

Report-HW1

(Abhishek Singhal)

Q1)

- a) The criteria for selecting a variable, as node in the tree, is minimizing the mean square error of the two splitting regions. We choose a variable and try all possible splits in the datum. Then we do this process for all variables. We select the variable(j) and split(s) s.t. the mean square error of the region to be divided, is minimum.

For tree of depth=1, it is optimal as we look at all the possible cases.

For trees of depth>1, we do recursive splitting. As naïve splitting is computationally expensive. No, it is not optimal as we have not searched all possible trees.

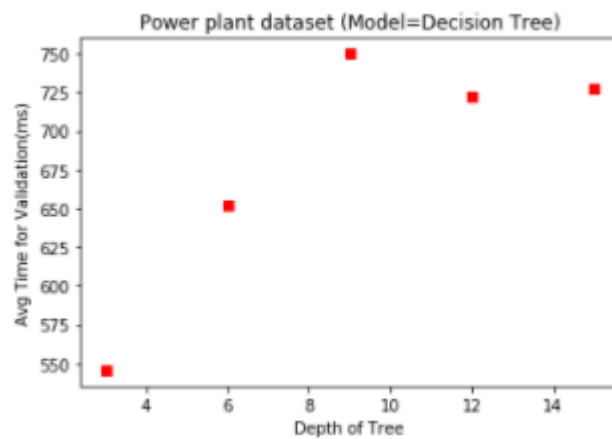
b) **Power Plant**

model chosen: max_depth=9,
criterion='mae'

out of sample error=0.160256966417

full training error=0.118336262681

Out of Sample error is close to Full
training error.



Mo del (de pth)	Error1	Error2	Error3	Error4	Error5	Avg. Error
3	0.21363307694986075	0.21213404111498263	0.21483569790940768	0.21206802787456452	0.2074905393728223	0.21203227664432758
6	0.17165978655988864	0.17703962125435538	0.1741445038327526	0.16561086898954705	0.17427250139372816	0.17254545640605437
9	0.15516345369080786	0.16125527909407672	0.16353107804878048	0.16053877177700346	0.16079624947735191	0.16025696641760409
12	0.16517684087743736	0.15946469616724745	0.16460551567944243	0.16203833275261331	0.15549862473867593	0.16135680204308328
15	0.16459260759052918	0.17032658989547039	0.16705502473867595	0.16769339268292685	0.17000531777003486	0.16793458653552742

c) **Indoor Localization**

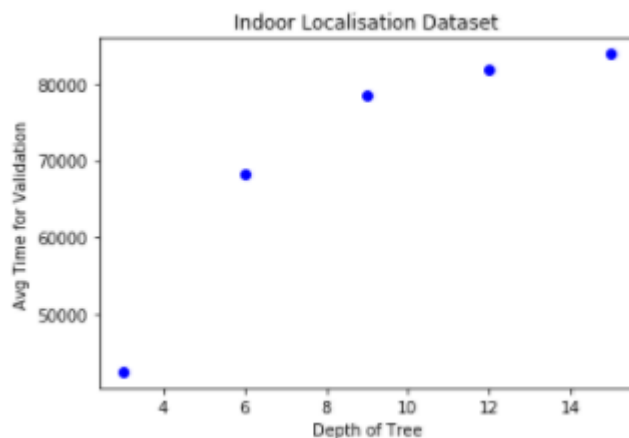
model chosen: max_depth=15

out of sample error=9.0142990917955039

full training error=7.48820330784

Out of Sample error is not very close to Full training error.

Model (depth)	Error1	Error2	Error3	Error4	Error5	Avg. Error
3	40.316933000376132	40.59329762738215	40.015660284926007	40.679415554176074	41.302495941183849	40.58156048160884
6	27.48125589568706	26.50964505692076	28.268997951718084	27.535988266804615	26.338258458490092	27.22682912592412
9	17.423034897881145	18.091903612274322	17.485970172686233	17.789486359794331	18.237538454602461	17.805586699447694
12	12.23399103435306	11.777981722291875	12.882466373651869	12.241183331640332	12.937726133182846	12.414669719023996
15	8.8434292245486468	9.0379293320586758	9.0100572337597189	9.0826730540506642	9.0974066145598194	9.0142990917955039



Q2)

- a) N samples. Each chunk of size =M.
 $k=N/M$ is the number of subsets of training dataset.
We will train (N-M) samples for k times.
Total samples trained = $(N-M) \times (N/M)$
Time taken $\sim N^2/M - N$
Complexity $\sim O(N^2/M)$

If $M=5$, complexity is $O(N^2/5)$.

If $M=N/2$, then running time is linear. $O(2N)$

- b) When we make M small, we actually increase k in k -fold cross validation. With this, the bias of estimated error will be small. Or we can say estimator can be fairly accurate.

Q3.KNN

a) Power Plant

model chosen $n=3$

out of sample error=0.22130739350784454

full training error= 0.361530806

Out of Sample error is close to Full training error.

Model (n_neighbors)	Error1	Error2	Error3	Error4	Error5	Avg. Error
3	0.5437207207520891	0.37095100255427843	0.19186524423285342	8.9421203028509396e-16	8.8916858516244428e-16	0.22130739350784454
5	0.53964015779944297	0.42937341086729358	0.33034079823502088	0.23170993457585784	0.11509728272017836	0.32923231683955867
10	0.58135280278551538	0.48551981532567046	0.4240396581049698	0.37354141827207804	0.31883918442028991	0.43665857578170469
20	0.60848649233983287	0.55284522499129218	0.5021451544356712	0.45413709683852982	0.42235784244704577	0.50799436221047434
25	0.58888568292479115	0.55646554952281435	0.51533930379006032	0.47546261017244379	0.44636510152731329	0.51650364958748463

b) Indoor Localization

model chosen $n=3$

out of sample error=1.60727811

full training error=2.2239872

Out of Sample error is not very close to Full training error.

Model (n_neighbors)	Error1	Error2	Error3	Error4	Error5	Avg. Error
3	3.4594528604563703	2.0966570466399213	1.3172279068795469	0.56479106005224922	0.59826170301449833	1.6072781154085174
5	3.8368866274322975	2.6133712489343037	2.0162536138092459	1.421127889768026	1.0011745483121846	2.1777627856512116
10	4.7261048981569713	3.3950221586822993	2.733493590441364	2.2961816310783743	1.9353148300747409	3.01722342168675
20	5.79761992	4.402627948	3.656813968	3.156505595	2.814353202	3.965584128

	70310937	8622118	1183669	3714785	7662208	4298745
25	6.06134445 60481442	4.814073318 698596	3.982288090 134583	3.468855283 3918527	3.090503718 6056132	4.283412973 375758

Q4. Linear Model

- a) The purpose of penalties in Ridge and Lasso regression is to control capacity. In several cases, all the attributes (p) of input data are not useful. The penalties makes coefficient β_i small for i^{th} attribute of \mathbf{X} . The larger the penalty (α), the smaller the value of β .
In ridge regression, β_s are made small but do not converge to zero. That is, no input attribute is ignored.
In lasso regression, β can be zero. This a selection of attributes occur.

b) Power Plant

Model chosen: alpha=0.01 (Lasso)

Out of sample error=0.18776822

Full training error=0.190764

Out of Sample error is close to Full training error.

Lasso (alpha)	Error1	Error2	Error3	Error4	Error5	Avg. Error
10^{-6}	0.189783 5245755 4071	0.193985380 69880425	0.191860269 21584834	0.190800045 36720838	0.190678764 92417468	0.19142159695 631528
10^{-4}	0.191755 0793772 2855	0.188871314 55626152	0.190059895 67954601	0.191086838 19736783	0.190679416 28063609	0.19049050881 8208
10^{-2}	0.184438 3544036 9816	0.186309279 54994064	0.187578892 74669534	0.189749688 57416435	0.190764887 68205507	0.18776822059 131071
1	0.266660 3571793 485	0.264340012 11203639	0.264054867 66376685	0.262189306 16776909	0.262453287 77656507	0.26393956617 989722
10	0.623296 2057636 6584	0.621489476 33336626	0.623656002 7055742	0.626516512 92788568	0.627689496 95425558	0.62452953893 694951

Ridge(alpha)	Error1	Error2	Error3	Error4	Error5	Avg. Error
10^{-6}	0.191957169 5636817	0.193288345 00600034	0.193177406 71953548	0.191956737 70276014	0.190678758 33731485	0.192211683 46585851
10^{-4}	0.187924906 46272467	0.191037361 19286774	0.190946324 74772946	0.189898398 29986058	0.190678758 34024595	0.190097149 80868567

10⁻²	0.194143618 25633891	0.190600843 53701598	0.190893287 68394551	0.190408498 08191673	0.190678758 63337849	0.191345001 23851911
1	0.186506793 27417996	0.188036059 98555781	0.189331269 88930352	0.191251371 17950047	0.190678787 94657855	0.189160856 45502406
10	0.187529663 66219898	0.189968380 54764123	0.190065329 63670314	0.191539792 73873339	0.190679054 4274605	0.189956444 20254745

c) Indoor Localization

Model chosen: alpha=10 (Ridge)

Out of sample error= 18.970032206799651

Full training error=18.6044577636

Out of Sample error is very close to Full training error.

Lasso (alpha)	Error1	Error2	Error3	Error4	Error5	Avg. Error
10⁻⁶	20.477080 49290879	20.63409129 5324243	20.05826991 8698816	19.87875618 2890246	19.93136286 5996984	20.19591215 1163817
10⁻⁴	20.151478 357417044	19.45512895 7219078	20.46018659 3805226	20.77735627 2540423	20.17055935 2385183	20.20294190 6673392
10⁻²	19.940839 87632904	20.43848729 9950857	20.04116876 6019013	20.19173058 3255904	20.16703379 3553983	20.15585206 3821762
1	21.441923 50345271	21.34193375 0531748	20.77250139 021557	20.99133301 1626896	21.23753469 1471264	21.15704526 9459637
10	35.476969 491427617	35.61380585 5557366	35.32975325 6371198	34.98271055 1830351	35.52238388 2598582	35.38512460 7557024

Ridge (alpha)	Error1	Error2	Error3	Error4	Error5	Avg. Error
10⁻⁶	19.204266 483693473	19.314761045 22444	18.48800472 0819927	19.133594596 324411	18.7533348 95918702	18.9787923483961 9
10⁻⁴	19.006484 640804544	18.824139184 243098	19.03859171 1638318	18.980103339 348965	19.0291399 46428252	18.9756917644926 3
10⁻²	19.039378 468528934	18.920775806 994008	19.19272837 7638382	18.986285827 973809	18.7647284 17721177	18.9807793797712 6
1	18.941915 583117915	18.815454068 655015	19.13698359 5819675	19.079018849 427069	19.0392400 45015809	19.0025224284070 97
10	18.874853 302804265	18.814607148 918189	19.08436171 7063626	18.752438823 573648	19.3239000 41638517	18.9700322067996 51

--	--	--	--	--	--	--