

# A big data approach to predicting grain crop yield

Presented by:  
Liana Johnson



THE UNIVERSITY OF  
**SYDNEY**

Sydney Institute of Agriculture

## ***Other team members:***

- *Edward Jones*
- *Thomas Bishop*
- *Niranjan Acharige*
- *Sanjeevani Dewage*
- *Patrick Filippi*
- *Sabastine Ugbaje*
- *Thomas Jephcott*
- *Stacey Paterson*
- *Brett Whelan*

# A problem

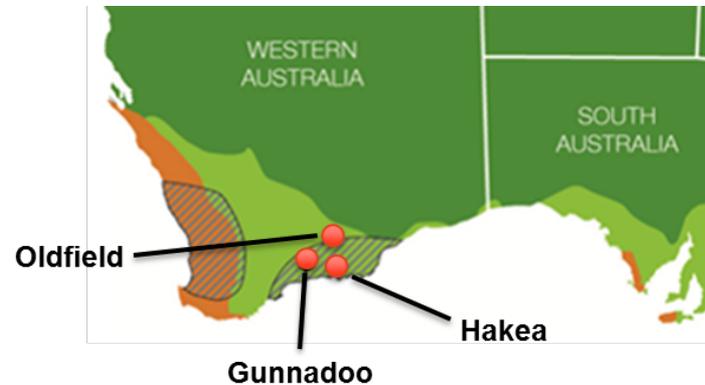
- Farms across Australia have large amounts of unused data
- This data may be difficult to utilise to make management decisions
  - Different formats and located in a variety of repositories
  - Different spatial and temporal resolutions



**How can we transform all of these disparate data streams into something useful, and then inform management decisions?**

# An opportunity

- Hackathon run by CSIRO and Lawson Grains
- Provided us with an abundance of spatial agricultural data (~15,000 ha)
  - Yield
  - EM and gamma surveys
  - Management data
  - Soil tests



- There is also a lot of publicly available spatio-temporal environmental data
  - Rainfall, soil map of Australia, remote sensing etc.

# Available spatial and temporal data

## Provided farmer data:

- Yield – 10 m (space & time)
- Radiometrics – 10 m (space)
- EM surveys – 10 m (space)
- Soil test results (space & time)

High spatial  
resolution  
data

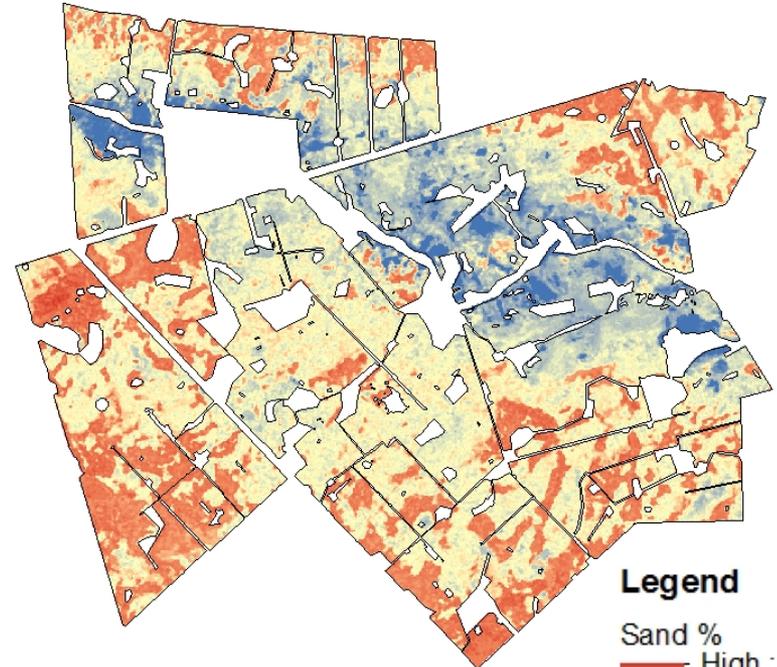
## Project-created data

- Clay and sand soil maps – 10 m (space)

## Publicly-available data:

- TERN – soil maps ~ 90m (space)
- NDVI – 250 m & 16 day (space & time)
- Rainfall forecasts – monthly (time)
- Rainfall received – 5km & daily (space & time)

Hakea Sand % 0-50 cm



Within-season  
measurements

# Approach: Our predictive model

- Using this farmer data and publicly-available datasets, we created a model to predict the yield for these three crops in the production rotation:
  - **Wheat**
  - **Barley**
  - **Canola**
- Modelling method: Machine learning (Random Forest)
  - Data-driven rather than mechanistic
- The idea is to use the data from all fields and years to predict yield within each individual field for a farm

# Modelling for decision support

- 3 different predictive yield models for 3 important time points to inform key management decisions:

## 1. APRIL MODEL

- » To provide suggestions for **variable sowing N rates**

## 2. JULY MODEL

- » To provide suggestions for **variable top-dress N rates**

## 3. SEPTEMBER MODEL

- » To determine **final yield prediction**

- More data becomes available as the season progresses

# Results – Map example

10 m resolution

Hakea - July 2015

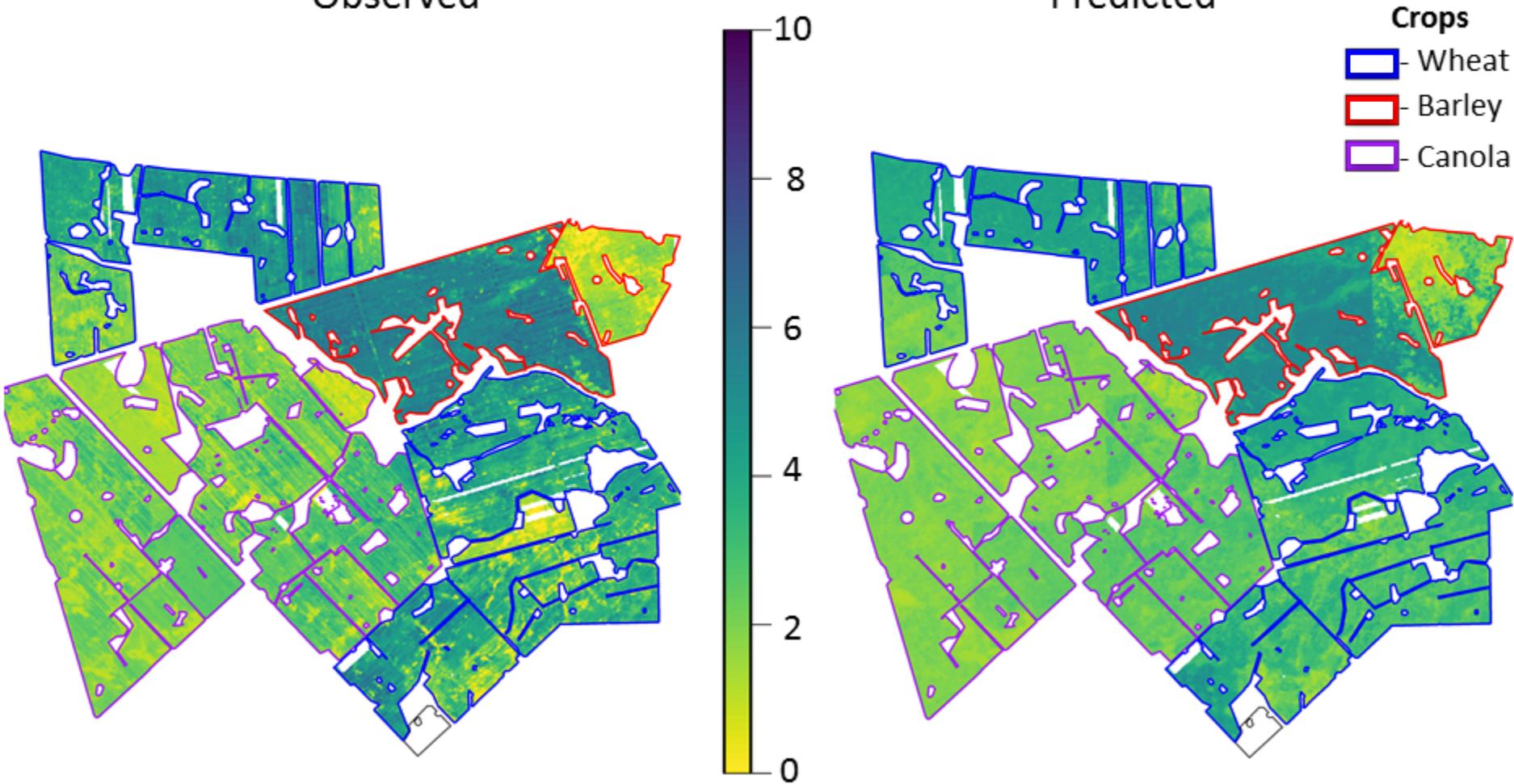
Observed

Yield (t/ha)

Predicted

Crops

- Wheat
- Barley
- Canola

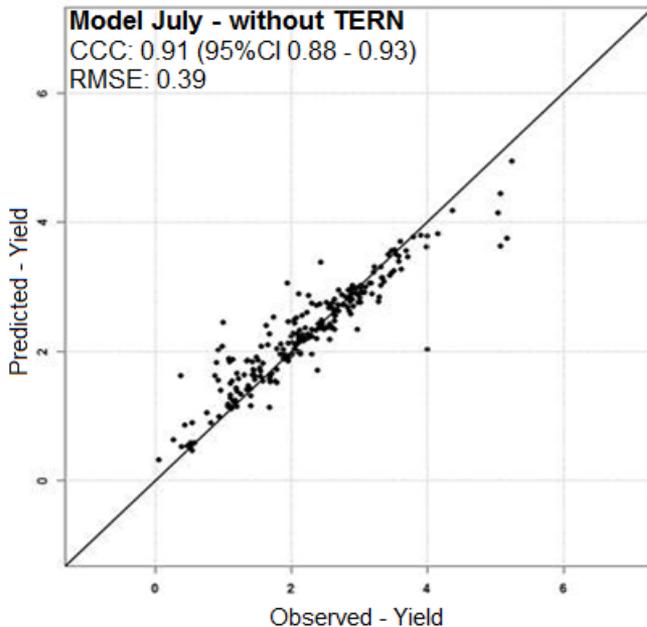


5,500 Ha

# Results – paddock resolution model assessment

1. Predict yield within a paddock, all years of previous yield data excluded
2. Predict yield within a paddock, previous yield data included

TIME	APRIL		JULY		SEPTEMBER	
CV Paddock Predictions	RMSE (t/ha)	LCCC	RMSE (t/ha)	LCCC	RMSE (t/ha)	LCCC
1) Without previous yield	0.64	0.19	0.63	0.20	0.62	0.27
2) With previous yield	0.42	<b>0.89</b>	0.39	<b>0.91</b>	0.36	<b>0.92</b>



**BEST MODELS**

- > LCCC of 1 characterises a perfect fit
- > Including previous data from the prediction paddock results in a better prediction
- > Models are very good

**We have a model that predicts yield, but how can we make this useful and user-friendly for growers and consultants to inform management decisions?**

# Answer: Our user interface

## USYD AgData Challenge



Select a paddock

Aggregate:

Farm:

Paddock:

Property:

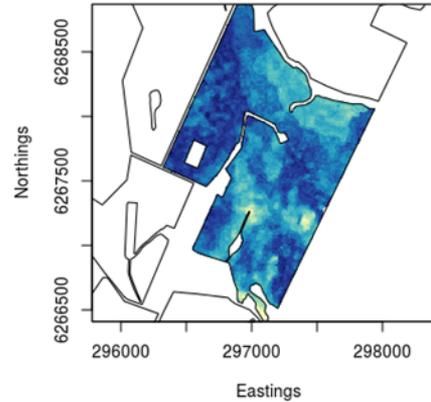
Select and run a model

Model:

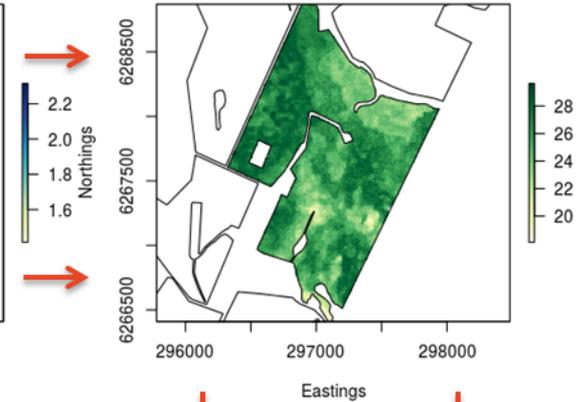
Protein target (%):

Commodity price (\$/tonne):

Predicted barley yield



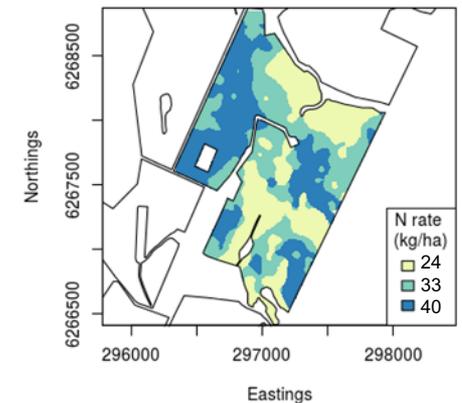
N required (kg/ha)



Comparison

Aggregate -	Oldfield
Farm -	Ranga
Paddock -	Doyle_Flats
Paddock area -	185 (ha)
<b>Model predictions:</b>	
Ave. yield -	2 (t/ha)
Protein target -	8 (%)
Mean N req. -	25.9 (kg/ha)
Total N required -	4788 (kg)

N zones (kg/ha)



# Conclusions

- We used large amounts of agricultural and environmental data to:
  - accurately predict wheat, barley and canola yield across a collection of farms
  - developed a user-friendly application for farmers to aid key management decisions

## What next?

- More data, more accurate predictions- model will improve over time
  - potential to integrate fine spatial resolution remote sensing (drones, Landsat, Sentinel etc.)
- With more consistent data collection it will be useful for:
  - Identifying yield gaps
  - Sowing seeding rate
  - Variable rate application – Lime, P, K & S
  - Futures contracts and market speculation

} For any cropping system