

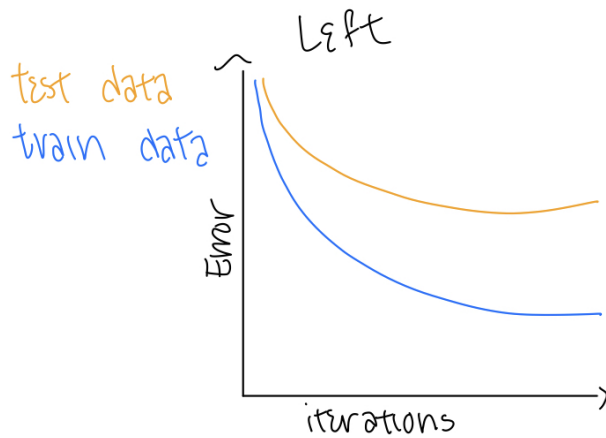
CS M148 Homework 2

Hanna Co

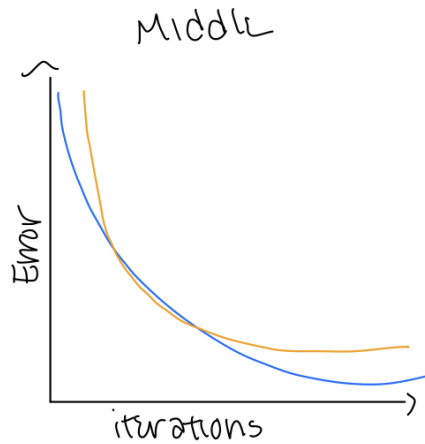
Due: February 16, 2021

1 Bias, Variance and Regularization

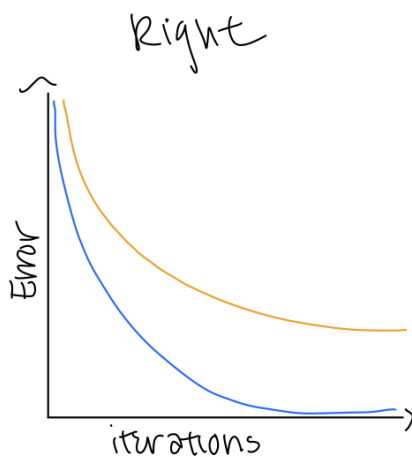
a) The graph on the far left has the largest bias and the lowest variance, and the graph on the far right has the largest variance and the lowest bias.



The graph on the left will perform bad on both the train and test data, because it's too generalized.



The graph in the middle will perform alright on both train and test data. This is because it's well fitted to the training data, but not overfitted to the point that it performs badly on test data.



The graph on the right will perform really good on the training data, but this is because it is overfitted to the training data. Thus, it will not perform well on the test data.

b) L1 regularization is used on a), because the coefficients are nullified fast, compared to b), where they are not nullified.

2 Maximum Likelihood View of Linear Regression

a) Since the likelihood function represents the probability of producing a particular sample, thus the equation is:

$$\prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{-1}{2}\left(\frac{y_i - \hat{y}_i}{\sigma}\right)^2\right)$$

b) We take the equation in 2a and change it into the log likelihood equation:

$$\begin{aligned} & \sum_{i=1}^n \ln\left(\frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y_i - \hat{y}_i)^2}{2\sigma^2}\right)\right) \\ & \sum_{i=1}^n \ln\left(\frac{1}{\sigma\sqrt{2\pi}}\right) - \frac{(y_i - \hat{y}_i)^2}{2\sigma^2} \\ & \sum_{i=1}^n -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{(y_i - \hat{y}_i)^2}{2\sigma^2} \\ & -\frac{n}{2} \ln(2\pi\sigma^2) - \sum_{i=1}^n \frac{(y_i - \hat{y}_i)^2}{2\sigma^2} \end{aligned}$$

We want to find $\text{argmax}(L)$:

$$\text{argmax}(L) = \text{argmax}\left(-\sum_{i=1}^n \frac{(y_i - \hat{y}_i)^2}{2\sigma^2}\right)$$

Since we will end up setting $\text{argmax}(L) = 0$, we can simply write

$$\text{argmax}(L) = \text{argmax}\left(-\sum_{i=1}^n (y_i - \hat{y}_i)^2\right)$$

Since $\text{argmax}(L) = \text{argmin}(-L)$,

$$\text{argmax}(L) = \text{argmin}\left(\sum_{i=1}^n (y_i - \hat{y}_i)^2\right)$$

Additionally, n is a constant, so we have

$$\text{argmax}(L) = \text{argmin}\left(\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2\right)$$

\hat{y}_i is our model, so we substitute \hat{y}_i with $\beta_0 + \beta_1 x_i$

$$\text{argmax}(L) = \text{argmin}\left(\frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2\right)$$

Notice that the right hand side of the above equation is the equation for MSE. Thus, maximizing log likelihood is equivalent to maximizing MSE.

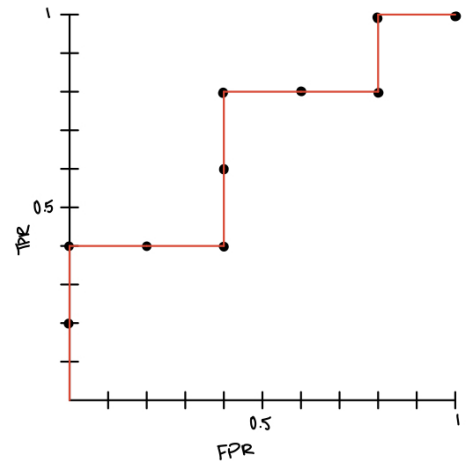
3 Classification Metric

a) ROC Curve

TPR	FPR	threshold	TPR	FPR	threshold
0.2	0	0.98	0.8	0.4	0.59
0.4	0	0.92	0.8	0.6	0.55
0.4	0.2	0.83	0.8	0.8	0.52
0.4	0.4	0.77	1	0.8	0.32
0.6	0.4	0.62	1	1	0.13

$$TPR = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{FP + TN}$$



b) $(0.4)(1) + (0.4)(0.6) + (0.2)(0.2) = 0.68$

c) Confusion Matrix

		True	
		Pos	Neg
Predicted	Pos	4	4
	Neg	1	1

d) $\text{Accuracy} = \frac{TP+TN}{TP+FP+TN+FN} = \frac{4+1}{4+1+4+1} = 0.5$

$\text{Precision} = \frac{TP}{TP+FP} = \frac{4}{4+1} = 0.5$

$\text{Recall} = \frac{TP}{TP+FN} = \frac{4}{4+1} = 0.8$

$\text{F1 Score} = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}} = 2 * \frac{0.5 * 0.8}{0.5 + 0.8} = 0.615$

e) Yes, we can improve the scores by changing our threshold to any value above 0.55, and less than or equal to 0.59. This changes two of our false positives to true negatives. Our new accuracy is 0.7, our new precision is 0.67, our new f1 score is 0.73, and our recall remains at 0.8.

4 4 K-Nearest Neighbors for Classification

	precision	recall	f1-score	support
0	1.00	0.50	0.67	6
1	1.00	1.00	1.00	3
2	0.50	1.00	0.67	3
3	0.33	0.20	0.25	5
4	0.00	0.00	0.00	0
5	1.00	1.00	1.00	2
6	1.00	1.00	1.00	2
7	1.00	0.67	0.80	3
8	0.75	1.00	0.86	3
9	1.00	0.50	0.67	2
10	1.00	1.00	1.00	4
11	1.00	1.00	1.00	3
12	0.33	0.25	0.29	4
13	1.00	1.00	1.00	2
14	0.33	1.00	0.50	2
15	0.50	0.50	0.50	2
16	1.00	0.80	0.89	5
17	0.60	1.00	0.75	3
18	1.00	1.00	1.00	1
19	1.00	0.67	0.80	3
20	1.00	1.00	1.00	5
21	1.00	1.00	1.00	3
22	0.50	0.20	0.29	5
23	1.00	1.00	1.00	1
24	1.00	1.00	1.00	4
25	1.00	1.00	1.00	2
26	1.00	1.00	1.00	2
27	1.00	0.80	0.89	5
28	1.00	1.00	1.00	3
29	1.00	1.00	1.00	4
30	1.00	1.00	1.00	2
31	1.00	0.75	0.86	4
32	1.00	1.00	1.00	1
33	1.00	1.00	1.00	5
35	1.00	0.75	0.86	4
36	1.00	1.00	1.00	2

a) The numbers are classes, and the lines include the precision, recall, f1 score and support for each class.

b) For the most part, yes, though there are a few classes where it doesn't perform well. For example, class=12 and class=39 do not perform well in terms of the reported metrics.

c) Yes, the results would be different if the lighting or angles of the faces varied more. This gives the model more variety to train on. The results could change with a different background – if the contrast was similar to the current background, the results would probably be very similar. However, if it was changed to a low contrast background, it would likely make the model worse, because it becomes harder to differentiate between the face and the background.

5 Logistic Regression

a) We have the equation $\ln\left(\frac{P(Y=1)}{1-P(Y=1)}\right) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$.

Since $X_1 = X_2 = 0$, this equation becomes $\ln\left(\frac{P(Y=1)}{1-P(Y=1)}\right) = \beta_0$.

$$\ln\left(\frac{P(Y=1)}{1-P(Y=1)}\right) = 3.$$

$$\frac{P(Y=1)}{1-P(Y=1)} = e^3$$

$$P(Y=1) = (1 - P(Y=1))e^3$$

$$P(Y=1) = e^3 - e^3 P(Y=1)$$

$$P(Y=1) + e^3 P(Y=1) = e^3$$

$$P(Y=1) * (1 + e^3) = e^3$$

$$P(Y=1) = \frac{e^3}{1+e^3}$$

$$P(Y=1) = 0.953$$

$$\frac{P(Y=1)}{1-P(Y=1)} = 20.086$$

The probability of the event $Y=1$ is 0.953, and the odds are 20.086.

b) A one unit increase in X_1 increases our log odds by 2 and our odds by a factor of e^2 . A one unit increase in X_2 decreases our log odds by 5 and our odds by a factor of e^{-5} .

c) Increasing $\beta_0, \beta_1, \beta_2$ will increase both our odds and log odds, while decreasing them would decrease our odds and log odds. Increasing β_0 by one will increase our log odds by 1 and our odds by a factor of e^1 , and decreasing β_0 will decrease our log odds by 1 and our odds by a factor of e^1 . As for β_1 and β_2 , assuming that X_1 and X_2 are kept the same, then increasing β_1 by one will increase our log odds by X_1 , and our odds by e^{X_1} and decreasing β_1 by one will decrease our log odds by X_1 . and our odds by e^{X_1} . Similarly, increasing β_2 by one will increase our log odds by X_2 and our odds by a factor of e^{X_2} , and decreasing β_2 by one will decrease our log odds by X_2 and our odds by e^{X_2} .

d) The formulation of our decision boundary is $\beta_0 + \beta_1 X_1 + \beta_2 X_2 = 3 + 2X_1 - 5X_2$. Points on the decision boundary are points where $P(Y=1) = P(Y=0)$. For example, the point (1, 1) is on the decision boundary.

e) The coefficients changing is indicative of multicollinearity. This is potentially problematic because it undermines the significance of a single variable. It can also give the illusion of statistical significance.

6 Logistic Regression with Interaction Term

a) The intercept indicates that on average, a mother who is 23 and made infrequent visits to the physician during the first trimester has a 0.52 less chance of having a baby with low birth weight. The age coefficient indicates that a one unit increase in age will, on average, produce 0.04 greater odds that a mother will have a baby with low birth weight. The frequency coefficient indicates that a mother who visits the physician frequently during the first trimester will on average, produce 0.47 less chance of having a baby with low birth weight. The age x frequency coefficient indicates that a one unit increase in age*frequency will on average result in 0.18 less chance of having a baby with low birth weight.

b) The model is $-0.52 + (0.04)(\text{Age}) + (-0.47)(\text{Frequency}) + (-0.18)(\text{Age} \times \text{Frequency})$. When a mother visits the physician frequently during the first trimester, the model is $-0.99 + (-0.14)(\text{Age})$. For this model, a one unit increase in the mother's age produces an average of 0.14 less chance of having a baby with a low birth weight. When a mother visits the physician infrequently during the first trimester, the model is $-0.52 + (0.04)(\text{Age})$, where a one unit increase in the mother's age produces an average of 0.04 greater chance of having a baby with a low birth weight.

c) The odds ratio is calculated as follows:

$$\ln\left(\frac{P(Y=1)}{P(Y=0)}\right) = -0.52 + (0.04)(\text{Age}) + (-0.47)(\text{Frequency}) + (-0.18)(\text{Age} \times \text{Frequency})$$

where we set age to a particular value, and take the quotient of when Frequency = 1 and Frequency = 0.

$$\text{Age 17: } \ln\left(\frac{P(Y=1)}{P(Y=0)}\right) = -0.52 + (0.04)(-6) + (-0.47)(\text{Frequency}) + (-0.18)(-6 * \text{Frequency})$$

$$\ln\left(\frac{P(Y=1)}{P(Y=0)}\right) = -0.52 + -0.24 + (-0.47)(\text{Frequency}) + (1.08)(\text{Frequency})$$

$$\ln\left(\frac{P(Y=1)}{P(Y=0)}\right) = -0.76 + (0.61)(\text{Frequency})$$

$$\text{Odds Ratio: } \frac{e^{-0.15}}{e^{-0.76}} = 1.840$$

To compute the odds ratio for the other ages, simply replace Age with 23-age.

Age	Odds Ratio	95% Confidence Interval
17	1.840	(0.705, 4.949)
23	0.625	(0.325, 1.201)
24	0.522	(0.262, 1.036)
25	0.436	(0.206, 0.916)
30	0.177	(0.050, 0.607)

d) An odds ratio of 1 indicates that both events have an equal probability of occurring. Thus, in our odds ratio table, the column indicates the probability of a mother who makes frequent physician visits during the first trimester have a baby with low birth weight, divided by the probability for a mother who did not make frequent visits. For example, for mothers age 17, the probability of a mother who makes frequent physician visits having a baby with low birth weigh has on average, 1.840 times the probability of a mother who does not make frequent physician visits, also age 17, having a baby with low birth weight. A number under 1 indicates that the event is less likely to occur, while an odds ratio greater than 1 indicates that the event is more likely to occur. An odds ratio that falls within the confidence interval indicates that the odds ratio is statistically significant.

e) The difference in probability is calculated as follows:

$$\frac{e^{-0.99-0.14(Age)}}{e^{-0.99-0.14(Age)} + 1} - \frac{e^{-0.52+0.04(Age)}}{1 + e^{-0.52+0.04(Age)}}$$

To compute the difference in probability for the other ages, simply replace Age with 23-age.

Age	Difference in Probability	95% Confidence Interval
17	0.144	(-0.788,0.393)
23	-0.102	(-0.197,0.088)
24	-0.138	(-0.232,0.046)
25	-0.173	(-0.315,-0.016)
30	-0.318	(-0.540,-0.092)

The difference in probability is a bit more self explanatory. It is the difference between the probability of a mother of a certain age that makes frequent physician visits having a baby with low birth weight and a mother of the same age but makes infrequent physician visits having a baby with low birth weight. The results are consistent with those from part c: mothers of age 17 that make frequent physician visits are more likely to have a baby with low birth weight compared of 17 year old mothers who don't make frequent

physician visits. For mothers of age 30, the difference in probability is the greatest compared to other ages in the table. However, for part c, all of the results were statistically significant, while not all results for difference in probability are statistically significant.

7 Multinomial Logistic Regression

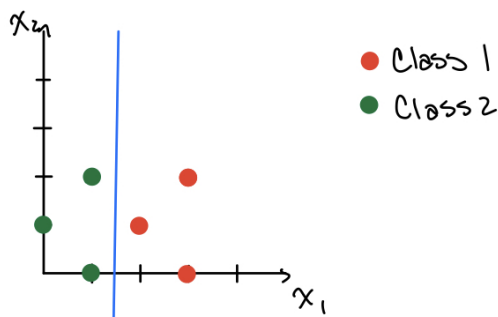
MNLogit Regression Results						
Dep. Variable:	PID	No. Observations:	944			
Model:	MNLogit	Df Residuals:	908			
Method:	MLE	Df Model:	30			
Date:	Tue, 15 Feb 2022	Pseudo R-squ.:	0.1648			
Time:	10:39:25	Log-Likelihood:	-1461.9			
converged:	True	LL-Null:	-1750.3			
Covariance Type:	nonrobust	LLR p-value:	1.822e-102			
PID=1	coef	std err	z	P> z	[0.025	0.975]
logpopul	-0.0115	0.034	-0.336	0.736	-0.079	0.056
selfLR	0.2977	0.094	3.180	0.001	0.114	0.481
age	-0.0249	0.007	-3.823	0.000	-0.038	-0.012
educ	0.0825	0.074	1.121	0.262	-0.062	0.227
income	0.0052	0.018	0.295	0.768	-0.029	0.040
const	-0.3734	0.630	-0.593	0.553	-1.608	0.861
PID=2	coef	std err	z	P> z	[0.025	0.975]
logpopul	-0.0888	0.039	-2.266	0.023	-0.166	-0.012
selfLR	0.3917	0.108	3.619	0.000	0.180	0.604
age	-0.0229	0.008	-2.893	0.004	-0.038	-0.007
educ	0.1810	0.085	2.123	0.034	0.014	0.348
income	0.0479	0.022	2.149	0.032	0.004	0.092
const	-2.2509	0.763	-2.949	0.003	-3.747	-0.755
PID=3	coef	std err	z	P> z	[0.025	0.975]
logpopul	-0.1060	0.057	-1.858	0.063	-0.218	0.006
selfLR	0.5735	0.159	3.617	0.000	0.263	0.884
age	-0.0149	0.011	-1.311	0.190	-0.037	0.007
educ	-0.0072	0.126	-0.057	0.955	-0.255	0.240
income	0.0576	0.034	1.713	0.087	-0.008	0.123
const	-3.6656	1.157	-3.169	0.002	-5.932	-1.399

PID=4	coef	std err	z	P> z	[0.025	0.975]
logpopul	-0.0916	0.044	-2.091	0.037	-0.177	-0.006
selfLR	1.2788	0.129	9.921	0.000	1.026	1.531
age	-0.0087	0.008	-1.031	0.302	-0.025	0.008
educ	0.1998	0.094	2.123	0.034	0.015	0.384
income	0.0845	0.026	3.226	0.001	0.033	0.136
const	-7.6138	0.958	-7.951	0.000	-9.491	-5.737
PID=5	coef	std err	z	P> z	[0.025	0.975]
logpopul	-0.0933	0.039	-2.371	0.018	-0.170	-0.016
selfLR	1.3470	0.117	11.494	0.000	1.117	1.577
age	-0.0179	0.008	-2.352	0.019	-0.033	-0.003
educ	0.2169	0.085	2.552	0.011	0.050	0.384
income	0.0810	0.023	3.524	0.000	0.036	0.126
const	-7.0605	0.844	-8.362	0.000	-8.715	-5.406
PID=6	coef	std err	z	P> z	[0.025	0.975]
logpopul	-0.1409	0.042	-3.343	0.001	-0.223	-0.058
selfLR	2.0701	0.143	14.435	0.000	1.789	2.351
age	-0.0094	0.008	-1.160	0.246	-0.025	0.007
educ	0.3219	0.091	3.534	0.000	0.143	0.500
income	0.1089	0.025	4.304	0.000	0.059	0.158
const	-12.1058	1.060	-11.421	0.000	-14.183	-10.028

- a) This is because we are keeping a variable fixed, but at different values, which gives different coefficients.
- b) No, there are no features that are insignificant across the board.
- c) age is only significant for $PID = 1$, $PID = 2$, and $PID = 5$.
- d) selfLR is significant for all the classes. logpopul, const, educ, income are significant for most the classes.

8 Support Vector Machine

a) Plot shown below

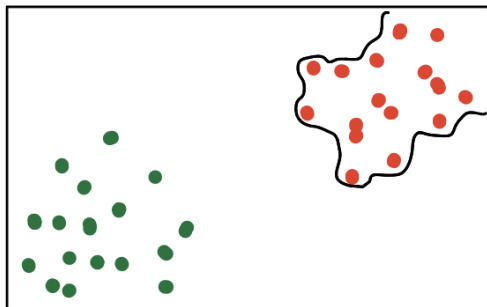


b) Points $(1, 0)$, $(1, 1)$, and $(2, 1)$ are support vectors. If $(2, 1)$ was removed, then the boundary would change.

c) A hard margin is used for linearly separable data, and does not allow for misclassifications. It tries to maximize the distance between data points and the boundary. Soft margin allows misclassification in hopes of achieving better generality, so it tries to minimize misclassification error. For this dataset, it doesn't matter whether use hard or soft margin, as they result in the same decision boundary.

d) The left sub-figure corresponds to SVM (linear), because the decision boundary is linear. The middle sub-figure corresponds to SVM with polynomial kernel, because the decision boundary follows a polynomial curve. Finally, the right sub-figure corresponds to SVM with RBF kernel with width (γ) equal to 1, because the RBF kernel because it has a radius around each of the classes.

8d)



8e)

