages long
- Loosely follows the ABNF grammar system
- Allows anyone to write a “language” that can interpret the grammar
- Comp-sci degree advised

[https://learn.microsoft.com/en-us/openspecs/office\\_file\\_formats/ms-xls/cd03cb5f-ca02-4934-a391-bb674cb8aa06](https://learn.microsoft.com/en-us/openspecs/office_file_formats/ms-xls/cd03cb5f-ca02-4934-a391-bb674cb8aa06)

# Abstract Syntax Tree - General



```
1. while b ≠ 0:
2.     if a > b:
3.         a := a - b
4.     else:
5.         b := b - a
6. return a
```

- Walks the program
- Breaks code into constituent parts
- Builds a tree where each leaf is an elementary component of the program
- The tree is language agnostic

# Abstract Syntax Tree – Excel Tokenizer

= SUM(A1:A10)/SQRT(10)

```
1 {
2   "py/object": "__main__.ProgramNode",
3   "type": "Program",
4   "body": [
5     {
6       "py/object": "__main__.ExpressionNode",
7       "type": "FunctionCall",
8       "name": "SUM(",
9       "params": [
10        {
11          "py/object": "__main__.OperandNode",
12          "type": "Operand",
13          "subtype": "RANGE",
14          "value": "A1:A10"
15        },
16        {
17          "py/object": "__main__.ParseNode",
18          "type": "FunctionClose",
19          "value": ")"
20        }
21      ]
22    },
23    {
24      "py/object": "__main__.ParseNode",
25      "type": "InfixOperator",
26      "value": "/"
27    },
28    {
29      "py/object": "__main__.ExpressionNode",
30      "type": "FunctionCall",
31      "name": "sqrt(",
32      "params": [
33        {
34          "py/object": "__main__.OperandNode",
35          "type": "Operand",
36          "subtype": "NUMBER",
37          "value": "10"
38        },
39        {
40          "py/object": "__main__.ParseNode",
41          "type": "FunctionClose",
42          "value": ")"
43        }
44      ]
45    }
46  ]
47 }
```

- Creates a JSON tree
- Each node contains
  - Type
  - Name
  - Then a combination of
    - Subtype
    - Value
    - Params

# Traverse and Translate

## The broad strokes method

- Pick a cell
- Create AST
- Walk the AST
- Pick R equivalents
- Rebuild codeline

## Caveats

- Cells picked from top left to bottom right.
- Each R equivalent needs to be written
- Need to store important syntactic data (variables, libraries, func type)

# Examples

```
('Inputparameters_F4', ['Inputparameters_F4 = Inputparameters_H4', ['Inputparameters_H4'], 'f'])  
( 'Inputparameters_F5', ['Inputparameters_F5 = Inputparameters_H5', ['Inputparameters_H5'], 'f'])  
( 'Inputparameters_F6', ['Inputparameters_F6 = Inputparameters_H6', ['Inputparameters_H6'], 'f'])  
( 'Inputparameters_F7', ['Inputparameters_F7 = Inputparameters_H7', ['Inputparameters_H7'], 'f'])  
( 'Inputparameters_F8', ['Inputparameters_F8 = Inputparameters_H8', ['Inputparameters_H8'], 'f'])  
( 'Inputparameters_G8', ['Inputparameters_G8 = Inputparameters_F8', ['Inputparameters_F8'], 'f'])  
( 'Inputparameters_F9', ['Inputparameters_F9 = Inputparameters_H9', ['Inputparameters_H9'], 'f'])  
( 'Inputparameters_G9', ['Inputparameters_G9 = Inputparameters_F9', ['Inputparameters_F9'], 'f'])
```

```
('Inputparameters_G16', ['Inputparameters_G16 = qnorm(runif(numberOfRuns)%sep%Inputparameters_H16%sep%Inputparameters_I16)', ['Inputparameters_H16', 'Inputparameters_I16'], 'f'])  
( 'Inputparameters_F17', ['Inputparameters_F17 = Inputparameters_H17', ['Inputparameters_H17'], 'f'])  
( 'Inputparameters_G17', ['Inputparameters_G17 = qnorm(runif(numberOfRuns)%sep%Inputparameters_H17%sep%Inputparameters_I17)', ['Inputparameters_H17', 'Inputparameters_I17'], 'f'])  
( 'Inputparameters_F18', ['Inputparameters_F18 = Inputparameters_H18', ['Inputparameters_H18'], 'f'])  
( 'Inputparameters_G18', ['Inputparameters_G18 = qnorm(runif(numberOfRuns)%sep%Inputparameters_H18%sep%Inputparameters_I18)', ['Inputparameters_H18', 'Inputparameters_I18'], 'f'])
```

# Code Generation and Code Cull

- Need to reorder code
- Remove code that should not exist
- Explicitly write code in the correct format
- Handle slippages (None vs NA)

```
def generate_code(self):  
    """  
    {'variable' : ["codeified string", ['list','of','contained','vars'],cell.data_type]}  
    """  
    with open(self.codefile,"w") as f:  
        f.write(f"{self.mandatoryCode}")  
  
        for item in self.culledcode.items():  
            if item[1][2] != "f":  
                f.write(f"{item[0]} = {item[1][0]}\n")  
            else:  
                item[1][0] = item[1][0].replace("%sep%", " ,")  
                f.write(f"{item[1][0]}\n")  
  
        f.write(f"{self.outputs.add_output_code()}")
```



# REEEVR - R and External Integration

- Creation of a REEEVR library
  - Contains converted functions to mimic Excel behaviour
  - Contains functions to allow validation
  - Contains functions to allow integration with other packages
- External integrations
  - Currently hooks into BCEA to generate CE/CEACs
  - Potential to hook into any R package if desired
  - Qian Xin currently investigating MLMC integration

# Decision Tree - Excel

<b>Inputs</b>	
Willingness-to-pay	20000
Analysis type	Probabilistic
Include MI?	Yes
<b>Outputs</b>	
ICER Warfarin vs aspirin	-1414.250708
INB Warfarin vs aspirin	-31115.37592

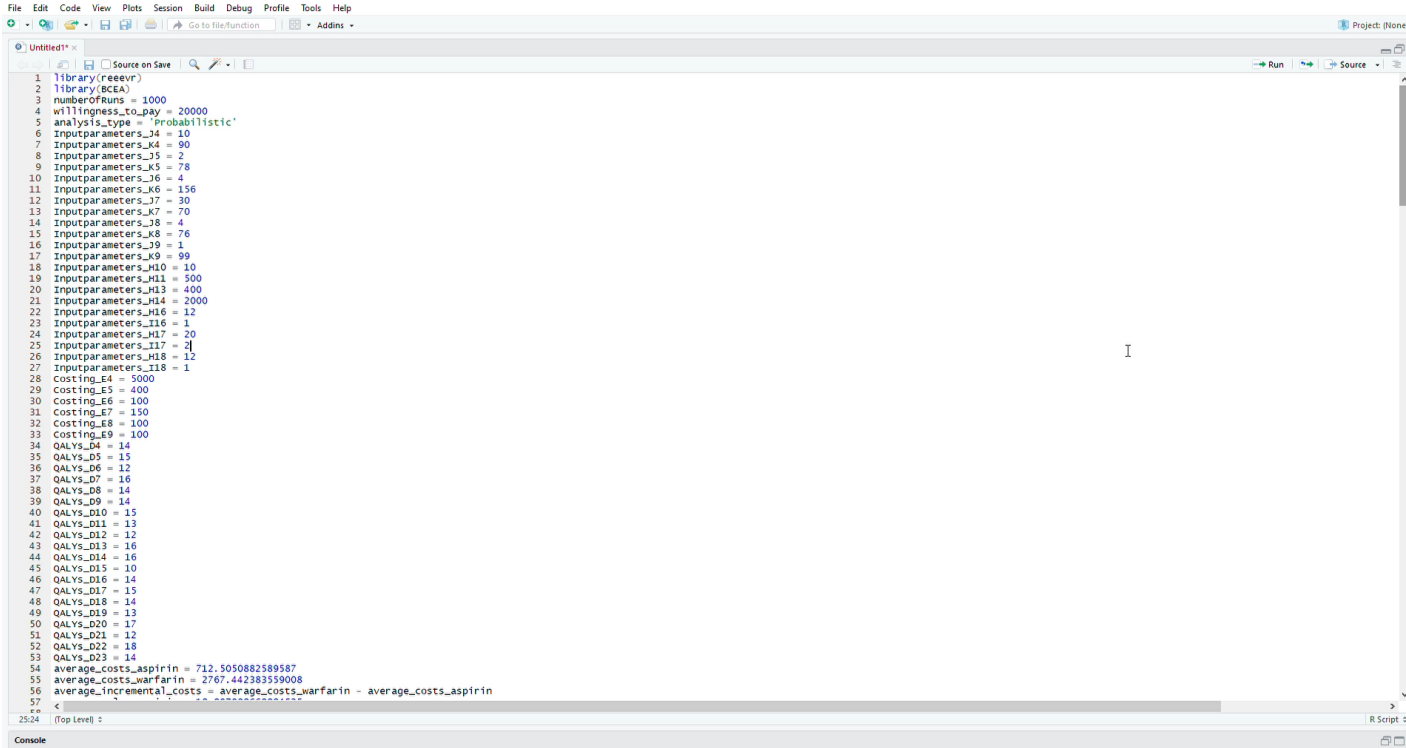
Notes

Base case includes MI but a sensitivity analysis is to set MI probabilities to zero

Run PSA

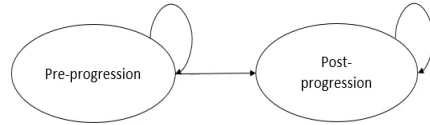
- Minimum viable model
- Contains only PSA
- Doesn't generate graphs
- Contains VBA

# Decision Tree – R (250 lines)



```
1 library(reemv)
2 library(BCEA)
3 numberOfRuns = 1000
4 willingness_to_pay = 20000
5 analysis_type = 'Probabilistic'
6 Inputparameters_j4 = 10
7 Inputparameters_k4 = 90
8 Inputparameters_j5 = 2
9 Inputparameters_k5 = 78
10 Inputparameters_j6 = 4
11 Inputparameters_k6 = 156
12 Inputparameters_j7 = 30
13 Inputparameters_k7 = 70
14 Inputparameters_j8 = 4
15 Inputparameters_k8 = 76
16 Inputparameters_j9 = 1
17 Inputparameters_k9 = 99
18 Inputparameters_j10 = 10
19 Inputparameters_k11 = 500
20 Inputparameters_j13 = 400
21 Inputparameters_k14 = 2000
22 Inputparameters_j16 = 12
23 Inputparameters_k16 = 1
24 Inputparameters_j17 = 20
25 Inputparameters_k17 = 2
26 Inputparameters_j18 = 12
27 Inputparameters_k18 = 1
28 costing_e4 = 5000
29 costing_e5 = 400
30 costing_e6 = 100
31 costing_e7 = 150
32 costing_e8 = 100
33 costing_e9 = 100
34 QALYS_d4 = 14
35 QALYS_d5 = 15
36 QALYS_d6 = 12
37 QALYS_d7 = 16
38 QALYS_d8 = 14
39 QALYS_d9 = 14
40 QALYS_d10 = 15
41 QALYS_d11 = 13
42 QALYS_d12 = 12
43 QALYS_d13 = 16
44 QALYS_d14 = 16
45 QALYS_d15 = 10
46 QALYS_d16 = 14
47 QALYS_d17 = 15
48 QALYS_d18 = 14
49 QALYS_d19 = 13
50 QALYS_d20 = 17
51 QALYS_d21 = 12
52 QALYS_d22 = 18
53 QALYS_d23 = 14
54 average_costs_aspirin = 712.5050882589587
55 average_costs_warfarin = 2767.44238359008
56 average_incremental_costs = average_costs_warfarin - average_costs_aspirin
57 <
58
59 25:24 (Top Level) 2
60
61 Console
```

# 2S Markov Model - Excel



**Assumptions:** The model consists of two health conditions: Pre-progression and Post-progression  
Patients initially in the "Pre-progression" state may either remain in that condition or progress to the "Post-progression" state.  
Patients who starts with or progressed to the "Post-progression" state can only remain in that condition  
Patients who progressed to "Post-progression" will stop the drug treatments  
No patients die within the 10 years time horizon

#### Cell colour codes

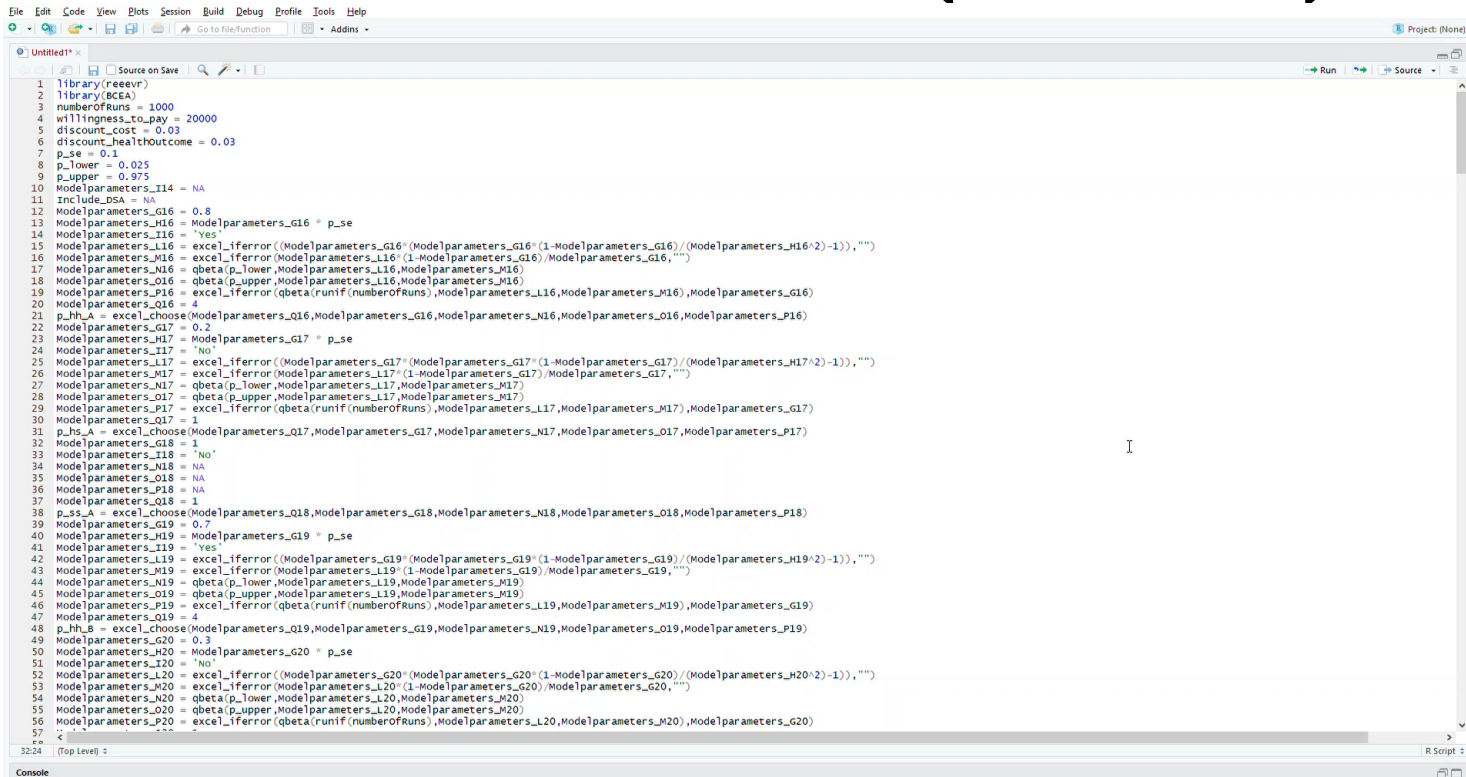
User modifiable input cells  
User modifiable drop-down menu  
Calculation  
Calculation - output by VBA  
Results

#### Tab colour codes

Description and settings  
Input parameters  
Calculations  
Results

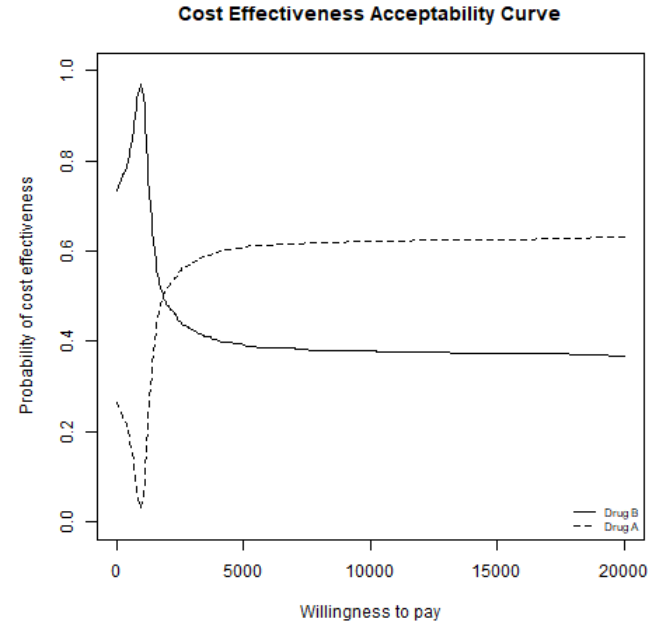
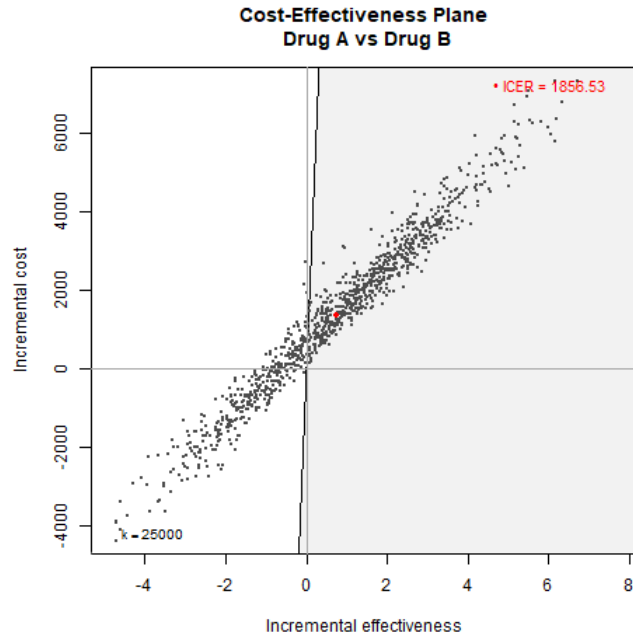
- 2 – State Markov Model
- PSA and DSA
- Excel sheet creates graphs
- Contains VBA

# 2S Markov Model – R – (400 lines)

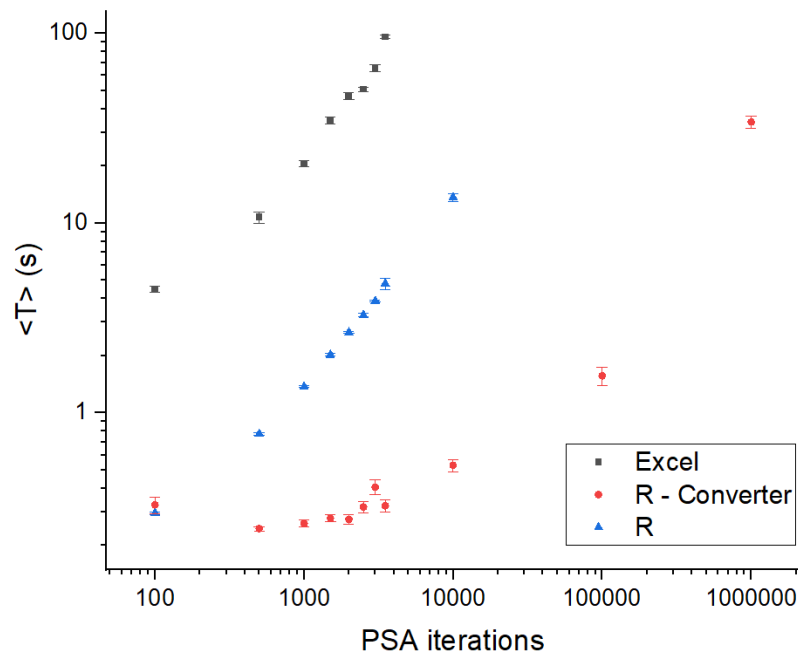
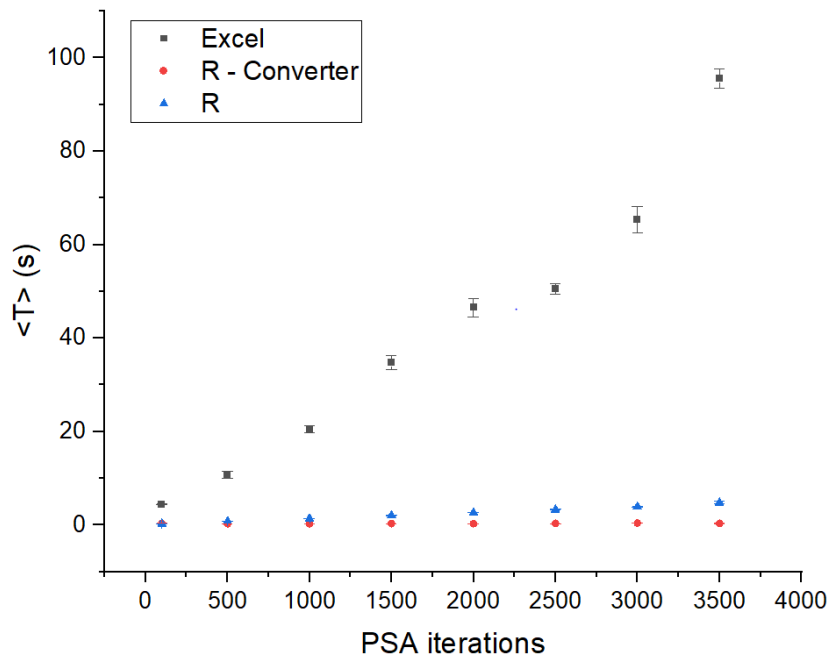


```
1 library(recover)
2 library(BCEA)
3 numberOfRuns = 1000
4 willingness_to_pay = 20000
5 discount_cost = 0.03
6 discount_healthoutcome = 0.03
7 p_se = 0.1
8 p_lower = 0.025
9 p_upper = 0.975
10 Modelparameters_I14 = NA
11 Include_DSA = NA
12 Modelparameters_G16 = 0.8
13 Modelparameters_M16 = Modelparameters_G16 * p_se
14 Modelparameters_I16 = "Yes"
15 Modelparameters_L16 = excel_1fferror((Modelparameters_G16*(1-Modelparameters_G16)/(Modelparameters_M16^2)-1)), "")
16 Modelparameters_M16 = excel_1fferror((Modelparameters_L16*(1-Modelparameters_G16)/(Modelparameters_G16)), "")
17 Modelparameters_M16 = qbeta(p_lower, Modelparameters_L16, Modelparameters_M16)
18 Modelparameters_O16 = qbeta(p_upper, Modelparameters_L16, Modelparameters_M16)
19 Modelparameters_P16 = excel_1fferror(qbeta(runif(numberOfRuns), Modelparameters_L16, Modelparameters_M16), Modelparameters_G16)
20 Modelparameters_O16 = 4
21 p_hh_A = excel_choose(Modelparameters_G16, Modelparameters_M16, Modelparameters_O16, Modelparameters_P16)
22 Modelparameters_G17 = 0.2
23 Modelparameters_M17 = Modelparameters_G17 * p_se
24 Modelparameters_I17 = "No"
25 Modelparameters_L17 = excel_1fferror((Modelparameters_G17*(1-Modelparameters_G17)/(Modelparameters_M17^2)-1)), "")
26 Modelparameters_M17 = excel_1fferror((Modelparameters_L17*(1-Modelparameters_G17)/(Modelparameters_G17)), "")
27 Modelparameters_M17 = qbeta(p_lower, Modelparameters_L17, Modelparameters_M17)
28 Modelparameters_O17 = qbeta(p_upper, Modelparameters_L17, Modelparameters_M17)
29 Modelparameters_P17 = excel_1fferror(qbeta(runif(numberOfRuns), Modelparameters_L17, Modelparameters_M17), Modelparameters_G17)
30 Modelparameters_O17 = 1
31 p_hs_A = excel_choose(Modelparameters_G17, Modelparameters_M17, Modelparameters_O17, Modelparameters_P17)
32 Modelparameters_G18 = 1
33 Modelparameters_I18 = "No"
34 Modelparameters_M18 = NA
35 Modelparameters_O18 = NA
36 Modelparameters_P18 = NA
37 Modelparameters_Q18 = 1
38 p_ss_A = excel_choose(Modelparameters_Q18, Modelparameters_G18, Modelparameters_M18, Modelparameters_O18, Modelparameters_P18)
39 Modelparameters_G19 = 0.7
40 Modelparameters_M19 = Modelparameters_G19 * p_se
41 Modelparameters_I19 = "Yes"
42 Modelparameters_L19 = excel_1fferror((Modelparameters_G19*(1-Modelparameters_G19)/(Modelparameters_M19^2)-1)), "")
43 Modelparameters_M19 = excel_1fferror((Modelparameters_L19*(1-Modelparameters_G19)/(Modelparameters_G19)), "")
44 Modelparameters_M19 = qbeta(p_lower, Modelparameters_L19, Modelparameters_M19)
45 Modelparameters_O19 = qbeta(p_upper, Modelparameters_L19, Modelparameters_M19)
46 Modelparameters_P19 = excel_1fferror(qbeta(runif(numberOfRuns), Modelparameters_L19, Modelparameters_M19), Modelparameters_G19)
47 Modelparameters_Q19 = 4
48 p_hh_B = excel_choose(Modelparameters_Q19, Modelparameters_G19, Modelparameters_M19, Modelparameters_O19, Modelparameters_P19)
49 Modelparameters_G20 = 0.3
50 Modelparameters_M20 = Modelparameters_G20 * p_se
51 Modelparameters_I20 = "No"
52 Modelparameters_L20 = excel_1fferror((Modelparameters_G20*(1-Modelparameters_G20)/(Modelparameters_M20^2)-1)), "")
53 Modelparameters_M20 = excel_1fferror((Modelparameters_L20*(1-Modelparameters_G20)/(Modelparameters_G20)), "")
54 Modelparameters_M20 = qbeta(p_lower, Modelparameters_L20, Modelparameters_M20)
55 Modelparameters_O20 = qbeta(p_upper, Modelparameters_L20, Modelparameters_M20)
56 Modelparameters_P20 = excel_1fferror(qbeta(runif(numberOfRuns), Modelparameters_L20, Modelparameters_M20), Modelparameters_G20)
57 <
58 <
```

# 2S Markov Model - Automated BCEA output



# Speed comparison for HIPS model



# Demonstration – HIPS model (2600 lines)



# Validation – show hips output

- Current state of Validation
  - Validation code automatically added to R output
  - Works by propagating known values through R to give known costs/utilities
  - Requires user to add equivalent validation to Excel model (change inputs to known default values)
- To be implemented
  - Validation via mean/median values to simplify Excel model changes
  - Suggestions welcome!

# Limitations and Caveats

- DSA may not be implemented
- If DSA implemented, likely to be median DSA, not mean DSA
- The full suite of Excel functions is not currently implemented
- The GUI requires additional usability work
- The R output, while “readable” is difficult to follow



# Future Work

- Unit and Integration testing
  - Currently utilise basic black box testing
  - Unit testing needed for Excel conversions
- Openpyxl updates to tokenizer (3.2?)
- Addition of missing Excel functions
- Testing on wider range of Excel models
  - Prior work guides development
- Integration of MLMC and BCEA more completely via optional commands

<https://xkcd.com/1421/>

