



عنوان : گزارش تمرین اول یادگیری ماشین

نگارنده : سحر داستانی اوغانی

شماره دانشجویی : ۹۹۱۱۲۱۰۸



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

دانشکده ریاضی و علوم کامپیوتر

این سوال هم به صورت کد حل شده و هم به صورت دستی. همانطور که در کد مشاهده می‌کنید، داده‌ها تولید شدند و مدل با تابع هزینه ذکر شده رسم شده است.

ابتدا داده‌ها توسط یک نویز گوسی با مرکز ۰ و میانگین ۲ تولید شدند و در فایلی با نام data ذخیره شدند. سپس برای جلوگیری از تفاوت موجود در رنج داده‌ها (scale) و عملکرد بهتر الگوریتم و نزدیک بودن آن به normal distribution، داده‌ها توسط feature scaling و با متد MinMaxScaler، نرمال شدند. برای درک تفاوت بهتر، داده‌های نرمال شده (عکس سمت راست) و

Unnamed: 0	X	T	Unnamed: 0	X	T		
0	0	-0.936773	-2.973252	0	0.000000	0.493907	0.511113
1	1	1.500040	-2.181905	1	0.010101	0.506233	0.511113
2	2	36.863630	45985.635813	2	0.020202	0.685114	0.534927
3	3	74.666853	399482.033794	3	0.030303	0.876335	0.717978
4	4	-60.206096	-229047.205798	4	0.040404	0.194104	0.392507
...
95	95	-95.046266	-885629.571085	95	0.959596	0.017871	0.052509
96	96	-3.566159	-76.751444	96	0.969697	0.480607	0.511075
97	97	-75.928006	-454948.248417	97	0.979798	0.114577	0.275529
98	98	-41.495118	-76568.334806	98	0.989899	0.288750	0.471465
99	99	44.262628	80802.423524	99	1.000000	0.722540	0.552956

داده‌های نرمال نشده (عکس سمت چپ) در این جا آمده

است. همانطور که مشاهده می‌کنید، scale داده‌ها بعد از

نرمال شدن به بالای صفر انتقال پیدا کرده است.

حال برای درج مدلی روی داده‌ها نیاز داریم آن‌ها را به دو گروه آموزش و آزمون تقسیم کنیم. برای این کار ۷۰ درصد داده‌ها را با عنوان داده‌ی آزمون و ۳۰٪ آن‌ها را با عنوان داده‌ی آموزشی برچسب زدیم. شکل مقابل، فرمت X, y داده‌های آموزش و X, y داده‌های آزمون را نمایش می‌دهد. که مشخص می‌کند از میان ۱۰۰ داده‌ی تولید شده در بالا، ۳۰ تا از آن‌ها برای آموزش مدل و ۷۰ تای دیگر برای تست به کار رفته است.

((30,)), (70,)), (30,)), (70,))

برای حل ادامه‌ی سوال‌ها نیاز است ۲ کار را فراهم کنیم:

- همانطور که در صورت سوال نیز از ما خواسته شده، باید تابعی با نام fi2 را بر روی متغیر X داده‌های خود اثر دهیم و سپس ادامه‌ی الگوریتم را بر روی داده‌های transformed شده پیاده کنیم. این کار نیز توسط PolynomialFeatures انجام گرفته است.
- باید رنج آلفا را نیز از 10^{-8} تا 10^4 فراهم کنیم.

برای مدل کردن داده‌ها در این مرحله از ماژول Ridge استفاده شده است. زیرا cost function در این ماژول SSE به همراه

regularization term می‌باشد. بردار w به ازای مقادیر مختلف لاندا (الف) بدست آمده است. برای مثال، بردار w برای الف با

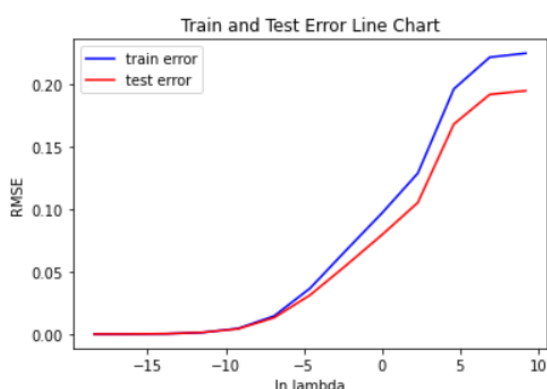
Alpha: 1e-05

RMSE: None

W= [0. 3.03404134 -5.16964706 1.55188358 1.876638 0.49052703
-0.38543553 -0.51434792 -0.23379231 0.09973821 0.24770773]

مقدار 1e-05 برابر است با:

از طرفی نمودار خطای RMSE برای لاندهای مختلف (که در اینجا بر حسب $\ln \lambda$ رسم شده است) به شرح مقابل است.



این نمودار برای لاندهای خیلی کوچک که به سمت صفر می‌روند،

فاقد جمله‌ی منظم ساز است. در جاهایی از نمودار که فاصله‌ی

میان خطای آموزش و آزمون یک حد منطقی‌ای است،

overfitting رخ داده است. و در کل میزان خطای آموزش بیشتر

از میزان خطای آزمون برای داده‌های تولید شده در بالا توسط نویز

گوسی است.

محاسبه RMSE را با اجرای ۱۰۰ بار الگوریتم بالا بر روی داده‌ها نیز انجام دادیم. مینیمم RMSE بر روی داده‌های تست زمانی رخ

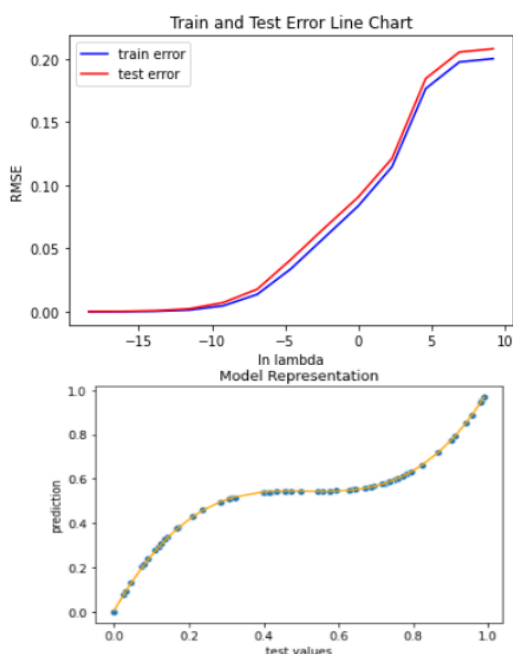
داده است که لاندا 1e-5 باشد.

	lambda	rmse_test	rmse_train
0	1.000000e-08	0.000058	0.000014
1	1.000000e-07	0.000203	0.000074
2	1.000000e-06	0.000715	0.000304
3	1.000000e-05	0.002310	0.001307
4	1.000000e-04	0.007228	0.004711
5	1.000000e-03	0.017732	0.013658
6	1.000000e-02	0.041244	0.033865
7	1.000000e-01	0.066314	0.058735
8	1.000000e+00	0.090635	0.083623
9	1.000000e+01	0.121158	0.114598
10	1.000000e+02	0.184332	0.176354
11	1.000000e+03	0.205326	0.197466
12	1.000000e+04	0.207899	0.200071

w نیز برای این مقدار به شرح زیر است:

Alpha: 1e-05
RMSE_test: 0.001086563925559874
W= [0. 3.05933674 -5.29306033 1.73477034 1.86017751 0.40639101
-0.42113947 -0.49719509 -0.20617355 0.10820386 0.24759529]

نمودار خطای آموزش و آزمون برای ۱۰۰ بار اجرای الگوریتم به نمودار زیر تغییر پیدا کرد:

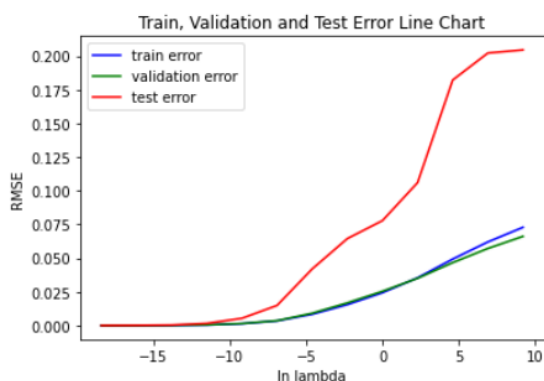


همانطور که مشاهده می‌کنید برای ۱۰۰ بار اجرای این الگوریتم، میزان خطای آزمون از خطای آموزش بیشتر شده و دو نمودار نزدیکی بیشتری به هم پیدا کردند. همین امر از **overfitting** جلوگیری کرده است. و میزان خطای کل نمودار کاهش پیدا کرده است. (از ۲۵ به ۲۰ به صورت تقریبی)

برای درک بهتر وضعیت مدل در بهترین لاند (یعنی $1e-05$)، به نمودار روبرو توجه کنید، مدل به خوبی بر روی داده‌ها **fit** شده است.

۳-الف-۲

با استفاده از ماژول **Kfold** مدل را تغییر داده و مدل **Kfold** را بر روی داده‌ها اعمال می‌کنیم. لازم به ذکر است عملیات قبل، از جمله‌ی تولید لاند، اعمال تابع **fi2** بر روی داده‌ها را نیز انجام می‌دهیم. سپس داده‌ها را به ۳ گروه **test, train, validation** تقسیم کرده و مدل **Ridge** را روی آن‌ها اعمال می‌کنیم. خطاهای **validation, test, train** را بدست آورده و آن‌ها را رسم می‌-

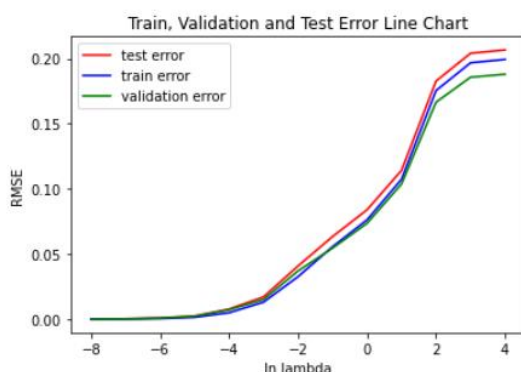


کنیم. در شکل مشخص است که با افزایش مقدار لاند، خطای **train** و **validation** دارای افزایش می‌باشند ولی خطای **test** دارای افزایش قابل تاملی است. خطای **validation, train** در حالت مجانبی قرار دارند. این بدان معنی است که چقدر مدلی که در نظر گرفتیم ساده‌تر از مدلی بوده است که دنبالش بودیم. در لندهای پایین، مقدار خطاها صفر شده است،

یعنی در این بازه لاند می‌توانیم تمام پارامترهای آزاد را دقیقاً مشخص کنیم.

این الگوریتم را ۱۰۰ بار بر روی داده‌ها انجام دادیم و نتیجه مطابق زیر حاصل شد:

	lambda	rmse_train	rmse_validation	rmse_test	همانطور که مشخص است، این الگوریتم
0	1.000000e-08	0.000016	0.000044	0.000111	
1	1.000000e-07	0.000078	0.000163	0.000253	
2	1.000000e-06	0.000344	0.000648	0.001340	زمانی که مقدار لاندا برابر با $1e-08$ بوده،
3	1.000000e-05	0.001363	0.001807	0.002616	
4	1.000000e-04	0.004909	0.007151	0.010468	
5	1.000000e-03	0.013166	0.015744	0.019924	بهترین خطا را بر روی داده‌های تست بدست
6	1.000000e-02	0.034980	0.041339	0.048257	
7	1.000000e-01	0.059182	0.059930	0.068661	آورده است.
8	1.000000e+00	0.081104	0.076624	0.089389	
9	1.000000e+01	0.112333	0.103359	0.120326	
10	1.000000e+02	0.177093	0.159224	0.183004	
11	1.000000e+03	0.198889	0.178015	0.203160	
12	1.000000e+04	0.201575	0.180314	0.205622	



نمودار ۳ خطا نیز در شکل زیر آمده است: این نمودار بیانگر

این است که با افزایش میزان داده‌ها میزان خطاها به هم

نزدیک‌تر شدند و با افزایش لاندا میزان خطای train,

validation از مقدار 0.050 به مقدار 0.20 افزایش یافته

است.

میزان میانگین خطا بر روی داده‌های validation برای

Alpha: $1e-08$

Mean_rmse_validation [0.2829484775408559, 0.22399044834884158, 0.21698283568734

605, 0.4298116498840433, 0.20962467955858827, 0.11736781583922078, 0.1168456074

9528422, 0.14216345957849008, 0.15523087162316498, 0.24340125066877513]

بهترین لاندا نیز به

شرح زیر است:

۳-ب

۵۰ مجموعه داده با اندازه‌های یکسان و مقادیر متفاوت در مرحله‌ی اول، توسط یک تابع درست شده است. ابتدا الگوریتم این

قسمت را برای ۲۰ داده و سپس برای ۲۰۰ داده آزمودیم.

در هر قسمت، لاندا را تعریف کرده و تابع $fi2$ را بر روی داده‌ها اثر داده‌ایم و مقادیر w را برای لانداهای مختلف حساب کردیم.

برای هر کدام از مجموعه‌های تولید شده در بالا و برای هر لاندا مقدار بایاس و واریانس را با تابع $bias_variance_decomp$

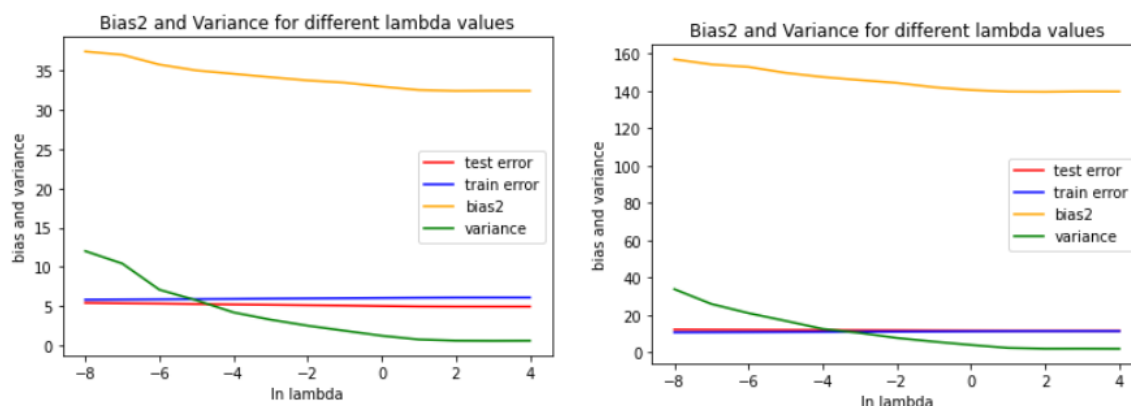
تولید کرده و میانگین خطای آزمون را به ازای مقادیر مختلف لاندا رسم کرده‌ایم.

در ادامه خلاصه‌ای از مقایسه دو خروجی (برای ۲۰ داده و ۲۰۰) را بررسی می‌کنیم:

n = 20						n = 200					
lambda	bias2	variance	rmse_test	rmse_train		lambda	bias2	variance	rmse_test	rmse_train	
0	1.000000e-08	37.444702	11.993828	5.393186	5.795533	0	1.000000e-08	157.040252	33.820659	12.171475	10.689674
1	1.000000e-07	37.023856	10.404760	5.337158	5.818018	1	1.000000e-07	154.304569	25.880113	12.081493	10.737805
2	1.000000e-06	35.776903	7.070255	5.302005	5.844100	2	1.000000e-06	152.973850	20.975604	12.010657	10.796246
3	1.000000e-05	35.028639	5.718313	5.240364	5.864766	3	1.000000e-05	149.759198	16.838562	11.924933	10.864589
4	1.000000e-04	34.599006	4.176251	5.186493	5.888916	4	1.000000e-04	147.563911	12.525667	11.853293	10.927325
5	1.000000e-03	34.163871	3.246315	5.143551	5.925189	5	1.000000e-03	145.912928	10.370205	11.805111	10.984993
6	1.000000e-02	33.754299	2.482615	5.092806	5.954407	6	1.000000e-02	144.395791	7.669827	11.759713	11.043435
7	1.000000e-01	33.472050	1.833321	5.049472	5.983139	7	1.000000e-01	142.077015	5.679347	11.674375	11.129422
8	1.000000e+00	32.967040	1.196933	4.986316	6.013483	8	1.000000e+00	140.571847	3.925590	11.613391	11.224739
9	1.000000e+01	32.535975	0.723253	4.932759	6.042040	9	1.000000e+01	139.719213	2.323354	11.579959	11.300226
10	1.000000e+02	32.410975	0.569986	4.915896	6.062902	10	1.000000e+02	139.606307	1.917477	11.573475	11.344972
11	1.000000e+03	32.432401	0.551593	4.914625	6.067633	11	1.000000e+03	139.831634	1.944614	11.574397	11.354718
12	1.000000e+04	32.418708	0.564835	4.914523	6.068167	12	1.000000e+04	139.757604	1.872290	11.574540	11.355813

همانطور که مشخص است، مقادیر بایاس، واریانس، خطای آزمون و خطای آموزش با افزایش داده، افزایش پیدا کرده است. این بدان معناست که با ثابت نگه داشتن پیچیدگی مدل و افزایش داده‌ها می‌توان از overfitting جلوگیری کرد.

نمودارهای بایاس و واریانس :



نمودار سمت راست، نتایج بایاس و واریانس را بر روی ۲۰۰ داده و نمودار سمت چپ، بر روی ۲۰ داده نشان می‌دهد. با افزایش مقدار داده، میزان واریانس و بایاس نیز افزایش یافته است. ولی فاصله میان آن دو و فاصله‌ی آن‌ها از نمودار خطای آزمون ثابت مانده است.

در انتها نیز گزارشی از مقدار w به ازای لاندای اول یعنی $1e-08$ آمده است:

```
Alpha: 1e-08
RMSE_test: 1.4954681076193808
W= [ 0. -174.86711582 1423.69011841 -4245.78107696
      2882.08750489 5789.7287718 -3845.98168043 -8921.3021034
      1186.25511 13427.64851686 -7526.49673701]
-----
Alpha: 1e-08
RMSE_test: 6.116031831721433
W= [ 0. 82.72103301 -1213.73684948 4501.31190234
      -2713.36745896 -9556.66623061 5776.77310387 15497.06424198
      -1796.45213017 -25039.46842805 14453.96000055]
```

پاسخ تمارین تشریحی

این خطای Validation هوای کمتر از خطای آموزش است.

نادرست نیست. داده های مسئله تنها با داده های آموزش، آموزش می بیند و توسط داده های val و test مدل آموزش قرار

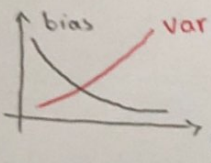
می بیند (استه می توان گفت در تئوری Model-Selection، خطای خطی مدل را با val تطبیق می دهیم پس به نوعی با آن

تیر آموزش می یابیم). در error، test و val برای این که بتواند تعمیمی خوبی از داده های مسئله با به یک order

قرار بدهد. بنابراین می توان رابطه ای میان آن ها بیان کرد و آنها می توان گفت بالا بردن train error مقدار داده یا حجم

به افتادن خطای مدل در رگرسیون همیشه سبب افزایش دارباز و کاهش بایاس می شود.

اگر برای این حجم این شرط را بگذاریم که خطای را ثابت در نظر بگیریم و بخواهیم خطای را بکاهش دهیم، این عمل به عنوان زیر



Model Complexity

ولی اگر بخواهیم Regularization Term را به آن اضافه کنیم، این کار

می تواند خیلی سبب زیاد را بگذارد و معیوس خط ها را محدود کند، بنابراین var کم می شود و حالت استثنای این مثال

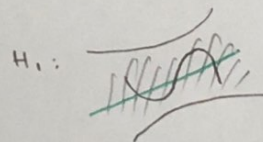
$H_0: f(n) = b$ مدل منفر

$H_1: f(n) = an + b$ مدل اول

زیر سؤال می شود. یعنی نادرست می شود.

بدون Regularization term

$H_0: \text{bias} = 0.5, \text{var} = 0.25$



$H_1: \text{bias} = 0.21, \text{var} = 1.69$

$H_1: \text{bias} = 0.23, \text{var} = 0.33 \rightarrow$ کم شود

Regularization term

$H_0 \rightarrow \text{bias} + \text{var} = 0.75$

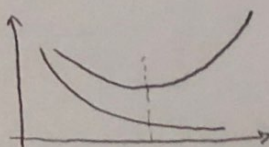
$H_1 \rightarrow \text{bias} + \text{var} = 1.90$

$H_1 + \text{Regularization term} \rightarrow 0.56$

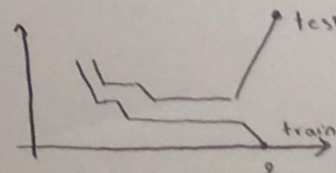
پس افزایش خطای مدل در رگرسیون همیشه سبب کاهش خطای آموزش و افزایش خطای آموزش می شود.

دقت خطای مدل افزایش می یابد، مدل روی داده های train تطبیق بیشتری پیدا می کند. با افزایش خطای و ثابت

مانند مقدار داده ها. overfit می شود، خطای داده های train منفر شده ولی روی داده های test



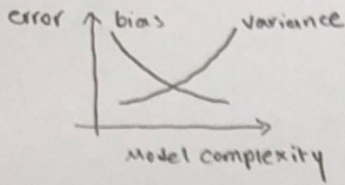
Polynomial degree



Polynomial degree

نت) در صورتی که مقدار بایاس بالا باشد و واریانس کم باشد، بایاس پیچیدگی مدل را کاهش داده.

نادرست است. در حالتی که مقدار بایاس بالا باشد واریانس کم، مدل دچار under fitting شده و این یعنی به خوبی نتوانسته داده‌ها را توصیف کند. بنابراین وقتی مقدار مقابله بایاس پیچیدگی مدل را افزایش داد تا مقدار error کم شود fit بهتری حاصل شود.



نت) کاهش فضای آموزش منجر به کاهش فضای آزمون می‌شود. درست است.

در حقیقت فضای آموزش کاهش می‌یابد، مدل fit بهتری روی داده‌ها می‌دهد و از طرفی دیگر فضای آزمون کاهش می‌یابد و این به ضرر است نه فضای آزمون کمتر شود زیرا در آن صورت overfit رخ داده و فضای آزمون به طور ناگهانی افت می‌یابد.

نت) در صورتی که واریانس بالا باشد، افزایش داده‌های آموزش همیشه منجر به بهبود بایاس و واریانس می‌شود.

درست. واریانس بالا منجر به overfit می‌شود، با افزودن داده‌های آموزش مدل بهتری می‌یابد و بایاس کاهش می‌یابد. اما واریانس همچنان بالا می‌ماند و این به ضرر است نه overfit کمتر رخ می‌دهد پس واریانس کاهش می‌یابد.

نت) در صورتی که بایاس بالا باشد، افزایش داده‌های آموزش هیچ کمکی به کاهش بایاس نمی‌کند و بایاس پیچیدگی مدل را افزایش داده.

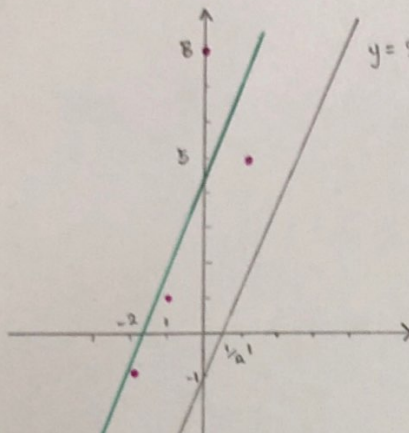
درست است. در صورتی که بایاس بالا باشد یعنی دچار underfit شده‌ایم پس یعنی پیچیدگی مدل در آن نقطه این قدر data کافی نبوده و train error بالاست.

با افزایش مقدار داده‌های آموزش این error بالاتر می‌رود پس بایاس پیچیدگی مدل را افزایش می‌دهد نه fit بهتری می‌دهد.

$$y = 2x - 1$$

$$\text{data} = \{(-1, 1), (0, 8), (1, 5), (-2, -1)\}$$

الف ← 1



رسم تابع و نقاط داده شده \Rightarrow مقدار رگرسیون:

$$y = 2.5x + 4.5$$

$$J(w) = \sum_{i=1}^n (y^{(i)} - w_0 - w_1 x^{(i)})^2$$

2) تابع هزینه SSE به شرح مقابل است:

مثل مسئله مدلی خط است ←

$$f(x^{(i)}; w) = w_0 + w_1 x^{(i)}$$

باید به سبب آردن خط به ازای کمینه شدن تابع هزینه از cost-function نسبت به پارامترهای آن (w_0, w_1) مشتق بگیریم. تعدادی برابر صفر قرار دهیم تا معادله پارامتر را به دست آوریم.

$$\left\{ \begin{aligned} \frac{\partial J(w)}{\partial w_0} &= \sum_{i=1}^n 2 (y^{(i)} - w_0 - w_1 x^{(i)}) (-1) = 0 \\ \frac{\partial J(w)}{\partial w_1} &= \sum_{i=1}^n 2 (y^{(i)} - w_0 - w_1 x^{(i)}) (-x^{(i)}) = 0 \end{aligned} \right.$$

$$\left\{ \begin{aligned} -2 \sum_{i=1}^n y^{(i)} + 2 \sum_{i=1}^n w_0 + 2 \sum_{i=1}^n w_1 x^{(i)} &= -2 \sum_{i=1}^n y^{(i)} + 2 w_1 \sum_{i=1}^n x^{(i)} + 2 n w_0 = 0 \\ -2 \sum_{i=1}^n x^{(i)} y^{(i)} + 2 \sum_{i=1}^n w_0 x^{(i)} + 2 \sum_{i=1}^n w_1 x^{(i)2} &= -2 \sum_{i=1}^n x^{(i)} y^{(i)} + 2 w_1 \sum_{i=1}^n x^{(i)2} + 2 w_0 \sum_{i=1}^n x^{(i)} = 0 \end{aligned} \right.$$

$$\left\{ \begin{aligned} -2(1+8+5-1) + 2w_1(-1+0+1-2) + 2 \times 4 w_0 &= 0 \\ -2\left(\frac{-1 \times 1}{-1} + \frac{0 \times 8}{0} + \frac{1 \times 5}{5} + \frac{-2 \times -1}{2}\right) + 2w_1(1+0+1+4) + 2w_0(-1+0+1-2) &= 0 \end{aligned} \right.$$

$$\left\{ \begin{aligned} -26 - 4w_1 + 8w_0 &= 0 \\ -12 + 12w_1 - 4w_0 &= 0 \end{aligned} \right. \rightarrow \left\{ \begin{aligned} -2w_1 + 4w_0 &= 13 \\ 12w_1 - 4w_0 &= 12 \end{aligned} \right. \begin{aligned} w_0 &= 4.5 \\ w_1 &= 2.5 \end{aligned}$$

$$y = 2.5x + 4.5$$

$$w^* = \operatorname{argmin}_{x,y} E_{x,y} [(y - w^T x)^2]$$

ب ←

$$R = E_x [x x^T], \quad C = E_{x,y} [x y]$$

(1)

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}, \quad x = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & x_1^{(2)} & \dots & x_d^{(2)} \\ \vdots & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix}, \quad w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}$$

$$\|a\|^2 = a_1^2 + \dots + a_d^2 \quad \leftarrow \text{مقدار}$$

$$\|y - xw\|^2 = (y^{(1)} - w^T x^{(1)})^2 + \dots + (y^{(n)} - w^T x^{(n)})^2 = \sum_{i=1}^n (y^{(i)} - w^T x^{(i)})^2 = J(w)$$

که به دنبال بهینه کردن $J(w)$ به فرم بالا هستیم

$$\nabla_w J(w) = 2x^T (y - xw) = 0$$

$$x^T y - x^T x w = 0 \rightarrow x^T x w = x^T y \rightarrow w = (x^T x)^{-1} x^T y = x^+ y$$

$$w^* = (x^T x)^{-1} x^T y = x^+ y$$

در نهایت که $J(w)$ بهینه است یعنی w^* داریم پس :

Pseudo inverse matrix

$$x^T \quad x \quad x = \begin{bmatrix} \quad \end{bmatrix} \quad \leftarrow \quad x = \begin{bmatrix} \quad \end{bmatrix} \quad \begin{matrix} \sim \\ \text{دور} \end{matrix}$$

در این مسئله که با آن روبرو هستیم، تعداد ویژگی‌ها محدود و تعداد Training data بسیار زیاد است. در این نوع مسائل کمترین مقایسه که $d+1$ به سادگی (معمولاً در وقت بالادار با این)

$$(2) \quad \text{مقدار} \quad w^* = \operatorname{argmin}_{x,y} E_{x,y} [(y - w^T x)^2] \quad \text{برای} \quad \text{Minimize} \quad w^* \quad \text{با} \quad \text{این} \quad \text{نسبت} \quad \text{به}$$

$$\nabla_w w^* = -2x (y - w^T x) = 0 \quad \leftarrow \text{که در این مسئله}$$

$$E_{x,y} [-2x (y - w^{*T} x)] = 0 \quad \leftarrow \text{در نهایت که} \quad w^* \quad \text{بهینه است عبارت بالا برابر صفر می‌شود یعنی}$$

$$E_{n,y} [(y - \hat{\omega}^T x)^2] = \underbrace{E_{x,y} [(y - \omega^{*T} x)^2]}_{\text{Structural error}} + \underbrace{E_n [(\omega^{*T} x - \hat{\omega}^T x)^2]}_{\text{Approximation error}} \quad (3)$$

$$\begin{aligned} E_{x,y} [(y - \hat{\omega}^T x)^2] &= E_{n,y} [(y - \omega^{*T} x + \omega^{*T} x - \hat{\omega}^T x)^2] \\ &= E_{n,y} [(y - \omega^{*T} x)^2] + E_n [(\omega^{*T} x - \hat{\omega}^T x)^2] + 2 \underbrace{E_{n,y} [(y - \omega^{*T} x)(\omega^{*T} x - \hat{\omega}^T x)]}_A \end{aligned}$$

: $\bar{A} = 0$ (معمولاً)

$$\begin{aligned} A &= E_{n,y} [(y - \omega^{*T} x)(\omega^{*T} x - \hat{\omega}^T x)] \\ &= E_{n,y} [(y - \omega^{*T} x)(\omega^* - \hat{\omega})^T x] \xrightarrow{a^T b = b^T a} \\ &= E_{n,y} [(y - \omega^{*T} x) x^T (\omega^* - \hat{\omega})] \\ &= (\omega^* - \hat{\omega}) E_{n,y} (y - \omega^{*T} x) = 0 \end{aligned}$$

$\leftarrow 0$ 2 خطه

(4) $\bar{A} = 0$ (معمولاً)

True error (خطای واقعی) distribution (توزیع)

Approximation error (خطای تقریب) (Training data) (داده های آموزشی)

خطای تقریب (Approximation error) خطای واقعی (True error) خطای تقریب (Approximation error)

خطای تقریب (Approximation error) خطای واقعی (True error) خطای تقریب (Approximation error)

الف) ابتدا نشان دهیم تابع هزینه پرسپترون را می توان به صورت $\sum_{i=1}^n \max(0, y^{(i)} w^T x^{(i)})$ نوشت، سپس تابع هزینه متعامد با آن را در حالت $y=1$ به حسب $x^T w$ رسم کنیم.
 نتیجه: تابع هزینه Logistic regression (در حالت $y=1$ به حسب $x^T w$) متناسب است و تابع این تفاوت را در عملکرد روش های مربوطه مشخص کنیم.

میانگین Cost function = پرسپترون به شرح زیر است:

$$J_P(w) = - \sum_{i \in M} w^T x^{(i)} y^{(i)}$$

که در این رابطه، M نشان دهنده نمونه های دسته N غلط دسته بندی شده اند (misclassified).

و T را نمونه هایی در تقسیم به N درست دسته بندی شده (True classified).

بنابراین Cost function که متعامد برابر با تابع هزینه غلط دسته بندی N درست دسته بندی N + تابع هزینه غلط دسته بندی غلط دسته بندی شده.

$$= \sum_{i \in T} w^T x^{(i)} y^{(i)} + - \sum_{i \in M} w^T x^{(i)} y^{(i)}$$

می دانیم مقدار تابع هزینه غلط دسته بندی N درست دسته بندی شده برابر منفی است و چون $\text{sign}(w^T x^{(i)})$ برای داده هایی که

غلط دسته بندی شده منفی است، پس عبارت $\sum_{i \in M} w^T x^{(i)} y^{(i)}$ مثبت می شود و می توان تابع هزینه Perceptron

را به صورت متعامد داریم که $\sum_{i \in M} w^T x^{(i)} y^{(i)}$ منفی باشد، حاصل منفی را در جمع این موارد

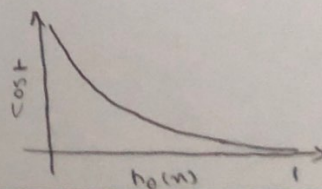
$$\sum_{i=1}^n \max(0, y^{(i)} w^T x^{(i)})$$

حاصل می شود $-\sum_{i \in M} w^T x^{(i)} y^{(i)}$

بنابراین Cost function = Logistic Regression به شرح زیر است:

$$J(\theta) = -\frac{1}{n} \left[\sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)})) \right]$$

$$J(\theta) = -\frac{1}{n} \sum_{i=1}^m \log(h_{\theta}(x^{(i)})) \rightarrow$$



در صورت $y=1$ داریم:

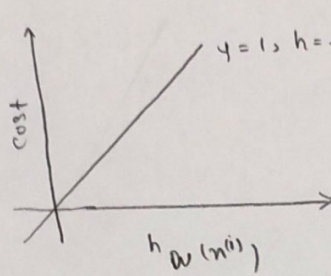
مثلاً اگر $y = 1$ ، تابع هزینه پرسپترون عبارت از این می باشد :

$$J_p(w) = - \sum w^T x^{(i)} y^{(i)}$$

به دلیل این که پارامترهای تابع هزینه Perceptron مثل

$$y = 1 \rightarrow J_p(w) = - \sum w^T x^{(i)}$$

$$= - \sum_{i=1}^n h_w(x^{(i)})$$



نتیجه : به دلیل این که error در logistic regression عواره مثبت است ، بنابراین اگر ما نسبت را از طریق perceptron ،

فقط نمونه های درست از برای w محاسبه می شود . معجزه در تقارن می بینیم . ممکن است الگوریتم در بین فکلی گلی استفاده و جواب بین را اعلام کند

1) نشان دهید نمونه‌های کاملاً جداپذیر خطی، جداپذیر خطی نیز هستند، اما برعکس آن لزوماً درست نیست.

2) نشان دهید نمونه‌های دو به دو خطی جداپذیر، لزوماً خطی جداپذیر نیستند. آیا بالعکس نمونه‌های خطی جداپذیر لزوماً دو به دو خطی جداپذیر نیز هستند؟

1) ابتدا تعریفی از هریت مسئله بهینه‌سازی می‌آوریم سپس آن را حل می‌کنیم.

$$\left. \begin{aligned} \mu_1 &= \frac{\sum_{x^{(i)} \in C_1} x^{(i)}}{N_1} \rightarrow \mu'_1 = W^T \mu_1 & \text{MAP شده‌ی نقاط کلاس 1 روی بردار} \\ \mu_2 &= \frac{\sum_{x^{(i)} \in C_2} x^{(i)}}{N_2} \rightarrow \mu'_2 = W^T \mu_2 & \text{MAP شده‌ی نقاط کلاس 2 روی بردار} \end{aligned} \right\} \begin{aligned} \mu'_1 &= \text{میانگین} \\ \mu'_2 &= \text{میانگین} \end{aligned}$$

$$\left. \begin{aligned} S_1'^2 &= \text{واریانس MAP شده‌ی کلاس 1} \\ S_1^2 &= \text{واریانس ضد داده‌ها در فضای اصلی است (برای داده‌های کلاس 1)} \end{aligned} \right\}$$

$$S_1^2 = \sum_{x^{(i)} \in C_1} \|x^{(i)} - \mu_1\|^2$$

$$S_1'^2 = \sum_{x^{(i)} \in C_1} (W^T x^{(i)} - W^T \mu_1)^2$$

دقیقاً برای کلاس 2.

$$J(w) = \frac{|\mu_1' - \mu_2'|^2}{\delta_1'^2 + \delta_2'^2} \rightarrow \max J(w)$$

راه اول: \rightarrow برای دستیابی به بیشترین دقت، باید به سمت μ_1' و μ_2' حرکت کرد.

در واقع هدف این بود که میانگین داده‌ها در MAP به سمت μ_1' و μ_2' حرکت کند. \rightarrow یعنی داده‌های μ_1' و μ_2' در سمت راست و چپ قرار بگیرند و به سمت μ_1' و μ_2' حرکت کنند.

$$|w^T \mu_1 - w^T \mu_2|^2 = w^T \overbrace{(\mu_1 - \mu_2)(\mu_1 - \mu_2)^T}^{S_B} w$$

S_B به ماتریس $d \times d$ است، چون بردار میانگین داده‌ها μ_1 و μ_2 به هم ضرب شده است.

$$\sum_{x^{(i)} \in C_1} (w^T x^{(i)} - w^T \mu_1)^2 + \sum_{x^{(i)} \in C_2} (w^T x^{(i)} - w^T \mu_2)^2$$

$$= w^T \left(\sum_{x^{(i)} \in C_1} (x^{(i)} - \mu_1)(x^{(i)} - \mu_1)^T \right) w + w^T \left(\sum_{x^{(i)} \in C_2} (x^{(i)} - \mu_2)(x^{(i)} - \mu_2)^T \right) w$$

$$= w^T \left[\sum_{x^{(i)} \in C_1} (x^{(i)} - \mu_1)(x^{(i)} - \mu_1)^T + \sum_{x^{(i)} \in C_2} (x^{(i)} - \mu_2)(x^{(i)} - \mu_2)^T \right] w = w^T \underbrace{S_w}_{S_1 + S_2} w$$

حال برای به دست آوردن w به سمت μ_1' و μ_2' حرکت کرد و برای به دست آوردن w به سمت μ_1' و μ_2' حرکت کرد.

$$J(w) = \frac{w^T S_B w}{w^T S_w w} \rightarrow \frac{\partial J(w)}{\partial w} = \frac{\frac{\partial w^T S_B w}{\partial w} \times w^T S_w w - \frac{\partial w^T S_w w}{\partial w} \times w^T S_B w}{(w^T S_w w)^2}$$

$$= \frac{(2 S_B w) \overbrace{w^T S_w w}^{\beta} - (2 S_w w) \overbrace{w^T S_B w}^{\alpha}}{(w^T S_w w)^2} = 0 \rightarrow \boxed{S_B w = \lambda S_w w} \rightarrow \lambda = \frac{\alpha}{\beta}$$

\downarrow جواب

التر متجه S_w معكوس S_B بالحد

$$S_B W = \lambda S_w W \rightarrow S_w^{-1} S_B W = \lambda W$$

$$S_B W = (\mu_2 - \mu_1) (\mu_2 - \mu_1)^T W \propto (\mu_2 - \mu_1)$$

$S_B W$ را با S_w بسط

$$\rightarrow S_B W = \alpha (\mu_2 - \mu_1) = S_w W \rightarrow W \propto S_w^{-1} (\mu_2 - \mu_1)$$

$$\rightarrow \boxed{W \propto (\mu_2 - \mu_1)}$$

توضیح دلیل چرا $\|W\| = 1$ \vec{u} \vec{v}

همچنین W را نزدیک به \vec{u} ، فاصله نزدیک و نزدیک به \vec{v} . پس چون W موازی این بردارها

$$\vec{u} \quad \vec{v} \quad \|W\| = W^T W = 1 \quad \text{در نظر میگیریم}$$

$$u_k = W^T u_k, \quad u_2 - u_1 = W^T (u_2 - u_1)$$

از W \leftarrow در این رابطه W را به W^T با استفاده از 2 جابجایی

ماتریس W را W^T :

$$L = W^T (\mu_2 - \mu_1) + \lambda (W^T W - 1)$$

$$\nabla L = \mu_2 - \mu_1 + 2\lambda W = 0 \rightarrow 2\lambda W = \mu_1 - \mu_2 \rightarrow W = \frac{1}{2\lambda} (\mu_1 - \mu_2)$$

$$\rightarrow W = \frac{-1}{2\lambda} (\mu_2 - \mu_1) \rightarrow \boxed{W \propto \mu_2 - \mu_1}$$