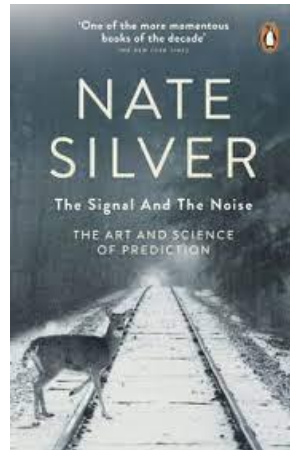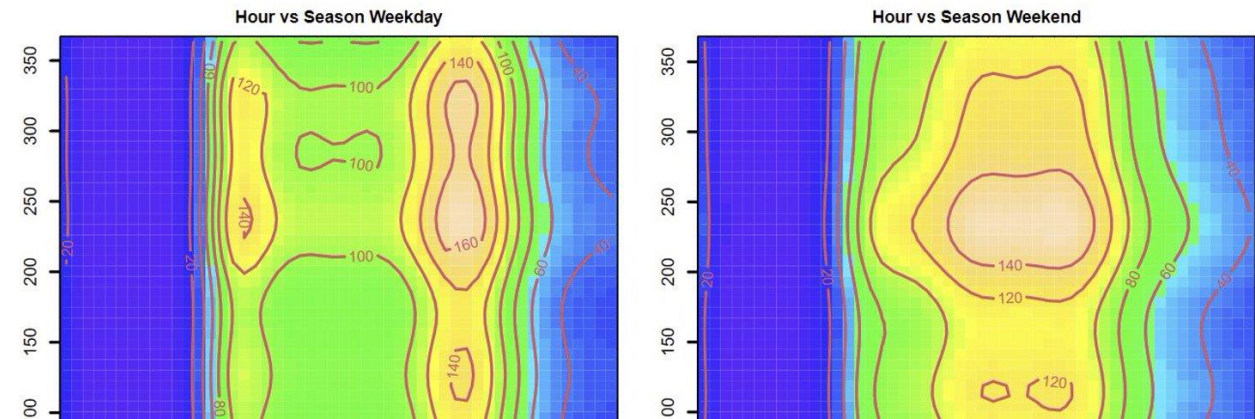# IEEE-CIS Predict-Optimize Technical Challenge

## **Richard Bean**

Centre for Energy Data Innovation

School of Information Technology and Electrical Engineering

**University of Queensland**

Australia

Newcastle Institute for Energy and Resources

**University of Newcastle**

Australia

(most of the time physically based in Newcastle, NSW)

- Ph.D. mathematics (UQ 2001, combinatorics)
- ROAM Consulting (now EY) 2007-2012
- AEMO (Australian Energy Market Operator) 2013
- Redback Technologies (2016-2019)
- University of Queensland (2019-2022)
  - Centre for Energy Data Innovation https://cedi.uqcloud.net/
- Australian and NZ Electricity Market regional and sub-regional demand at ROAM/AEMO – "macro" forecasting
- Individual buildings/solar/distribution transformers at Redback/UQ from inverter or smart meter data - "micro" forecasting
- Cybersecurity – localization of houses with ERA5 solar / load data
- ROAM – simple quadratic programming for modelling NEM bidding (COIN-OR)
- Battery/inverter scheduling at Redback – linear programming
- Combinatorics / graph theory – 0-1 integer programming (CPLEX, BonsaiG, COIN-OR, Gurobi)
- Bike sharing forecasting with GAMs and ERA5 data emph. explainability >> error rate ~ energy
- Classical cryptanalysis – pattern recognition (closely connected)

- The most important step! Reproducible code
- Find the approach that gives the lowest MASE for each time series on phase 1
- R script change PHASE value to 2 and rerun

```
rm(list=ls())

PHASE <- 1

FLIST <- c("phase_1_data.tsf","phase_2_data.tsf")

PDAY <- c(31,30)
PMONTH <- c(10,11)

DAYS <- PDAY[PHASE]
PERIODS <- DAYS * 24 * 4
HOURS <- DAYS * 24
HOUR1 <- HOURS - 1
FIRSTPERIOD <- paste("2020-",PMONTH[PHASE],"-01 00:00:00",sep="'
```

Data Replication & Reproducibility

PERSPECTIVE

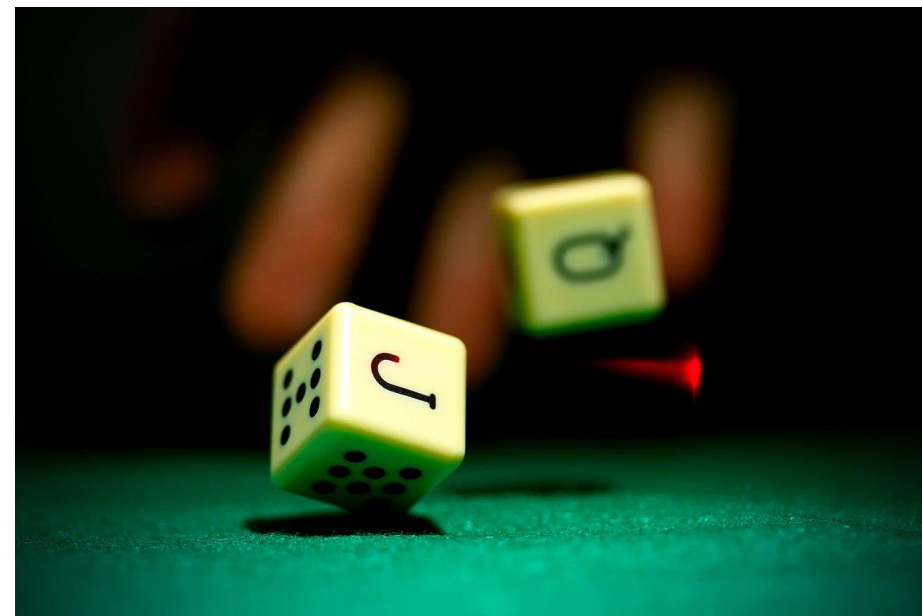**Reproducible Research in Computational Science**

Roger D. Peng



**Fig. 1.** The spectrum of reproducibility.

If you think the competition is just pure skill you won't enter Phase 2 but if you think luck is involved you'll definitely just run your Phase 1 model on Phase 2. i.e. it's better for the competitors and competition organizers if they believe luck is involved.

Pedro Domingos @pmddomingos · Sep 30
Considering that **random forests** have many layers and beat **deep learning** in most applications, maybe we just need to rebrand them as deep forests and they'll be the next big thing.
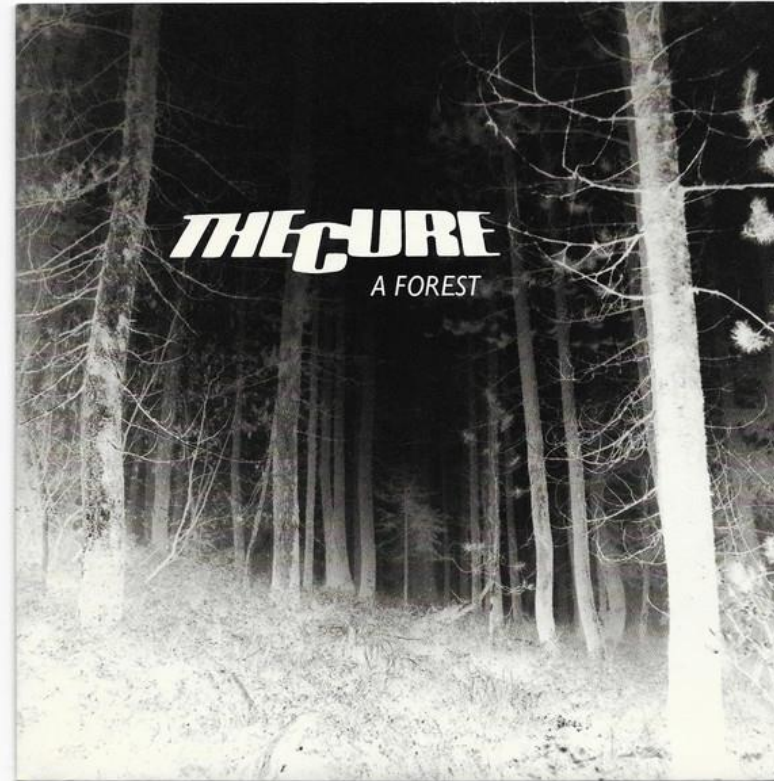
35    95    731

And the forests will echo with laughter. Does anybody remember forests?

Led Zeppelin - Stairway To Heaven - Seattle 07-17-1977 Part 18
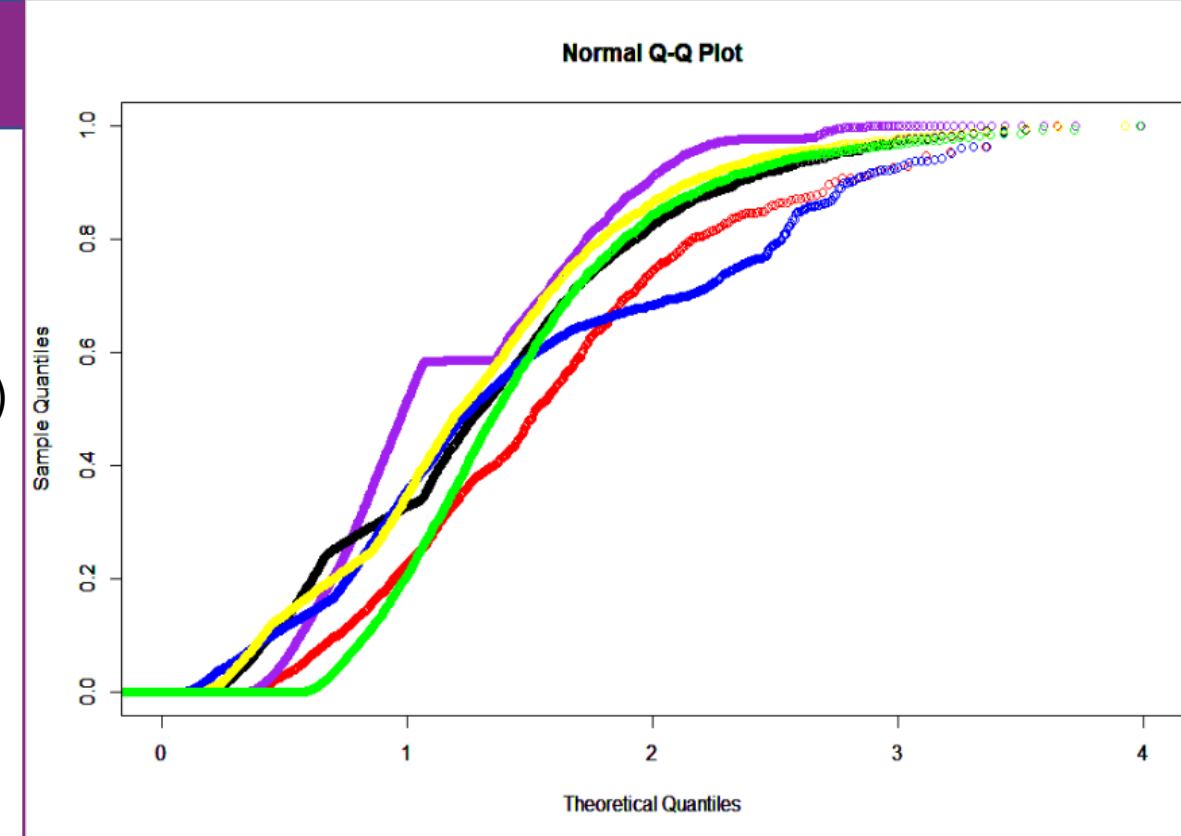5,304,014 views · 19 Dec 2014    22K    1.7K    SHARE    SAVE

Led Zeppelin

THE CURE
A FOREST

The Cure

Deep inside the forest is a door into another land

Thom Pace

5

- Quantile regression forest – forecast **median** to minimise MAE

- (i.e. sum of deviations from actual value)

- Most important parameter to tune – *"mtry"*

- Training against individual phase 1 time series (without overfitting)

- Each hour gets 4 random forests (each quarter hour)

- Choosing building start months of 2020 (Building 0,1,3,6)

- Removing building outliers

- Choosing solar start months (Solar1 has some cumulative data)

- Predictor variables: ECMWF vars lead/lag 3h, day of week, day of year etc

- Public holiday – 23 October Grand Final holiday excluded from training

- Building4 and Building5 set to median values of Oct 2020 (1 and 19 kW)

- Forecast groups of buildings and solar together with normalization (critical, but mentioned by organizers *"cross-learning across time series"*)

- Using BOM daily and ECMWF 1 hour data together (critical … is this surprising?)

- Solar0 and Solar5 thresholding hugely improves MASE (critical)



Normal Q-Q Plot

Almost all my forecast MASE improvement came after Phase 1 data was released.
Obviously lots of room to improve Solar0/5 still

*"Progress usually comes from many small improvements; a change of 1% can be a reason to break out the champagne."*
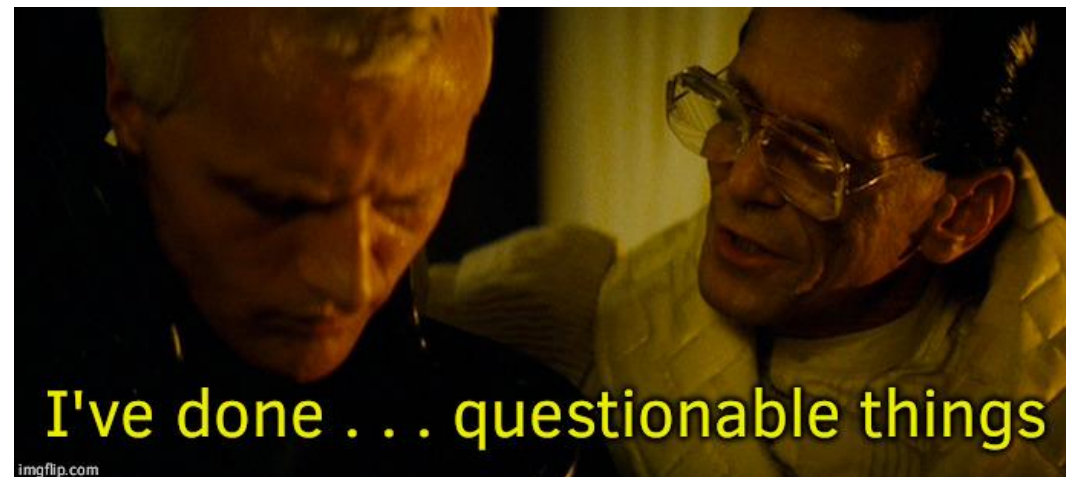
| Case | MASE Phase 1 | MASE after tuning |
|------|--------------|-------------------|
| Building0 | 0.4301 | 0.3859 |
| Building1 | 0.6115 | 0.4251 |
| Building3 | 0.3310 | 0.2913 |
| Building4 | 0.5637 | 0.5637 |
| Building5 | 1.0370 | 0.8383 |
| Building6 | 0.7676 | 0.7336 |
| Solar0 | 0.8479 | 0.6558 |
| Solar1 | 0.4619 | 0.3619 |
| Solar2 | 0.5251 | 0.4139 |
| Solar3 | 0.5910 | 0.4990 |
| Solar4 | 0.5624 | 0.4219 |
| Solar5 | 0.8559 | 0.6092 |
| Mean | 0.6320 | 0.5166 |

Questioning how provided ERA-5 data was derived.

Inverse distance weighting (exponent 2) of four ERA-5 points (0.25 degrees).

Lots of subtleties e.g. exponent choice in IDW, losing wind speed dir/quant nuances.

BOM three points – 8.3, 3.2, 16.1 km away
ERA5 four pts – 21.1, 15.2, 20.7, 14.5 km
ERA5-Land three points – 11.5, 2.9, 10.3 km
MERRA-2 four pts - 15.2, 46.9, 44.6, 63.1 km
JRA-55 one pt – 440 metres

- Could Phase 1 forecast have been improved with extra data (NWP, AEMO etc) or a different approach? (using AEMO data might be a bit circular)

  Yes, but not by large amounts

  - AEMO price and demand data (had to download 3 files for competition Phase 1 & 2) is *half hourly* – is microgrid subject to wholesale price? Price/Demand improves B0/B6 forecast!
  - AEMO Rooftop PV Actual data from NemWeb is *half hourly*
  - ERA5 precipitation data – e.g. ILSPF "Instantaneous large-scale surface precipitation fraction"
  - ERA5-Land data is 0.1 degrees – but only 3 points to interpolate from
  - Other solar vars for PVLib: FDIR ~ GHI, SSRDC, CDIR to derive DNI, DHI etc. Diffuse radiation.
  - Wind direction
  - JRA-55 has 3-hourly data grid point 400 m from Monash
  - NASA MERRA-2 1h data - SWGNT ~ SSRD
  - GFS reanalysis data (3-hourly) is painful to process
  - PvOutput.org has many nearby points (5 min data, $15 donation for 1 year access) or Solar Analytics
  - WeatherMan/Solcast approach – derive solar installation parameters from data, resimulate

- Solving the model as a MIP is much easier than solving the MIQP.

- Almost all of the submitted solution depends on first deriving the best MIP solution possible (i.e. minimizing the recurring load or minimizing the recurring + once-off load) and only then solving as an MIQP

- Gurobi 9.1.2 (laptop phase 1, UQ HPC phase 2)

- Various papers about "Predict+Optimize" problem but Phase 1 and leaderboard seem to indicate no close relationship between forecast result and cost. Complex problem, competition issues, limited time

- **Conservative** is just choosing the lowest recurring load and lowest recurring + once off load and evaluating cost using a naive or flat forecast. This was probably the winning approach for cost in Phase 1, as some competitors had winning results with no forecast, or a poor forecast, but seemed pointless to me as the organizers said quality of forecast should contribute to results in phase 2.

- **Forced discharge** forbids any charging in peak hours, and forces at least one of the two batteries to be discharging in every peak period.

- **No forced discharge** forbids any charging in peak hours, but the MIQP solver decides whether to discharge or do nothing in those hours.

- **Liberal** allows charging in peak, but the maximum of recurring + once off + charge effect for each period is limited to the maximum of recurring + once off load over all periods. This is to avoid nasty surprises when the solver thinks that a period has low underlying load and schedules a charge (due to a low price in that period) but then accidentally increases the maximum load over all periods, which can be very costly.

- **Very liberal** allows charging over peak and does not attempt to control the maximum of recurring + once off + charge effect. This would be the best approach if the forecast was perfect.
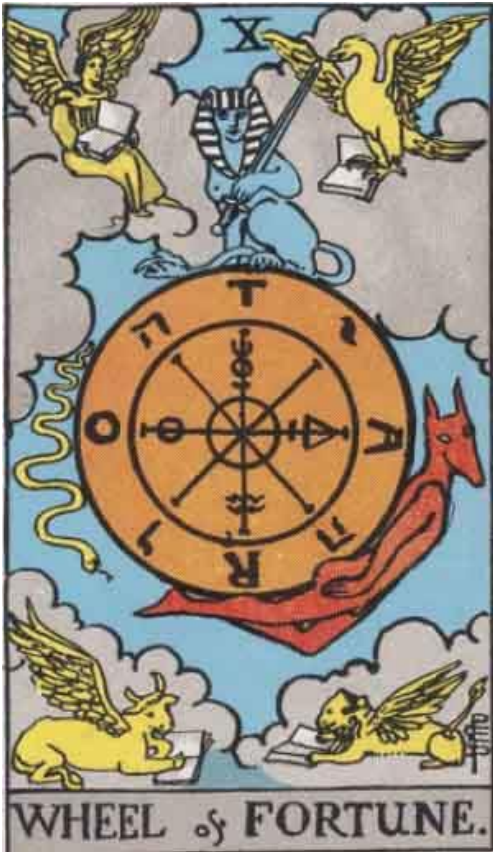
Estimated total cost (3 November) -- $261,906

| Case | Estimated Cost | Actual Cost |
|---|---|---|
| small0 | 26681 | |
| small1 | 26233 | |
| small2 | 26251 | |
| small3 | 26452 | |
| small4 | 26107 | |
| large0 | 26265 | |
| large1 | 26666 | |
| large2 | 25389 | |
| large3 | 26010 | |
| large4 | 25849 | |
| Total | 261906 | |

Only Large2/Large4 had the once-off load in, all activities, in peak.

The estimated cost is very different from the real cost.



WHEEL of FORTUNE.

- Random forest – 4 models for each hour

- Use daily BOM solar data + ECMWF hourly data + temporal variables

- Train buildings and solar together in groups

- Thresholding two solar series


- Arrays approach with 0-1 Mixed Integer Program (MIP)

- First minimize recurring and recurring + once-off load, then solve MIQP

- "No forced discharge" approach chosen from 5 approaches

Contact information

**R.Bean1@uq.edu.au**