

Depth model selection outputs using LMM and wild bootstrap

July 27, 2015

1 Variables

Latitude, Longitude, Elevation, Tmax, Tmin, DelTT, DMI, Nino3.4;

u-wind at 200, 600 and 850;

v-wind at 200, 600, 850;

10 indices of Madden-Julian Oscillations: 20E, 70E, 80E, 100E, 120E, 140E, 160E, 120W, 40W, 10W;

Teleconnections: North Atlantic Oscillation (NAO), East Atlantic (EA), West Pacific (WP), East Pacific-North Pacific (EPNP), Pacific/North American (PNA), East Atlantic/Western Russia (EAWR), Scandinavia (SCA), Tropical/Northern Hemisphere (TNH), Polar/Eurasia (POL);

Solar Flux;

Land-Ocean Temperature Anomaly

Total 35 variables. Medians of precipitation and all these variables are taken across stations and year, and log of precipitation is modeled as LMM with all variables as fixed effects and yearwise random effect.

2 Bootstrap scheme

We are assuming the model $Y = X\beta + Z\gamma + \epsilon$. We do bootstrap on the random effect and residuals separately, i.e. bootstrap samples are calculated as:

$$Y_b = X\beta + Z\gamma_{b_1} + \epsilon_{b_2} = X\beta + ZU_{b_1}\gamma + U_{b_2}\epsilon$$

Where diagonals of U_{b_1} are mean 0 variance τ_k^2 and U_{b_2} are mean 0 variance τ_N^2 : k and N being number of classes and number of samples, respectively. Number of bootstrap samples is 1000.

3 Model selection

The model selection procedure is as given below:

1. Compute depth-based model criterion C_n for the full model;
2. Drop a predictor, compute C_n for the reduced model: repeat for all predictors;
3. Collect the predictors dropping which causes a *decrease* in C_n . Build final model on these predictors.

Criterion values for all drop-1 models as well as full model (`<none>`) are given in the table below and sorted in increasing order. Also (experimental) since C_n is an expected value and we are comparing whether one mean is less than another, I thought of calculating p -values for each drop-1-vs.-full model comparison, by t -test with less than type alternative. They are given in third column of the table.

Scheme I: $\tau_n^2 = n^{0.1}$

	DroppedVar	Cn	pValue
1	- TMAX	0.1311765	0.000000e+00
2	- ELEVATION	0.1567890	0.000000e+00
3	- TempAnomaly	0.1744005	5.720136e-254
4	- del_TT_Deg_Celsius	0.2107309	2.157130e-84

5	- Nino34	0.2407754	2.446939e-08
6	- v_wind_850	0.2409109	2.168531e-08
7	- POL	0.2434999	9.084717e-06
8	- SolarFlux	0.2437457	2.064007e-05
9	- EPNP	0.2440997	2.423362e-05
10	- X120W	0.2443009	7.781675e-06
11	- LONGITUDE	0.2453031	2.852942e-04
12	- u_wind_850	0.2461120	9.066987e-04
13	- TNH	0.2471282	3.579881e-03
14	- EA	0.2477579	8.112529e-03
15	- u_wind_600	0.2479662	1.029812e-02
16	- LATITUDE	0.2503148	1.000400e-01
17	- u_wind_200	0.2506321	1.297290e-01
18	- NAO	0.2519853	2.910135e-01
19	- DMI	0.2520626	2.973663e-01
20	- X20E	0.2523445	3.369725e-01
21	<none>	0.2532943	1.000000e+00
22	- v_wind_200	0.2537721	5.800021e-01
23	- EAWR	0.2545490	7.012957e-01
24	- v_wind_600	0.2554020	8.148190e-01
25	- WP	0.2557333	8.499471e-01
26	- TMIN	0.2558768	8.606282e-01
27	- X160E	0.2559359	8.735409e-01
28	- PNA	0.2569771	9.399411e-01
29	- X140E	0.2580453	9.789425e-01
30	- X120E	0.2615192	9.997137e-01
31	- SCA	0.2616763	9.997404e-01
32	- X40W	0.2623092	9.999299e-01
33	- X70E	0.2630892	9.999868e-01
34	- X100E	0.2648561	9.999996e-01
35	- X10W	0.2655411	9.999997e-01
36	- X80E	0.2732275	1.000000e+00

Scheme II: $\tau_n^2 = n^{0.2}$

	DroppedVar	Cn	pValue
1	- TMAX	0.1594159	1.924099e-318
2	- ELEVATION	0.1831007	5.087412e-209
3	- TempAnomaly	0.1950830	1.610033e-149
4	- del_TT_Deg_Celsius	0.2272400	2.230674e-32
5	- v_wind_850	0.2445276	5.883643e-05
6	- LONGITUDE	0.2461506	7.222295e-04
7	- Nino34	0.2463974	1.399145e-03
8	- SolarFlux	0.2474990	6.114477e-03
9	- u_wind_850	0.2482293	1.402247e-02
10	- POL	0.2488754	2.575858e-02
11	- u_wind_600	0.2494258	4.762298e-02
12	- EPNP	0.2495537	4.843258e-02
13	- EAWR	0.2512024	1.698848e-01
14	- LATITUDE	0.2514350	2.025272e-01
15	- EA	0.2515075	2.070797e-01
16	- TMIN	0.2525158	3.550617e-01
17	- v_wind_200	0.2525735	3.644007e-01
18	- NAO	0.2526962	3.821308e-01

```

19      - v_wind_600 0.2527427 3.920411e-01
20      - X120W 0.2528025 3.941054e-01
21      - TNH 0.2530542 4.405466e-01
22      - u_wind_200 0.2530875 4.486618e-01
23      <none> 0.2533959 1.000000e+00
24      - PNA 0.2537365 5.583152e-01
25      - WP 0.2558643 8.517429e-01
26      - X20E 0.2562980 8.970100e-01
27      - DMI 0.2564950 9.070440e-01
28      - X160E 0.2595774 9.955072e-01
29      - SCA 0.2597392 9.960893e-01
30      - X140E 0.2601329 9.977884e-01
31      - X120E 0.2605941 9.986716e-01
32      - X70E 0.2618301 9.997547e-01
33      - X40W 0.2624374 9.999073e-01
34      - X100E 0.2633385 9.999840e-01
35      - X10W 0.2641683 9.999959e-01
36      - X80E 0.2733044 1.000000e+00

```

4 Discussion

- Top 5-6 variables are as expected;
- EPNP teleconnection and 120W MJO (X120W) are both selected in the model. Both deal with the same longitude region... are they related?
- Interesting variables: Solar Flux and Polar/Eurasia teleconnection (POL): an indicator of Eurasian snow cover;
- Temperature Anomaly has a huge effect. If we do not include this variable, some MJOs are selected, particularly 80E and 40W get selected consistently (Indian ocean and Atlantic oscillations?) across bootstrap schemes. But including temp anomaly makes MJOs almost expendable;
- A thing about p -values: don't really know how much they are useful here, but when I tried the method for linear model selection on a simulated dataset with 25 predictors from Charlie's webpage (<http://www.stat.umn.edu/geyer/5102/exampl/select.html>) all variables with true non-zero coeffs had p -values < 0.05 . Here is its output ($\tau_n^2 = n^{0.2}$):

	DroppedVar	Cn	pValue
1	- x2	0.2051023	1.698015e-42
2	- x3	0.2164894	1.192099e-10
3	- x4	0.2169777	1.211069e-09
4	- x1	0.2209519	1.723179e-04
5	- x5	0.2215058	6.844795e-04
6	- x20	0.2224559	4.012763e-03
7	- x21	0.2252643	1.730813e-01
8	- x9	0.2257273	2.498471e-01
9	<none>	0.2268667	1.000000e+00
10	- x22	0.2279793	7.439175e-01

11	-	x17	0.2285479	8.406974e-01
12	-	x10	0.2285719	8.426074e-01
13	-	x6	0.2290502	8.971758e-01
14	-	x13	0.2292558	9.147554e-01
15	-	x19	0.2293505	9.295542e-01
16	-	x8	0.2294200	9.298908e-01
17	-	x16	0.2299119	9.591960e-01
18	-	x24	0.2304502	9.795701e-01
19	-	x23	0.2304744	9.806469e-01
20	-	x25	0.2305981	9.846541e-01
21	-	x18	0.2309955	9.905178e-01
22	-	x14	0.2311176	9.935712e-01
23	-	x7	0.2314548	9.958976e-01
24	-	x15	0.2322827	9.991981e-01
25	-	x11	0.2324040	9.993293e-01
26	-	x12	0.2327577	9.996462e-01

In the truth, coeffs for x1, x2, x3, x4, x5 are 1, others are 0.

5 Estimation

5.1 Full model

Fixed effect $R^2 = 0.613$, random effect $R^2 = 0.657$.

Summary

```
Linear mixed model fit by REML ['lmerMod']
Formula: log(PRCP + 1) ~ LATITUDE + LONGITUDE + ELEVATION + TMAX + TMIN +
  del_TT_Deg_Celsius + DMI + Nino34 + u_wind_200 + u_wind_600 +
  u_wind_850 + v_wind_200 + v_wind_600 + v_wind_850 + X20E +
  X70E + X80E + X100E + X120E + X140E + X160E + X120W + X40W +
  X10W + NAO + EA + WP + EPNP + PNA + EAWR + SCA + TNH + POL +
  SolarFlux + TempAnomaly + (1 | year)
Data: rainsmall
```

REML criterion at convergence: 3669.2

Scaled residuals:

Min	1Q	Median	3Q	Max
-3.8546	-0.6103	0.0413	0.6686	5.3352

Random effects:

Groups	Name	Variance	Std.Dev.
year	(Intercept)	0.1255	0.3543
Residual		0.9866	0.9933

Number of obs: 1254, groups: year, 35

Fixed effects:

	Estimate	Std. Error	t value
(Intercept)	2.25454	0.06613	34.09
LATITUDE	-0.20304	0.13208	-1.54
LONGITUDE	0.22292	0.06087	3.66
ELEVATION	-1.05131	0.07017	-14.98
TMAX	-1.36674	0.05006	-27.30

TMIN	0.04238	0.08863	0.48
del_TT_Deg_Celsius	0.51886	0.09999	5.19
DMI	0.10713	0.09945	1.08
Nino34	-0.21959	0.13555	-1.62
u_wind_200	0.12262	0.07379	1.66
u_wind_600	0.16603	0.07233	2.30
u_wind_850	-0.26071	0.09036	-2.89
v_wind_200	-0.03265	0.04539	-0.72
v_wind_600	0.03697	0.05347	0.69
v_wind_850	-0.29750	0.05613	-5.30
X20E	-0.41108	0.50330	-0.82
X70E	-0.28442	0.54449	-0.52
X80E	-0.29888	1.02230	-0.29
X100E	-0.21397	0.50682	-0.42
X120E	-0.12062	0.41540	-0.29
X140E	-0.27430	0.39681	-0.69
X160E	-0.38897	0.51338	-0.76
X120W	-0.89048	0.84230	-1.06
X40W	0.14006	0.44374	0.32
X10W	-0.11153	0.46846	-0.24
NAO	-0.08291	0.10330	-0.80
EA	-0.13854	0.12878	-1.08
WP	0.07865	0.10909	0.72
EPNP	-0.18460	0.14379	-1.28
PNA	-0.04457	0.10634	-0.42
EAWR	0.04610	0.09443	0.49
SCA	-0.02448	0.12695	-0.19
TNH	0.15647	0.13868	1.13
POL	0.18843	0.12327	1.53
SolarFlux	-0.14216	0.09213	-1.54
TempAnomaly	0.51812	0.13323	3.89

5.2 Reduced model

Fixed effect $R^2 = 0.606$, random effect $R^2 = 0.649$.

Summary

Linear mixed model fit by REML [`'lmerMod'`]

Formula: `log(PRCP + 1) ~ LATITUDE + LONGITUDE + ELEVATION + TMAX + del_TT_Deg_Celsius + DMI + Nino34 + u_wind_200 + u_wind_600 + u_wind_850 + v_wind_850 + X20E + X120W + NAO + EA + EPNP + TNH + POL + SolarFlux + TempAnomaly + (1 | year)`
 Data: rainsmall

REML criterion at convergence: 3656.4

Scaled residuals:

Min	1Q	Median	3Q	Max
-3.7498	-0.6258	0.0296	0.6806	5.3763

Random effects:

Groups	Name	Variance	Std.Dev.
year	(Intercept)	0.1209	0.3478
Residual		0.9847	0.9923

Number of obs: 1254, groups: year, 35

Fixed effects:

	Estimate	Std. Error	t value
(Intercept)	2.25535	0.06512	34.63
LATITUDE	-0.15626	0.12676	-1.23
LONGITUDE	0.25969	0.05225	4.97

ELEVATION	-1.08153	0.03886	-27.83
TMAX	-1.36575	0.04072	-33.54
del_TT_Deg_Celsius	0.48347	0.09609	5.03
DMI	0.09566	0.06879	1.39
Nino34	-0.28544	0.10602	-2.69
u_wind_200	0.09826	0.07097	1.38
u_wind_600	0.16279	0.07041	2.31
u_wind_850	-0.24310	0.08651	-2.81
v_wind_850	-0.27575	0.04669	-5.91
X20E	-0.01084	0.10114	-0.11
X120W	-0.18115	0.12545	-1.44
NAO	-0.11622	0.07539	-1.54
EA	-0.15751	0.07692	-2.05
EPNP	-0.22894	0.09024	-2.54
TNH	0.23562	0.08035	2.93
POL	0.12736	0.09691	1.31
SolarFlux	-0.09754	0.07652	-1.27
TempAnomaly	0.49256	0.08593	5.73

5.3 Full model vs. reduced model ANOVA

	Df	AIC	BIC	logLik	deviance	Chisq	Chi	Df	Pr(>Chisq)
mod.final	23	3620.9	3739.0	-1787.4	3574.9				
mod.full	38	3627.0	3822.1	-1775.5	3551.0	23.847	15		0.06774 .

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1									

6 Prediction

We now use both the model with all variables and reduced set of variables to evaluate out-of-time prediction performance of resulting models. We use a rolling prediction scheme, in which for each of the years in 2003-2012, we use the previous 25 year's data to build the training model and test it to obtain predictions for yearly median rainfall for all 36 stations at the testing year.

In general predictions from full models are erratic, both in bias and MSE comparisons with the true values. Predictions from reduced models are more stable consistently. There seems to be slight positive bias from the reduced models, which are due to zero median precipitation values, which are always predicted positive.

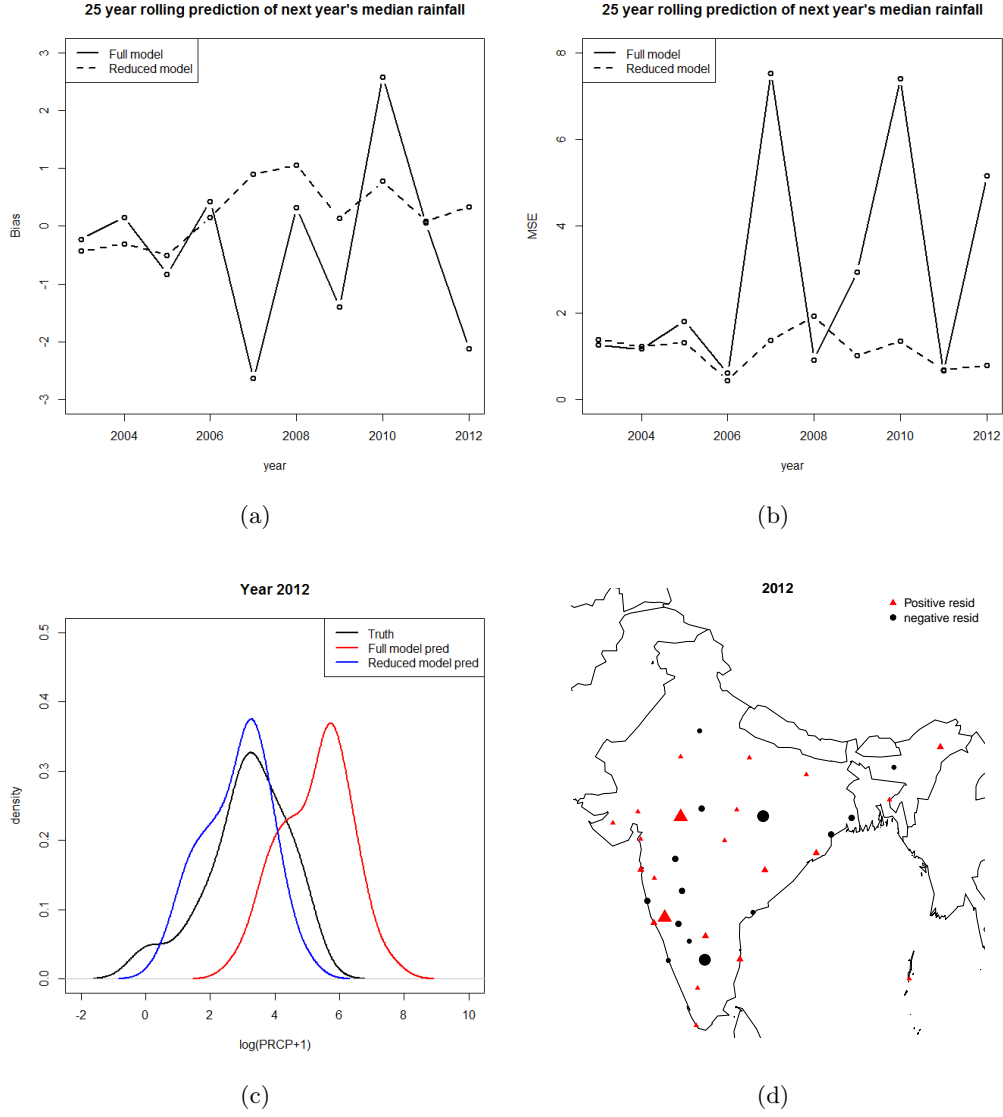
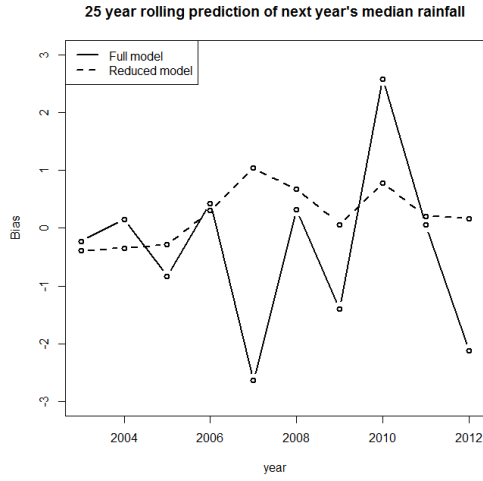
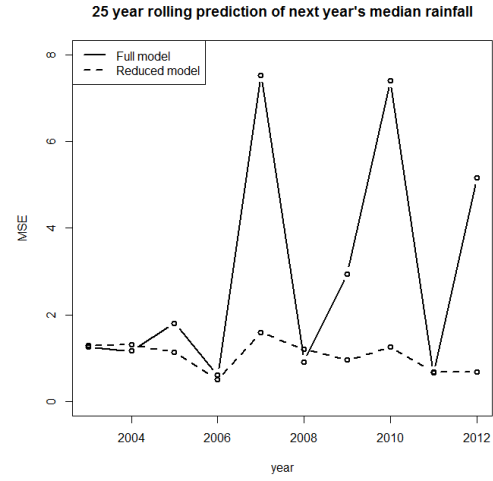


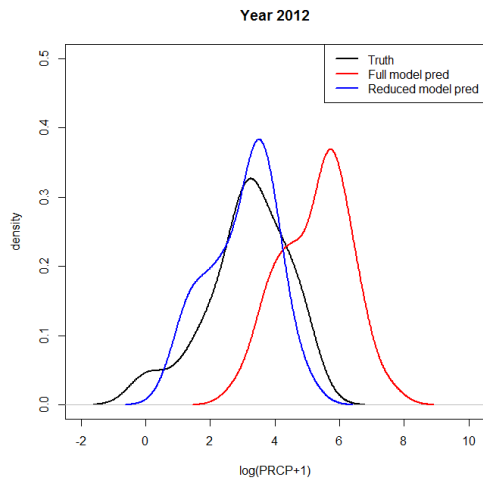
Figure 1: Comparing full model rolling predictions with reduced models: (a) Bias across years, (b) MSE across years, (c) density plots for 2012, (d) stationwise residuals for 2012



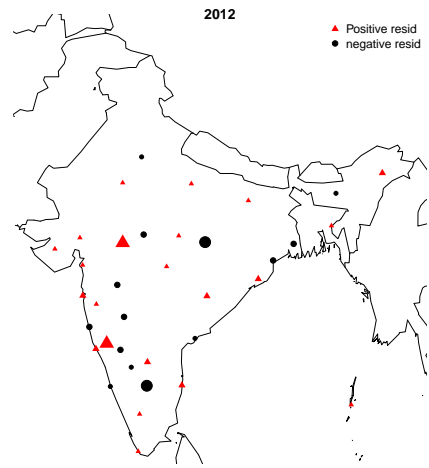
(a)



(b)



(c)



(d)

Figure 2: Comparing full model rolling predictions with p-value reduced models: (a) Bias across years, (b) MSE across years, (c) density plots for 2012, (d) stationwise residuals for 2012