

Homework 1 Report - PM2.5 Prediction

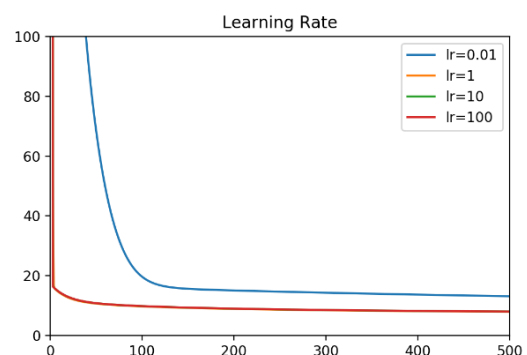
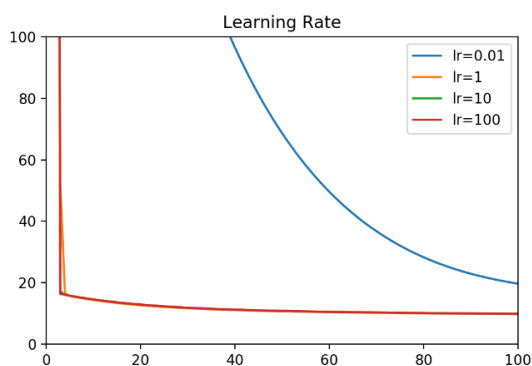
學號：b04505026 系級：工海三 姓名：蔡仲閔

1. (1%) 請分別使用每筆 data9 小時內所有 feature 的一次項（含 bias 項）以及每筆 data9 小時內 PM2.5 的一次項（含 bias 項）進行 training，比較並討論這兩種模型的 root mean-square error（根據 kaggle 上的 public/private score）。

Submission and Description	Private Score	Public Score
output_1.csv a few seconds ago by b04505026_zmtsai add submission details	8.37096	8.39316
output_18.csv 7 minutes ago by b04505026_zmtsai add submission details	7.43071	7.44861

從上圖結果可以看到，使用所有 feature 的 model 得到較佳的結果，而僅使用 PM2.5 的 model 則有不小的落差。和老師上課提到的一樣(寶可夢分類)，這與我們預期的結果相同，若使用較多的 feature，在沒有 noise 的前提下，應該可以 train 出較佳的結果。

2. (2%) 請分別使用至少四種不同數值的 learning rate 進行 training（其他參數需一致），作圖並且討論其收斂過程。



上面兩個分別用了 learning rate=0.01,1,10,100(以下稱 lr)進行作圖，並取不同的 iteration 做觀察，可以看到 lr=0.01 時明顯收斂的比較慢，甚至在 iteration 到 500 次時仍然沒辦法到達其他組的水準，而其他三組的行為則相當接近，即便放到 100 次內也沒有看出太明顯的區別。。

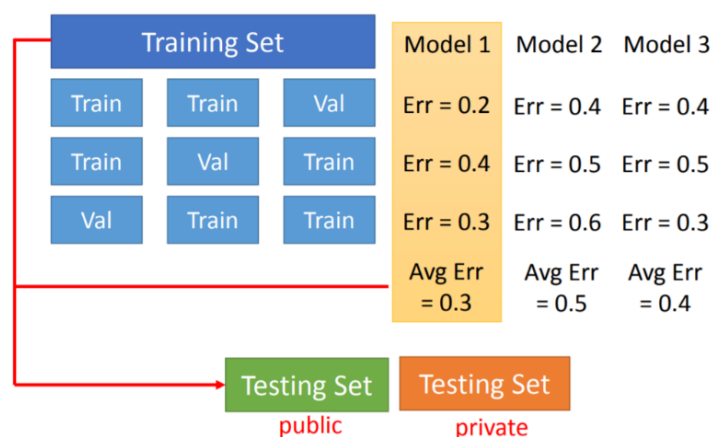
3. (1%) 請分別使用至少四種不同數值的 regularization parameter λ 進行 training（其他參數需一至），討論其 root mean-square error（根據 kaggle 上的 public/private score）。

Submission and Description	Private Score	Public Score
output_ita_2.csv 3 hours ago by b04505026_zmtsai add submission details	12.26638	12.20687
output_ita_1.csv 3 hours ago by b04505026_zmtsai add submission details	11.09248	11.20569
output_ita_05.csv 3 hours ago by b04505026_zmtsai add submission details	10.06456	10.27533
output_ita_0.csv 3 hours ago by b04505026_zmtsai add submission details	7.63631	7.79871

上圖從下到上 lamda 分別從 0, 0.5, 1 到 2，我們可以觀察隨著 lamda 增加在相同的 iteration 中收斂的速度是越來越慢，這也印證 lamda 雖會讓我們得到一個 smooth 的線，但若是 lamda 過大則會造成 under fitting 的結果。

4. (1%) 請這次作業你的 best_hw1.sh 是如何實作的？（e.g. 有無對 Data 做任何 Preprocessing？Features 的選用有無任何考量？訓練相關參數的選用有無任何依據？）在這次作業讓我對 machine learning 有了第一步的認識，從剛開始完全無從著手，到後來看著 Sample code 慢慢推敲並改善，其中讓我順利通過 strong baseline 的方法是資料的挑選，首先我將資料分組，進行 training 並找出較佳的 model。

N-fold Cross Validation



在過程中發現某些月份出現異常的極端值，上網查完資料後，確定那是不應出現的數值，便將資料從 data set 中去除，便順利通過 strong baseline。之後更發現某些月份的資料出現異常的連續的 0，

#	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	
4342	##### 大星	CO	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	
4343	##### 大星	NNBNC	0.13	0.13	0.13	0.1	0.08	0.07	0.06	0.04	0.15	0.16	0.16	0.21	0.2	0.18	0.2	0.2	0.21	0.26	0.26	0.19	0.14	0.11	0.12	0.14			
4344	##### 大星	NO	0.6	0.5	0.4	0.6	0.6	0.8	0.7	1.7	3.5	4.7	5.2	4.9	6	4.7	4.1	3.5	3.4	2.5	2.3	1.5	1	0.8	0.7	1			
4345	##### 大星	NO2	15	15	14	12	10	9	11	17	16	16	15	18	19	18	20	21	24	28	27	22	18	15	15	12			
4346	##### 大星	NOx	16	15	15	12	11	9.8	12	18	20	21	20	23	25	23	24	25	26	31	29	23	19	16	16	13			
4347	##### 大星	O3	15	13	14	15	19	22	21	18	20	22	24	20	20	19	19	20	20	16	16	21	27	27	24				
4348	##### 大星	PM10	9	7	9	11	10	7	8	11	10	4	1	5	15	29	54	74	97	110	122	112	100	79	78	76			
4349	##### 大星	PM2.5	14	18	26	36	12	4	5	1	4	2	3	7	12	26	32	52	59	75	70	60	53	45	41	34			
4350	##### 大星	RAINFAIR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR
4351	##### 大星	SO2	71	72	70	72	72	75	75	71	68	65	65	66	68	67	66	67	67	67	67	67	65	60	62	65			
4352	##### 大星	SO2	1.9	1.8	1.5	1.2	1.1	0.9	0.9	1.3	1.3	1.5	1.5	1.6	2	1.7	1.8	2.1	2.4	2.6	2	1.8	1.9	1.7	1.7	1.6			
4353	##### 大星	TBIC	1.9	2	1.9	1.9	1.9	1.9	1.9	1.9	1.9	1.9	1.9	2	2	2	2	2	2	2	2	2	1.9	1.9	1.9	1.9			
4354	##### 大星	WQI_RR	390	4.1	19	336	303	18	26	384	384	1.2	323	312	320	356	311	301	30	31	18	25	39	30	21	48			
4355	##### 大星	WQI_D	340	48	33	359	318	329	307	311	306	28	317	315	300	319	27	151	234	38	35	339	44	28	49	88			
4356	##### 大星	WQI_QI	2.7	3.2	2.9	3.3	2.9	2.5	2.8	3	3.5	3.9	3.5	4.5	4.7	3.2	3.3	4.1	2.8	3.9	3.4	2.7	3.3	2.5	2.5	1.6			
4357	##### 大星	WQI_RR	1.2	0.5	0.9	1.1	1.1	0.6	0.9	1.3	0.9	1.2	1.2	2.2	1	1.8	1.3	0.7	0.5	0.9	1.1	1	1.3	1.1	1.2	1.5			
4358	##### 大星	AMBI_TB	9.9	9.8	9.6	9.7	9.7	10	10	10	12	13	0	0	0	0	0	0	0	0	11	11	11	10	10	11			
4359	##### 大星	CH4	1.8	1.8	1.8	1.8	1.7	1.7	1.7	1.7	1.7	1.7	0	0	0	0	0	0	0	0	-0.2	-0.2	-0.2	0	0	0			
4360	##### 大星	CO	0.04	0.36	0.33	0.32	0.31	0.33	0.35	0.45	0.51	0.51	0	0	0	0	0	0	0	0.61	0.7	0.76	0.61	0.55	0.6	0.67			
4361	##### 大星	NNBNC	0.1	0.08	0.07	0.07	0.06	0.07	0.08	0.11	0.13	0.14	0	0	0	0	0	0	0	0.03	0.03	0.03	0	0	0	0			
4362	##### 大星	NO	0.6	0.5	0.5	0.8	0.9	0.7	0.9	1.6	3.2	4.6	0	0	0	0	0	0	0	0.7	1.4	2.7	2.3	1.5	1.9	1.8			
4363	##### 大星	NO2	11	10	16	6.7	6.2	6.8	5.8	12	15	15	0	0	0	0	0	0	0	6.9	33	35	28	25	28	31			
4364	##### 大星	NOx	12	11	9.1	7.5	7.1	7.5	11	13	18	20	0	0	0	0	0	0	0	7.6	33	37	30	27	30	33			
4365	##### 大星	O3	26	31	35	35	34	31	27	25	24	20	0	0	0	0	0	0	0	1	2.8	12	36	18	13	9.9			
4366	##### 大星	PM10	71	84	86	47	35	60	63	52	47	44	0	0	0	0	0	0	0	56	18	73	78	76	66	66			
4367	##### 大星	PM2.5	33	40	40	43	31	21	12	24	34	3	0	0	0	0	0	0	0	34	47	47	42	32	25	26			
4368	##### 大星	RAINFAIR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR	NR
4369	##### 大星	SO2	64	60	60	60	60	59	58	58	56	42	0	0	0	0	0	0	0	0	57	59	59	60	61	62			
4370	##### 大星	SO2	1.9	1.8	1.8	1.8	1.8	1.8	1.8	1.8	1.8	1.9	1.9	1.9	1.9	1.9	1.9	1.9	1.9	1.9	1.5	6.3	5.3	5.3	3.4	3.6			
4371	##### 大星	TBIC	1.9	1.8	1.8	1.8	1.8	1.8	1.8	1.8	1.8	1.9	1.9	1.9	1.9	1.9	1.9	1.9	1.9	1.9	-0.1	-0.1	-0.1	0	0	0			
4372	##### 大星	WQI_RR	31	305	39	399	43	26	35	36	10	15	0	0	0	0	0	0	0	327	337	331	336	12	238	92			
4373	##### 大星	WQI_D	26	36	36	36	28	47	356	40	380	227	0.5	0	0	0	0	0	0	338	326	316	33	42	272	86			
4374	##### 大星	WQI_QI	2.7	3	2.7	3	2.4	1.9	3.1	3	2.4	0.1	0	0	0	0	0	0	0	2.2	2.3	2.4	1.8	0.8	0.8	0.7			
4375	##### 大星	WQI_RR	1.1	1.2	1.3	1	1.1	1.1	1.3	1.7	1.2	0.8	0	0	0	0	0	0	0	0.8	0.7	1.1	1	0.6	0.9	0.3			
4376	##### 大星	AMBI_TB	10	10	9.5	9.2	9.5	9.3	9.2	9.5	12	16	19	22	27	24	23	22	20	18	18	17	16	15	15	15			
4377	##### 大星	CH4	1.8	1.8	1.8	1.8	1.7	1.7	1.7	1.7	1.7	1.7	0	0	0	0	0	0	0	0	-0.2	-0.2	-0.2	0	0	0			
4378	##### 大星	CO	0.6	0.55	0.5	0.48	0.55	0.61	0.62	0.7	0.81	0.8	0.6	0.48	0.43	0.45	0.48	0.57	0.68	0.86	0.96	0.96	0.96	0.96	0.91	0.81			
4379	##### 大星	NNBNC	0.1	0.08	0.07	0.07	0.06	0.07	0.08	0.11	0.13	0.14	0	0	0	0	0	0	0	0.03	0.03	0.03	0	0	0	0			
4380	##### 大星	NO	1.9	1.5	1.4	0.9	0.8	1.5	4.6	10	14	11	6.8	3.3	1.7	1	1	1.2	1	0.7	1.6	8.1	9.8	11	8.9	11			
4381	##### 大星	NO2	31	27	35	23	25	26	26	25	24	27	26	18	15	13	15	19	21	30	41	47	43	38	37	36			
4382	##### 大星	NOx	33	28	26	28	26	28	30	33	31	37	33	21	17	14	16	21	28	31	42	55	53	82	46	46			

此時再將這些資料去除後，便得到較佳的結果。這次的經驗也讓我明白資料整理及挑選的重要性，而對於該領域的 domain knowledge 也是相當重要，若是拿到 data 便開始 train 往往會得到不好的結果。