Applications of Logistic Regression in Protein Contact Map Prediction

Sritika Chakladar^{1, *, +}, Rahul Alapati^{1, +}, Vineet Nayak^{1, +}, Anuj Gupta^{1, +} and Austin Ream^{2, +}

¹Graduate Student, Computer Science and Software Engineering, Auburn University, 36830, USA ² Undergraduate Student, Computer Science and Software Engineering, Auburn University, 36830, USA

ABSTRACT

Prediction of protein contact map is of greater importance since it can improve the prediction of protein 3D structure. Protein contacts contain key information for the understanding of protein structure and function. This report lays out the implementation of Logistic Regression a non-linear regression algorithm, that has been used to train a binary logistic model to estimate the probability of contacts based on the position specific scoring matrices generated by PSI-BLAST. Two variants of gradient ascent namely: Batch & Stochastic and MCLE estimation have been used to learn the weight vectors during model training. Accuracies of top L/10, L/5 and L/2 predicted contacts are calculated to evaluate the classification performance of the model.

1. Introduction

To understand the mechanism of protein functions at the molecular level, it is usually necessary to determine their structures. However, determining 3D protein structure by experimental methods are high-cost, time-consuming and challenging. Fortunately, protein contact map or residue-residue contact map prediction simplifies the problem of protein structure prediction. The contacts in a protein sequence are predicted in the form a contact map, and then used in 3D protein structure prediction.

Contact Map of a protein, is a binary symmetric matrix and is a useful and simplified representation of the 3D protein structure. Contact Map contains important structural information of a protein, especially the backbone information and hence is crucial in predicting the 3D protein structure. The protein contact map prediction refers to the problem of predicting the spatial closeness of residue pairs in a folded protein structure. Contacts occurring between

^{*}szc0098@auburn.edu

these authors contributed equally to this work

sequentially distant residues, i.e. long-range contacts, impose strong constraints on the 3D structure of a protein and are particularly important for structural analysis, understanding the folding process and predicting the 3D structure. Even a small set of correctly predicted long-range contacts can be useful for improving ab initio structure prediction for proteins without known templates.

Logistic regression is a classification algorithm which is used to assign observations to a discrete set of classes. It transforms the input training dataset using a logistic function to return a probability which is mapped to discrete response classes based upon the prediction threshold. The learning algorithm uses regression coefficient to reduce the prediction error. It discovers the best value of the coefficient based on the prediction error and iteratively updates them until the stopping criteria is met.

Gradient ascent optimization is an iterative algorithm for finding the maximum of the cost function by taking steps proportional to the positive of the gradient at the current point. It has two variants: Batch and Stochastic Gradient Ascent. Batch gradient ascent calculates the gradient of the cost function by using the sum of the cost of each sample in every iteration. On the other hand, stochastic gradient ascent is a stochastic approximation of the gradient ascent optimization algorithm and it's an iterative method for maximizing the cost function over a single instance in the dataset. We use the gradient ascent optimization algorithm to update the regression coefficient vector to reach convergence quickly.

In this project, we train a binary logistic model to estimate the probability of contacts based on the position specific scoring matrices generated by PSI-BLAST. We implement both the variants of gradient ascent namely: Batch & Stochastic and MCLE estimation to learn the weight vectors during model training. We use these models to predict the residue-residue contacts in a protein sequence.

2. Methods

The protein sequences in Multi-FASTA format, the actual intra residue distances (d) from the rr file and the position specific scoring matrices have been used to train and test the binary logistic model. Gradient ascent optimization algorithms and MCLE Estimation have been used to learn the weight vector. The learned weight vector is used in residue-residue contact prediction in a protein sequence.

The input protein sequences and their respective actual intra residue distances (d) classes have been curated into non-overlapping sets of Training (75%) and

Test (25%) datasets using the simple random sampling without replacement.

Now, PSIBLAST and the nr database have been used to generate train and test PSSMs for protein sequences in train and test datasets, respectively.

BLAST (Basic Local Alignment Search Tool) is a sequence similarity search method, in which a query protein is compared to protein sequences in a target database to identify regions of local alignment and report those alignments that score above a given score threshold. Position-Specific Iterative (PSI)-BLAST is a protein sequence profile search method that builds off the alignments generated by a run of the BLASTp program.

A sample PSSM for a sequence is shown below in figure 1:

>sequence

TIKVLFVDDHEMVRIGISSYLSTQSDIEVVGEGASGKEAIAKAHELKPDLILMDLLMEDMDGVEATT QIKKDLPQIKVLMLTSFIEDKEVYRALDAGVDSYILKTTSAKDIADAVRKTSRGESVFEPEVLVKMR NRMKKRAELYEMLTEREMEILLLIAKGYSNQEIASASHITIKTVKTHVSNILSKLEVQDRTQAVIYAF QHNLIQ

	Α	R		N	D	C	Q	E	G	H	L	L	K	M	F	P	S	T	W	Y	V
1 T	-1	L	-3	-2	-3	-3	-2	-2	-3	-3	-3	-8	-	2 -3	3 -4	-3	3	6	-4	-3	
2 I	-3	3	-5	-6	-6	-3	-5	-5	-6	-6	7	(-	5 2	2 -2	-5	-5	-2	-5	-4	
3 K	-3	3	7	-2	-3	-5	0	-2	-4	-3	-5	-5	;	4 -4	4 -5	-4	-1	1	-5	-4	
4 V	-2	2	-5	-5	-6	-3	-5	-5	-6	-6	3	1	-	5 -:	L -3	-5	-4	-2	-5	-4	
5 L	-2	2	-5	-6	-6	-3	-4	-5	-6	-5	1		j -	5 3	3 -1	-5	-4	-3	-4	-3	
6 F	-3	3	-5	-6	-6	-3	-5	-5	-6	-5	5	4	-	5 -:	1 1	-5	-5	-3	-4	-3	
7 V	4	ļ	-4	-5	-5	1	-4	-4	-4	-5	0	-2	-	4 -2	2 -4	-4	-3	-2	-5	-4	
8 D	-4	ı	-4	-1	. 8	-6	-2	0	-4	-3	-6	-6	j -	3 -6	5 -6	-4	-2	-3	-7	-6	
9 D	-4	ı	-4	-1	. 8	-6	-3	-1	-4	-3	-6	-6	j -	3 -6	5 -6	-4	-2	-3	-7	-6	
10 H	-3	3	-2	-2	-3	-5	6	-1	-4	9	-5	-5	-	2 -3	3 -4	-4	-3	-3	-5	-1	
11 E	3	3	-1	-2	0	-4	1	4	-2	3	-4	-8	-	2 (-4	2	. 0	-1	-5	-4	
12 M	-3	3	-4	-5	-6	-4	-4	-5	-6	-5	2		- ا	4	7 -2	-5	-4	-3	-4	-3	
13 V	-2	2	-5	-5	-6	-3	-5	-5	-6	-5	2	1		5 :	1 0	-5	-4	-3	-5	-3	
14 R	-4	ı	8	-3	-4	-6	-1	-2	-5	-3	-5	-8	3	0 -4	4 -5	-5	-3	-3	-5	-4	
15 I	1	L	2	-3	-1	4	3	0	-2	-3	1)	1 3	3 -4	-4	-1	1	-5	-4	
16 G	-1		-5	-3	-4	-5	-4	-4	7	-4	-6	-6	j -	4 -5	5 -6	-4	-2	-4	-5	-5	
17 I	-4	ļ	-5	-6	-6	-4	-5	-5	-6	-5	1		- ا	5 () 6	-6	-5	-3	-3	-1	
18 S	1		5	-2	-4	-3	0	-1	-1	-3	-2	-8	3	2 -2	2 -5	-4	2	-1	-5	-4	
19 S	3	3	-4	-3	-3	-3	-2	-3	-1	-4	-1		-	3 5	5 1	-4	2	1	-4	-1	

Figure 1. Sample PSSM for a sequence.

The PSSM contains 20 values for each residue in a protein sequence. For a protein sequence of N residues, there will be N * 20 PSSM Values.

Now, for feature generation a sliding window of 5 around central residue is used, as shown in Figure 2.

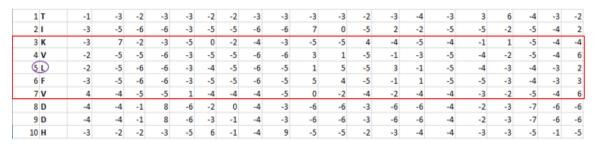


Figure 2. Sliding window of 5 around the central residue 'L'.

For the residue L, we consider the PSSM values of K & V above it and the PSSM values of F & V below it to generate a feature vector. Therefore, there will be 20 * 5 = 100 PSSM values for each non-terminal residue like L. For terminal residues like T that do not have one or more neighbors on either side, rows containing 20 values of -1 are used for feature generation.

The N * 20 PSSM values are now converted into N * 100 Feature vectors, as shown below in Figure 3.

20 PSSM values for residue 'L':



100 Feature vector for residue 'L':

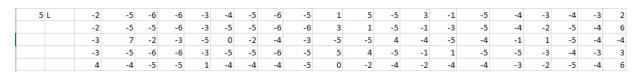


Figure 3. Feature generation from N * 20 PSSM Values.

Now, we combine the 100 PSSM values of each pair of residues (i, j), where j > i + 5, in a protein sequence i.e. for a protein sequence of length N, we will have N * 200 PSSM values.

The actual intra residue distance (d) of the two residues (i, j) from the corresponding .rr files has been used to generate the class labels as 1 if (i, j) is in contact i.e. d < 8, 0 otherwise. Also, all the residue-residue which are not in contact with a class label 0, have been considered in feature generation.

A sample training instance is shown in Figure 4. Each training instance consists of i, j, 200 PSSM values and the corresponding class label for the i, j residue pair.

A sample training instance:

```
1
33
(-0.0255238280445, -0.0255238280445, -0.0255238280445......-0.153142968267, -0.127619140223, -0.0255238280445)
1
```

Figure 4. Feature generation (i, j, 200 PSSM Values, class label).

2.1. Logistic Regression Learning on Training Set:

Initially, we implemented the Stochastic gradient ascent based optimization algorithm and MCLE estimation to learn the weight vector of Logistic Regression using the training set.

In Stochastic Gradient Ascent, we begin by selecting a single training instance for each iteration using simple random sampling with replacement and initially all the weights are set to zero. In every iteration, we calculate a prediction based on the logistic function using the training instance and the weight vectors as follows:

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

where, P(z) = Probability estimate z = input to the function ($w0 + \sum_i w_i x_i$)

The weight vectors are updated in every iteration using the MCLE estimation as follows:

$$w_i \leftarrow w_i + \eta X_i (Y - P(Z))$$

where,

 w_i : Weights η : learning rate

X_i: Features from the selected training instance

Y: Prediction

P(z) = Probability estimate Assumption : X_0 = 1 for W_0

The learning algorithm uses gradient ascent to discover best values of regression coefficients based on the error in prediction on training data in an iterative process, until convergence is reached. Given a dataset, we want to find the parameter vector 'w' which maximizes the likelihood.

During the training, we use a variable learning rate (η) which follows an exponential function and varies from 0.045 to 0.009405, i.e. the learning rate is high at higher gradient value and reduces as gradient decreases and we are closer to the maxima.

The above process is repeated until the error in our prediction doesn't change over a few iterations, i.e. change in the weights is less than $\epsilon = 0.0005$ for 10 iterations.

Also, we have implemented mini batch gradient ascent optimization algorithm in combination with MCLE estimation to learn the weight vector of Logistic Regression using the training set.

In mini batch gradient ascent, we begin by selecting a subset of the large training dataset, say, 500 training instances using random sampling with replacement and repeat the above mentioned process of learning weights.

2.2. Challenges faced in training:

1. **Normalization of dataset:** We observed that, the PSSM values for the 150 protein sequences were not normalized, due to which weights learned using gradient ascent optimization were too large and non-converging with the use of general learning rates and stopping criteria.

Hence, we normalized our features vectors using Euclidean distances. Using Euclidean distances, our feature vector values were normalized in the range of 0 and 1.

 Batch Gradient Ascent: The implementation of batch gradient ascent on the whole training dataset was computationally intensive and time taking. Hence, we used mini batch gradient ascent where we learn the weights on a randomly selected small subset of the training dataset.

2.3. Logistic Regression Classification on Test Set:

For classification, we use the weights learned as a part of training the binary logistic model and the logistic sigmoid function to classify the residue-residue pairs in the protein sequence as 1 if in contact, 0 otherwise.

The residue-residue pairs are classified as in contact (1) if the probability of the prediction is greater 0.5, not in contact (0) otherwise.

The prediction is calculated using the following equation:

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

where, P(z) = Probability estimate z = input to the function ($w0 + \sum_i w_i x_i$)

Accuracy:

The 'top' L/10, L/5, L/2 predicted contacts are used to evaluate the classification performance averaged over the proteins in the test dataset, where L is the length of the protein sequence and 'top' contacts are the contacts predicted with high probabilities.

3. Results

The following is a sample output showing the residue-residue contacts in a protein sequence in RR format. The contacts are sorted in the non-increasing their probabilities:

```
KTRWTREEDEKLKKLVEQNGTDDWKVIANYLPNRTDVQCQHRWQKVLNPE
8 17 0 8 0.510280046958
8 50 0 8 0.510017163785
9 17 0 8 0.509908357755
9 50 0 8 0.509645466791
7 17 0 8 0.509619764289
7 50 0 8 0.5093356867476
8 36 0 8 0.5093345703421
8 22 0 8 0.509339780959
8 48 0 8 0.509103081724
8 21 0 8 0.509068859846
9 36 0 8 0.508968064837
8 29 0 8 0.508968064837
8 29 0 8 0.508807248786
9 48 0 8 0.508679136455
7 36 0 8 0.5086685373931
7 22 0 8 0.508667136455
7 36 0 8 0.508687398746
9 21 0 8 0.508697136455
7 36 0 8 0.508887248786
9 31 0 8 0.5088873931
7 22 0 8 0.5088873931
7 22 0 8 0.5088873931
7 22 0 8 0.5088873931
8 29 0 8 0.50889834607826
7 48 0 8 0.508435518584
7 32 0 8 0.508435518584
7 32 0 8 0.508412321872
7 21 0 8 0.508438517369
10 17 0 8 0.508232047148
9 33 0 8 0.508162870742
7 29 0 8 0.508146894409
23 50 0 8 0.507969124657
8 20 0 8 0.507969124657
8 20 0 8 0.507969124657
8 20 0 8 0.50796974241431
8 45 0 8 0.507695762295
9 23 0 8 0.507590390664
8 14 0 8 0.507590390664
8 14 0 8 0.507590390664
```

The following are the accuracies of our logistic model trained using **Stochastic**Gradient Ascent:

Protein	L10	_	L5	_	L2		TP		FP	FN
1gzc.rr	82.6086956522	%	91.4893617021		93.2773109244			26535		726
1jos.rr	90.0	%	95.0	%	98.0	%	0	4296	0	169
5ptp.rr	95.4545454545	%	97.7272727273	%	98.1981981982	%	0	22783	0	653
1dbx.rr	93.333333333	%	93.333333333	%	97.3333333333		0	10053	0	387
1lo7.rr	92.8571428571	%	96.4285714286	%	98.5714285714	%	0	8758	0	287
1w0h.rr	94.7368421053	%	97.3684210526	%	98.9690721649	%	0	17363	0	403
1aba.rr	87.5	%	94.1176470588	%	97.6744186047			3181	0	140
1kw4.rr	83.333333333	%	92.3076923077	%	97.0588235294	%	0	1881	0	72
1h4x.rr	90.9090909091	%	90.9090909091	%	94.5454545455	%	0	5265	0	195
1i1j.rr	90.0	%	95.2380952381	%	98.1132075472	%	0	4795	0	255
1hfc.rr	93.333333333	%	96.7741935484	%	97.4358974359	%	0	11137	0	339
1tif.rr	85.7142857143	%	93.333333333	%	94.7368421053	%	0	2369	0	116
1g2r.rr	88.88888889	%	88.88888889		95.7446808511			3757		159
2mhr.rr	90.9090909091	%	95.652173913	%	96.6101694915		0	6179	0	149
1pch.rr	87.5	%	94.1176470588		95.4545454545			3218		185
1bdo.rr	75.0	%	87.5		90.0			2583		192
2cua.rr	91.6666666667	%	95.8333333333	%	98.3606557377			6473	0	313
1c52.rr	84.6153846154	%	92.3076923077		96.9230769231			7639		236
1aoe.rr	94.7368421053	%	97.3684210526		98.9583333333			16988		403
1ag6.rr	88.88888889	%	94.7368421053		97.9591836735			4134		237
1ku3.rr	83.333333333	%	91.6666666667		93.333333333			1483	0	57
1cjw.rr	81.25	%	84.8484848485		86.7469879518			12531	0	349
1a3a.rr	92.8571428571		96.5517241379		98.6111111111		0	9405	0	325
1mk0.rr	77.777777778		89.4736842105		91.666666667			4005	0	181
1ek0.rr		%	96.9696969697	%	98.8095238095			12807	0	396
1i4j.rr	90.9090909091		95.4545454545	%	98.1818181818			5230	0	230
1i5g.rr	92.8571428571		92.8571428571		95.8333333333			9291	0	300
1xdz.rr	95.652173913		97.8723404255		98.3193277311		0	26483	0	545
1dqg.rr	84.6153846154	%	88.4615384615	%	91.0447761194		0	7910	0	346
2arc.rr		%	96.875	%	98.75			11742	0	348
1fqt.rr		%	85.7142857143		90.7407407407			5075	0	281
1jbk.rr	94.444444444		97.2972972973		96.8085106383			16481	0	355
1a70.rr	88.88888889	%	94.7368421053	%	97.9166666667		0	3966	0	220
1i58.rr	94.444444444		97.2972972973		98.9361702128		0	16476	0	360
1g9o.rr	77.777777778		83.333333333		88.88888889		0	3472	0	183
1i1n.rr	95.4545454545		97.7272727273		99.1071428571			23301	0	570
1brf.rr		%	90.0		84.6153846154			1037	0	91
1atl.rr		%	97.5		99.0	%		18462	0	453
	89.0723292642 %									
Average $L/5 = 9$										
Average $L/2 = 9$										
	cy positive cont									
Average accurac	cy negetive cont	acts = 96	.9693035835 %							
						·				

We achieve an **accuracy of 0% in stochastic gradient ascent** because our model predicts all the probabilities as 0's. This is due to the fact that the weight bias W_0 , dominates all the other 200 weights trained by our model. The following are the values of the weights learned using stochastic gradient ascent:

14.862120852093954, -0.3068340424467362, -0.5645102324261682, -0.5919559154062827, -0.6359521475922265, -0.8748815468892239, -0.4930407319865049, -0.5192237281350535, -0.6559887144199089, -0.61351842392457 99, -0.5000423028157083, -0.5093203095213102, -0.47421024502506327, -0.5062396524322125, -0.7225858119307782, -0.8418707560299706, -0.4067223814493749, -0.3087036292066165, -1.1022047140919125, -0.703765894 6550015, -0.38782066586006514, -0.2673175014104722, -0.6081460744947806, -0.7260299514299328, -0.7701833988239132, -0.784040343754791, -0.5674813421965329, -0.6592163318764798, -0.782373449386169, -0.67560 08830884719, -0.38810283709551325, -0.49573283156700765, -0.5565408819552907, -0.4233849129404532, -0.6866249506606149, -0.942730478345319, -0.44457529304206433, -0.3395950090799228, -1.0289972891813013, -0 .6775480823792572, -0.2345973516195558, -0.2971783680596982, -0.8202607171059093, -0.8419641149305267, -0.9571630252278982, -0.7106034838502338, -0.7275304927400609, -0.8300606933011238, -0.8313607918335759 80421754, -0.6383827883754704, -0.14296756621174211, -0.29955405587812983, -0.635774700304924, -0.658293127991988, -0.7056795720337099, -0.8623305508327862, -0.5360030479949217, -0.5702834623994521, -0.8112 168939519353, -0.6992666127062951, -0.45951235989296424, -0.5644537034243267, -0.5729699179522403, -0.4934420229057796, -0.6840522383234774, -0.9300190289102966, -0.43464555588488873, -0.36374657173909397, -1.058444680522129, -0.658876917953763, -0.3209963641638313, -0.2725481201379408, -0.5370703675909478, -0.5364056080441187, -0.6381152991298793, -0.8911546351416636, -0.43835700520997056, -0.480397695014773 14, -0.6522588864940227, -0.579370385899744, -0.6270229994223799, -0.6553513221862953, -0.43171975400025075, -0.5527034201451749, -0.7612109075704906, -0.8715080557530678, -0.37045135973126847, -0.376450570 7895551, -1.0266825870434748, -0.6477946573351032, -0.5401732888750778, -0.22935157981328855, -0.46163552422059634, -0.52494073112021, -0.5186123352469275, -0.851482757921263, -0.4034664522620266, -0.374858 11248634715, -0.5689470643496338, -0.5524490695347273, -0.6026261077040717, -0.632293033114486, -0.36927585116671724, -0.5206724295636918, -0.819106086846733, -0.8601878114164373, -0.38846414899706544, -0.3 7008558029451216, -1.1052898489781595, -0.6782881331165902, -0.4230513615614961, -0.2655168652825049, -0.4549173250861444, -0.5328051732449073, -0.5595455418112547, -0.8933479746895162, -0.41824476310192416 -0.4393369603708005, -0.7394503204376078, -0.5726206122446678, -0.4878254376395328, -0.5412159530039465, -0.4012026308896511, -0.5091434567054046, -0.7006775692842923, -0.8051664131632604, -0.3900219575859059, -0.3275434512708515, -0.998391559878311, -0.5997998555174555, -0.3539545353238802, -0.2907433797572786, -0.7202704133415879, -0.8016494779228029, -0.9045629211139206, -0.7524811502714239, -0.624026804 4696639, -0.7398345845164525, -0.8517072885428768, -0.7428198322940265, -0.2470007821114372, -0.34659313103375006, -0.6760083630081659, -0.37164790146589183, -0.5746064563552234, -0.8757203789313087, -0.553 1760684414949, -0.42990009397329665, -1.074797296649172, -0.6326278010709231, -0.17466282511730544, -0.3222467172900929, -0.5455829994134067, -0.6244041197361697, -0.6801017685153092, -0.8533423867813058, . 4887570633403971, -0.5413205621046804, -0.8045174102910962, -0.5821368955369175, -0.4170819589112943, -0.4776618509231465, -0.48170448958841927, -0.4812423178598619, -0.649461295994092, -0.847928911841034 . -0.46195745159873736, -0.3412650266012254, -1.0258541407272794, -0.6127633521247882, -0.276505842989548, -0.24313758932377916, -0.4924094211931457, -0.4989806611584737, -0.49196347432096, -0.9664212764 498331, -0.39230155155475627, -0.4041978919394768, -0.6187359622471781, -0.601739510463208, -0.5972147387224545, -0.6516094742849202, -0.3838024658387346, -0.5649166860074479, -0.8627911391225757, -0.604882 2856379218, -0.3666018421248239, -0.3996698383437451, -1.1263752770859805, -0.726085039110758, -0.48267986513611094]

The value of W_0 is 14.06, whereas all the other weights are negative values.

The following are the accuracies of our logistic model trained using mini Batch Gradient Ascent:

rotein	L10		L5		L2		TP	TN	FI	,	FN
gzc.rr	0.0		0.0		0.0	%	123	18401		134	603
os.rr	0.0		0.0	%	0.0	%	17	2893		403	152
tp.rr	0.0		0.0	%	0.0	%	84	17166	5	517	569
lbx.rr	0.0	%	0.0		0.0	%	43	7340		713	344
lo7.rr	0.0	%	0.0	%	0.0	%	25	7009	1	749	262
ν0h.rr	0.0	%	0.0	%	1.03092783505	%	62	12179		184	341
aba.rr	0.0		0.0	%	0.0	%	11	2130	10	951	129
(w4.rr	0.0	%	0.0	%	0.0	%	10	1358		23	62
14х.гг	9.09090909091	%	4.54545454545		1.81818181818	%	28	3642	10	523	167
i1j.rr	0.0	%	0.0	%	0.0	%	59	2903	1	392	196
hfc.rr	0.0	%	0.0	%	1.28205128205	%	57	7751	3:	386	282
tif.rr	0.0	%	0.0	%	0.0	%	14	1291	10	978	102
27.77	0.0	%	0.0	%	0.0	%	28	2129	10	528	131
mhr.rr	9.09090909091	%	8.69565217391		3.38983050847	%	32	4020		159	117
och.rr	0.0	%	0.0	%	0.0	%	39	2114		104	146
odo.rr	0.0	%	0.0	%	0.0	%	31	1826		57	161
cua.rr	0.0	%	0.0	%	0.0	%	60	4308	2	165	253
52.гг	0.0	%	0.0	%	0.0	%	79	4190	3-	149	157
oe.rr	0.0	%	0.0	%	0.0	%	81	11403		585	322
ıg6.rr	0.0	%	10.5263157895	%	8.16326530612	%	43	2694	1	140	194
u3.rr	16.666666667	%	8.33333333333	%	6.6666666667	%	15	822	6	51	42
jw.rr	0.0	%	0.0	%	0.0	%	47	9160	3:	371	302
3a.rr	0.0	%	0.0	%	1.38888888889	%	46	5941	3-	164	279
nk0.rr	0.0	%	0.0	%	0.0	%	21	2804	1	201	160
ek0.rr	0.0	%	0.0	%	0.0	%	45	9123	3	584	351
L4j.rr	0.0	%	9.09090909091	%	7.27272727273	%	40	3508	1	722	190
i5g.rr	0.0	%	0.0	%	0.0	%	46	6073	3:	218	254
kdz.rr	0.0	%	0.0	%	0.0	%	79	19530	6	953	466
lgg.rr	0.0	%	0.0	%	0.0	%	83	4299	3	511	263
arc.rr	0.0	%	0.0	%	0.0	%	68	8445		297	280
qt.rr	0.0	%	0.0	%	1.85185185185	%	40	3754	1	321	241
jbk.rr	5.555555556	%	2.7027027027	%	1.06382978723	%	74	10201	6	280	281
70.гг	0.0	%	0.0	%	0.0	%	34	2735	1	231	186
L58.rr	0.0	%	0.0	%	0.0	%	81	10909	5	567	279
90.гг	0.0	%	0.0	%	0.0	%	29	2091	1	381	154
.1n.rr	0.0	%	0.0	%	0.892857142857	' %		92	16499	5802	478
rf.rr	0.0	%	0.0	%	0.0	%	15	762	2	75	76
tl.rr	0.0	%	0.0	%		%	44	14056		106	409
	0 = 1.06326422116 9	%									
	= 1.15511493778 %										
	= 0.916344167371	*									
	uracy positive con		= 1.61633159153 %								
	uracy negetive con										

We achieve an **accuracy of 1.06% in mini Batch gradient ascent**. This is due to the fact that the weight bias W_0 , dominates all the other 200 weights trained by our model. Due to the dominating W_0 value, our model is forced to predict more number of False Positives and less number of True Positives. The following are the values of the weights learned using stochastic gradient ascent:

2.0.00496227714559144, 0.00023081380806411265, 0.00023081398086411265, 0.00023081398097027235179193, 0.00010952053675379975, 0.0003399377492147533, 0.00019976305561127214, 0.0002000590641590747, 0.0002006433654922637, 0.0002308138359702825]
[-0.0004105138032970742, 2.5155995853066576-05, 0.0002296455335513812, 0.0002234525611637124, 0.000310929224467962, 0.0001634060723262653, 0.0001340985334390932, 0.00013713375695645252, 5.9467149377
5764556-05, 1.900261315817915-05, 0.0002471208617760445, 5.8232272377113356-05, 0.000169260807060667, 3.935695623900316-05, 4.909973375778516-05, 6.982776466852820-05, 4.801338823970666-05, 4.0001749277715946-05, 9.74803812795364-05, 9.74803812795364-05, 0.0002478780955727-06, 0.00024980860225567, 0.000476796091127715, 0.00040570953737, 0.0005720347783908, 0.00017923447783908, 0.00017918884444498304, 0.000128095373739, 0.00057277846-05, 0.0008777859460855211, 0.000100052634407579, 0.000477878749977, 0.00047778787879, 0.00047778789409578, 0.00087787879497878, 0.00087778789409578, 0.00087778789409578, 0.00087778789409578, 0.00087778789409578, 0.00087778789409578, 0.00087778789409578, 0.00087778789409578, 0.00087778789409578, 0.00087778789409578, 0.00087778789409578, 0.000877878789409578, 0.000877878789409578, 0.000877878789409578, 0.000877878789409578, 0.000877878789409578, 0.000877878789409578, 0.000877878789409578, 0.000877878789409578, 0.000877878789409578, 0.000877878789409578, 0.000877878789409578, 0.00087787894095789409578, 0.0008778789409578, 0.000877894095789409578, 0.00087787894095789409578, 0.0008778947894095789409578, 0.0008778945789409578, 0.0008778945789409

4. Challenges due to Class Imbalance

We observed that, the number of non-contact pairs in the training dataset largely outnumbered the number of contact pairs. Due to this class imbalance, initially our model was incorrectly trained to predict non-contact pairs.

Hence, we balanced our dataset with equal number of contact and non-contact pairs. We randomly selected non-contact pairs to equate them to the number of contact pairs.

Before class balancing, our overall prediction accuracy was 97% and both of our models (stochastic & batch) were only predicting non-contact pairs. Due to class balancing, both of our models improved, and the model trained using mini batch gradient ascent gained a contact accuracy of 1.06%.

5. Discussion

We performed different experiments as discussed above to improve the accuracy of our binary logistic model. Our contact accuracy was the 0 % and our model was only predicting non-contacts due to class imbalance. When the contact and non-contact pairs in our training dataset were balanced, our accuracy increased to 1.06 %. These experiments outlay the importance of class imbalance, feature normalization and mini batch gradients.

5. References

- 1. Wikipedia contributors. (2018, March 21). Protein structure. In Wikipedia, The Free Encyclopedia. Retrieved 19:23, March 26, 2018, from https://en.wikipedia.org/w/index.php?title=Protein structure&oldid=831674619
- 2. Pauling, L., Corey, R. B., & Branson, H. R. (1951). The structure of proteins: two hydrogen-bonded helical configurations of the polypeptide chain. Proceedings of the National Academy of Sciences, 37(4), 205-211.
- 3. Alberts, B., Johnson, A., Lewis, J., Raff, M., Roberts, K., & Walter, P. (2002). The shape and structure of proteins.
- 4. Anfinsen, C. B. (1973). Principles that govern the folding of protein chains. Science, 181(4096), 223-230.
- 5. Zhang, H. (2004). The optimality of naive Bayes. AA, 1(2), 3.
- 6. Schmidler, S. C., Liu, J. S., & Brutlag, D. L. (2000). Bayesian

- segmentation of protein secondary structure. Journal of computational biology, 7(1-2), 233-248.
- 7. Robles, V., Larrañaga, P., Peña, J. M., Menasalvas, E., Pérez, M. S., Herves, V., & Wasilewska, A. (2004). Bayesian network multi-classifiers for protein secondary structure prediction. Artificial Intelligence in Medicine, 31(2), 117-136.
- 8. Di Lena, P., Nagata, K., & Baldi, P. (2012). Deep architectures for protein contact map prediction. Bioinformatics, 28(19), 2449-2457.
- 9. Tress, M. L., & Valencia, A. (2010). Predicted residue—residue contacts can help the scoring of 3D models. Proteins: Structure, Function, and Bioinformatics, 78(8), 1980-1991.