

# S-MAP Workflows on and off the Cluster

Robbie Price

Pascal Omondiagbe

Presenting an overview of technical work to support scientists and data managers undertake repeatable spatial soils research

# OutLine

 Robbie and Pascal, Staff Diagram

#### **SMAP** quick overview

General overview

#### **Digital Soil Mapping**

- Method Overview
- DSM Process
- S-map Ecosystem
- Covariate Process
- DSM modelling Process
- Segmentation Process
- Post processing/Upload

#### Workflow Development

- Requirement
- Retrospective
- Workflow Principle
- Development Process
- DSM Workflow
- Performance/Re-evaluate
- S-Map Workflow Summary
- □ Conclusion

### **SMAP**



# The digital soil map for New Zealand



**EXPLORE MAPS** 







FIRST TIME HERE?



# Soils – multi scale response to geology, landform, climate and time





# Pedologists





# Soil-Landscape Model





Otor\_28 + Kinle\_1.1 + Aroha 3.1

Allophanic hill complex - varying depths of Taupo Pumice and weathered tephras over weathered ignimbrite

Allophanic soils from composite tephras over pumice alluvium

Otor\_19.1

Pumice soils from Taupo Pumice alluvium

Taup\_82.1

Moes 2.1 + Ngah 9.1

Allophanic soils from post 60ka composite tephras Otor 28.1



Report penerated: 11-Feb-2014	for hts lands lands we	meanth courts	
The information shared describes	The best of success success	tax of the specified and is a dapth of 1 matrix, and should	id not be the
pimary source of data when ma	Ring land use decisions on	individual farms and publicity.	
firau (Otorohanga_2	8.1)	Fan	ily: Otorchang
Key physical properties			
Depth class integritelity)		Deep (> 1 m)	
Testure profile		StyLoam	
Putertial rooting depth		Untireland	
Rooting barrier		No significant barrier within 1 m.	
Topsol storiness		Stanatass	
Topsoil day range		12-10%	
Drainage siless		Well drained	
Available in root cone		Unlimited	
Permeability profile		Rapid Over Moderate	
Depth to slowly perseable to	arizon.	No stowly permeable horizon	
Permusbility of showest horiz	-	Muderate (4 - 72 mm/h)	
Profile available water	(0-10001 or tool carrier)		
	- d- ton or not same)	Very high (157 mm)	
	d-Xoron a noi sane	Very high (80 mm)	
Dry balk density, topsoil		0.76 (glom3)	
Dry bulk density, substall		0.86 (glow2)	
Depth to hard rock		No hard rock within 1 m	
Depth to soft rock		No soft rock within 1 m	
Depth to stony layer class		No significant story layer within 1 m	
Key chemical properties			
Toppost Production		High (82%)	
and the second se		1.00000000	
About this publication			
<ul> <li>The Homaton and Becides</li> <li>For Setter internation on Indust</li> </ul>	the typical always properties of last entry, control Landsone Resea	the specified with	
<ul> <li>Advise should be assign from an</li> </ul>	Ant lant use experts before ma	Erg besters or individual fame and particular.	
- The internation sheet is not	rael by Landure Resonant	or an 'as it' and 'as making have and whost any wa	
· Locare Present eat no		makes whose setting registrion and express solution	
designed biosciences' and when the	er sauset to a year of this Solition	e	Enderson
the state of the s			
b   I and rate Barray			Waike

# ... Soil Map

- •Art
- Not data driven
- Not suitable for modern demands
- •Regional planning max
- Definitely not farm scale
- •S-Map plenty not like



# ... And so digital soil mapping was born



# DSM (Digital Soil Mapping) – developed by Pedologists and slotted into the normal manual way of doing things...

# S-Map & Soil Mapping Ecosystem





## **Covariate Process**

LiDAR elevation model, Climate data, Geology layers

□ Covariates derived using existing tools

- ArcGIS
- GRASS
- SAGA
- R
- C/C++
- Python

□ New covariate algorithms

• R – e.g Landscape Pattern (neighbourhood elevation range)

## **Workflow Progress**

□ Wrap existing non ESRI models in R covariate library

- Standardised way of writing
- Definitive Implementations as reusable functions
- · Metadata generation embedded into the workflow

□ Assess alternative Implementations of ESRI models

# **DSM Modelling Process**

□ Expert feel for Area being mapped (Pedologists)

- Site Data to describe soils and locate geospatially
- □ Sites + Covariates + Model
  - Currently using an over-fitting model: Random Forest

□ Iterative process to develop and test models and covariates

- Choice of covariates
- Field Sampling Strategy

### **Workflow Progress**

- □ Site data "management" in Excel
- R based using Knitr to create evaluation documentation at run time
- □ Workflow to introduce type and play options
- Scope for parallel processing by tiling assessed but not implemented

# **Segmentation Process - Raster to Vector**

□ This area most considered an Art - tells a story

- In Definitive Raster Soil model/s
- □ Out Representative vector (linework) map
- □ Vector map is an inexact representation correct at the attribute level.
- □ Removal of small areas (minimum polygon size)
- □ Smoothing of boundaries

## **Workflow Progress**

- □ Re-implemented ArcGIS concept in GRASS
- □ Workflow runs and works
- □ Segmentation needs complete rewrite for science-art reasons

# **Post Processing Issues**

## Curation

- Line work manually linked back to soils
- Yet to gain an understanding of the process

# Edge Matching

- Joining different datasets along boundaries without obvious artefacts
- Traditionally done manually at end of mapping
- Attempting to put at beginning of process

## **Final Curation**

- Scripts to prepare data for upload
- Burning in data that didn't model well
- ArcGIS based



## Retrospective

□ King said I was daft to build a castle in a swamp...

https://www.youtube.com/watch?v=aNaXdLWt17A

#### □ DSM is Evolving Science

- No previous projects are significantly workflow compatible
- Spatial processing often quite different
- Highly dependent on the ESRI suite for GIS component
- Metadata or documentation for intermediate steps inadequate to replicate work
  - Even where it exists
- Some OS geospatial processing methods not particularly good

□ Pedologists require flexibility into the future

□ First attempt of Segmentation was done using the previous NeSI cluster (New hardware, software, rules... basically new system)

# Someone wake Pascal up...



## **Requirement- The Brief as Given**

□ Turnkey end-to-end workflow for each DSM project

□ Standardised Covariates with national coverage and reproducible

#### The Brief we wrote ourselves

- □ Standardise covariate algorithms
- □ Attempt to document prior DSM work
- □ Create workflow system for current and future wor0k
- Identify non workflowable aspects of SMAP work
  - Are there other approaches?
- Find holes in the existing methodology
- Identify which steps of the workflow requires HPC
- □ Code management integrated into the system and documentation

# **Workflow Principles**

- Single language entry-point
  - □ R most of our analytical code
- •One code base for all platforms
  - Same code to run on local workstation as NeSI
  - □ Tell it where to run
- •R6 Classes for wrapping models
- Generate metadata for each output
  - □ ISO standard XML
- Home project location for all project data (File System)

# **DEVELOPMENT PROCESS**

- □ Coding Standardisation
- I Metadata
- □ R6
- □ Workflow template

# **DSM Workflow –Initial plan**



# **DSM Workflow – Design**





#### **DSM Workflow** – Covariate Process

UWrap existing non ESRI models (e.g. Saga, Grass) in R library

□ Previous Grass and Saga R library exits ?

R6 Classes for wrapping models (easier to use by The DSM guys)

□ Process Locally or on the cluster

□ Provide easy template to DSM guys

#### **DSM Workflow** – Covariate Process (R6 Class)

```
## Topographic Wetness Index #####
sagaAlgorithm <-function(){</pre>
 sm=smapDSM::runSaga()
 sm$path ="C:/Program Files/saga-6.3.0/saga_cmd.exe"
 smSprojectDir ="D:/Projects/smap-dsm/data/"
 sm$rasterFile="test_DEM.tif"
 sm$rasterFolder="D:/Projects/smap-dsm/data/"
 sm$resultPath="result"
 sm$outputVector="final"
 sm$sagaInit()
 sm$metaDataFirstName="pascal"
 sm$metaDataSurName="pascal"
 sm$isCluster=F
 #specify saga algorithm name
 sm$sagaAlgorithm = sm$saga$ta_hydrology$saga_wetness_index
 #get paramters
 sm$sagaParameters()
 sm$parameters$`AREA (Type: Grid (output):default value:NA)`="area"
 sm$parameters$`SLOPE (Type: Grid (output):default value:NA)`="slope"
 sm$parameters$`SUCTION (Type: Floating point:default value:10.000000)`=10.0
 sm$parameters$`AREA_TYPE (Type: Choice:default value:1)`=1
 sm$parameters$`SLOPE_TYPE (Type: Choice:default value:1)`=1
 sm$parameters$`SLOPE_OFF (Type: Floating point:default value:0.100000)`=0.1
 sm$parameters$`SLOPE_WEIGHT (Type: Floating point:default value:1.000000)`=1.0
 sm$process()
```

#### **DSM Workflow** – Covariates Process

```
#wetness
covariate=sagaAlgorithm(),
```

#### #geomorphons covariate2=grassAlgorithm()

)

```
clean(destroy = TRUE)|
make(singlePlan, parallelism = "future", jobs = 2,force = TRUE)
config<- drake_config(singlePlan,verbose = 3)
outdated(config)
vis_drake_graph(config)</pre>
```







# DSM Modelling Using Random Forest (Previous Approach)

```
#-- -DEM derived
lstGrds[04] <- "P:\\Projects\\SL1705_Franklin_SMap\\Data\\GIS\\Grids\\DJP_backup30_3_2017\\F5ranklin DEm\\TerrainAtts\\img
\\eudist.img"
lstGrds[05] <- "P:\\Projects\\SL1705_Franklin_SMap\\Data\\GIS\\Grids\\DJP_backup30_3_2017\\F5ranklin DEm\\TerrainAtts\\img
\\eudist.img"
lstGrds[06] <- "P:\\Projects\\SL1705_Franklin_SMap\\Data\\GIS\\Grids\\DJP_backup30_3_2017\\F5ranklin DEm\\TerrainAtts\\img
\\lelementx.img"
lstGrds[07] <- "P:\\Projects\\SL1705_Franklin_SMap\\Data\\GIS\\Grids\\DJP_backup30_3_2017\\F5ranklin DEm\\TerrainAtts\\img
\\lengthx.img"
lstGrds[08] <- "P:\\Projects\\SL1705_Franklin_SMap\\Data\\GIS\\Grids\\DJP_backup30_3_2017\\F5ranklin DEm\\TerrainAtts\\img
\\lengthx.img"
```

```
#--- run the random forest
```

```
#rfResult <- randomForest(frmFields, data = tblInput[tblTrainingData,],ntree=500,mtry= 5)</pre>
```

```
rfResult <- randomForest(frmFields, data = tblTrainingData, ntree=500,mtry= 5)</pre>
```

#### plot(rfResult)



#### DSM Workflow – Modelling Via SmapDSM Package



)

# **Segmentation Workflow - Initial Design**

Manually Upload segmentation script to cluster
Manually Run segmentation process
(you get the picture) Download Result back **INot Efficient ? ILots of passwords and keys**

# Segmentation Workflow Re-design



#### **DSM Workflow : End-to-End**



### **Performance/Re-evaluate - Workflow**

Ben Roberts (NeSI) provides report on usageEvaluate the performance

□ Redesign segmentation script

- Auto-generate individual tiles
- Auto-generate SLURM script on the fly
- Process each tile job on different node
- R SLURM job to stitch job together

# $\bigcirc$

# S-Map Workflow - Summary

Dream of an end-to-end workflow not fully realised (almost there, (not)

Chose a single language of entry

An over arching code protocol

Designing a flexible coding system for evolving science

Individual components of Workflow being built

Code management crucial



# Acknowledgements

- Linda Lilburne funding and driving
- James Barringer preliminary segmentation, feedback
- Scott Fraser, Sharn Hainsworth DSM work guinea pigs (photos we stole from Scott's presentations)
- Ben Roberts (and other NeSI support)