1a) 1. What is the difference, between increasing ciacle by 2 and decreasing it triangle by 2?
ciacle by 2 and decreasing et triangle by 2?
Solution: There is no difference between increasing
Solution: Theore is no difference between increasing circle by 2 and decreasing triangle by 2.
When Ociale is increased by 2,
$p(solid circle) = \underbrace{e(\theta_{circle} + \theta_{solid}) \cdot e^2}_{e(\theta_{circle} + \theta_{solid}) \cdot e^2 + e(\theta_{circle} + \theta_{$
e(Beiscle+ Osolid).e2+ e(Ociscle+Ostriped).e2+
elletriangle + Asolid) + ellerriangle +  Astriped)
p(striped cir.
When Otriangle is decreased by 2
p(solid circle) = e (Ocircle + Osolid)  e (Ocircle + Osolid) + e (Ocircle + Ostriped) +
p (Ocircle + Osolid) + e (Ocircle + Ostriped) +
e (Otriangle + Osolid) e-2 + e (Otriangle + Ostripe)e
The state of the s

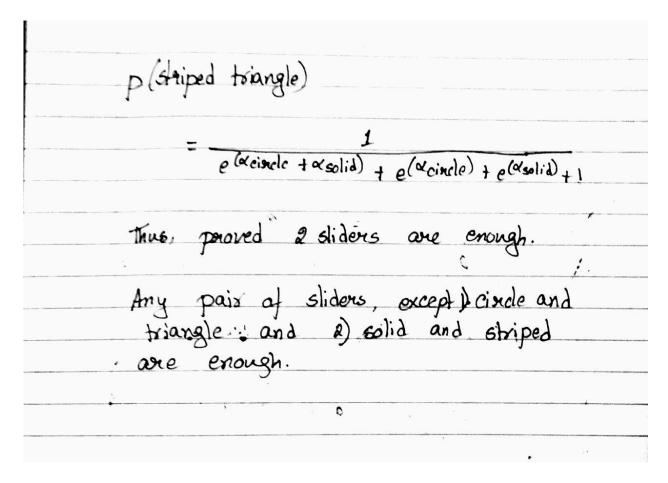
	$p(\text{solid circle}) = \frac{e(\text{Ocircle} + \text{Osolid}) \cdot e^2}{e(\text{Ocircle} + \text{Osolid}) \cdot e^2 + e(\text{Ocircle} + \text{Osolid})}$ $e(\text{Otriangle} + \text{Osolid}) + e(\text{Otriangle} + \text{Osolid})$
	which is same as when beiscle is increased by 2.
<b>2</b> .	The same holds for other probabilities  What happens if you increase both  circle and triangle by the same amount?
	Nothing happens if we increase both circle and triangle by the same arount.
	Let's assume we increase both ocircle and Ottiangle by 2,

then,  p(solid circle)' = e (Ocircle + Osolid). e <sup>2</sup> p(solid circle)' = e (Ocircle + Osolid). e <sup>2</sup> + e(Ocircle + Ostriped).e <sup>2</sup> e(Ocircle + Osolid). e <sup>2</sup> + e(Oricle + Ostriped).e <sup>2</sup> e(Ocircle + Osolid)  e(Ocircle + Osolid)  e(Ocircle + Osolid) + e(Oricle + Ostriped) +  e(Oricle + Osolid) + e(Oricle + Ostriped) +  e(Oricle + Osolid) + e(Oricle + Ostriped)  = p(solid circle)  Same holds for all probabilities  3. I Does it help at all to have 4 sliders,  on could you do just as well with 2 of them?  Which 2?  Solution: There is no need for 4 sliders. Everything can be done, with 2 sliders.		
p(solid circle)' = e (Ocircle + Osolid). e <sup>2</sup> e (Ocircle + Osolid). e <sup>2</sup> + e(Ocircle + Ostriped). e <sup>2</sup> e (Ocircle + Osolid). e <sup>2</sup> + e(Otriangle + Ostriped). e <sup>2</sup> e (Ocircle + Osolid)  e (Ocircle + Osolid)  e (Ocircle + Osolid)  e (Ocircle + Osolid) + e(Otriangle + Ostriped) +  e (Otriangle + Osolid) + e (Otriangle + Ostriped)  = p(solid circle)  Same holds for all probabilities  3. I Does it help at all to have 4 sliders,  on could you do just as well with 2 of them?  Which 2?  Solution: There is no need for 4 sliders. Everything  can be done with 2 sliders.		then,
e (θcircle + Osolid). e² + e(θcircle + Ostriped). e².  e(θtriangle + Osolid). e² + e(θtriangle + Ostriped). e²  e (θcircle + Osolid)  e (θcircle + Osolid)  e (θcircle + Osolid)  e (θtriangle + Osolid) + e(θtriangle + Ostriped)  = p(solid circle)  same holds for all probabilities  3. I Does it help at all to have 4 sliders,  on could you do just as well with 2 of them?  Which 2?  Solution: There is no need for 4 sliders. Everything can be done, with 2 sliders.		
e (Ocircle + Osolid)  e (Ocircle + Osolid) + e (Otriangle + Osolid) +  e (Otriangle + Osolid) + e (Otriangle + Osolid)  = p(solid circle)  Same hotels for all probabilities  3. It Does it help at all to have 4 sliders,  on could you do just as well with 2 of them?  Which 2?  Solution: There is no need for 4 sliders. Everything can be done with 2 sliders.		p/solid circle) = e (Ocircle + Osolid). e2
e (Ocircle + Osolid)  e (Ocircle + Osolid) + e (Otriangle + Osolid) +  e (Otriangle + Osolid) + e (Otriangle + Osolid)  = p(solid circle)  Same hotels for all probabilities  3. It Does it help at all to have 4 sliders,  on could you do just as well with 2 of them?  Which 2?  Solution: There is no need for 4 sliders. Everything can be done with 2 sliders.		e c(Ocircle + Osolid). e2 + p(Ocircle + Ostriped), p2
e (Ocircle + Osolid)  e (Ocircle + Osolid) + e (Otriangle + Osolid) +  e (Otriangle + Osolid) + e (Otriangle + Osolid)  = p(solid circle)  Same hotels for all probabilities  3. It Does it help at all to have 4 sliders,  on could you do just as well with 2 of them?  Which 2?  Solution: There is no need for 4 sliders. Everything can be done with 2 sliders.	619.5	e (Otriangle + Osolid). P2+ o (Otriangle+
e (Ocircle + Osolid)  e (Ocircle + Osolid) + e (Otriangle + Osolid) +  e (Otriangle + Osolid) + e (Otriangle + Osolid)  = p(solid circle)  Same hotels for all probabilities  3. It Does it help at all to have 4 sliders,  on could you do just as well with 2 of them?  Which 2?  Solution: There is no need for 4 sliders. Everything can be done with 2 sliders.		Oskijæd)ć
e (Ocircle + Osolid) + e(Ocircle + Oshiped) + e (Otriangle + Osolid) + e(Otriangle + Osolid) + e(Otriangle + Osolid)  = p(solid circle)  Same holds for all probabilities  3. If Does it help at all to have 4 sliders, on could you do just as well with 2 of them? Which 2?  Solution: There is no need for 4 sliders. Everything can be done with a sliders.		
Same he has for all probabilities  3. It Does it help at all to have 4 sliders, on could you do just as well with 2 of them? Which 2?  Solution: There is no need for 4 sliders. Everything can be done, with a sliders.		
Same he has for all probabilities  3. It Does it help at all to have 4 sliders, on could you do just as well with 2 of them? Which 2?  Solution: There is no need for 4 sliders. Everything can be done, with a sliders.		e (Otriangle + Osolid) + e (Otriangle + Osolid)
Same he has for all probabilities  3. It Does it help at all to have 4 sliders, on could you do just as well with 2 of them? Which 2?  Solution: There is no need for 4 sliders. Everything can be done, with a sliders.		The first that the following
Same holds for all probabilities  3. P Does it help at all to have 4 sliders, on could you do just as well with 2 of them? Which 2?  Solution: There is no need for 4 sliders. Everything can be done with 2 sliders.		= p(solid ciscle)
3. If Does it help at all to have 4 sliders, on could you do just as well with 2 of them? Which 2?  Solution: There is no need for 4 sliders. Everything can be done with 2 sliders.		<u> </u>
3. If Does it help at all to have 4 sliders, on could you do just as well with 2 of them? Which 2?  Solution: There is no need for 4 sliders. Everything can be done, with 2 sliders.		Same holds for all probabilities
Solution: There is no need for 4 sliders. Everything can be done with 2 sliders.		
Solution: There is no need for 4 sliders. Everything can be done with 2 sliders.	3.	I Does it help at all to have 4 sliders,
Solution: There is no need for 4 sliders. Everything can be done with 2 sliders.		on could you do just as well with 2 of them?
Solution: There is no need for 4 sliders. Everything can be done with 2 sliders.	100	Which 2?
	Solution:	There is no need for 4 sliders. Everything
		can be done with a sliders.
Marine Har Land Land Land Prince (1)		
	er, A.	or His course the same the same of parts (b)

p(solid cixele) p (Asolid + Orisela) - decircle + Osolid) + e (Ocircle + Ostriped) + e (Omangle + Osolid) + e (Otriangle + Ostripod) p(striped circle) elacirale + Ostriped) o (Ocircle + Osdid) + o (Ocircle + Ostriped) + c (Otriangle + Osta) + e (Otriangle + Ostriped) p (solid triangle) e (Otriangle + Osolid)
e (Otriangle + Osolid) + e (Otriangle 1 Ostriped) + e (Otriangle 1 Ostriped) + elevingle + Ostriped) · p(striped triangle) e (Otriangle + Oskiped) p (Ociacle + Osolid) + e (Ociacle + Ostriped) + e (Otriongle + Osolid) + elamangle + Ostro

	Dividing numerator and denominator by
	Dividing numerator and denominator by elotriangle + Ostriped),
1 - 11 - 11	p(solid circle)
House and the second se	e (Osolid + Ocircle - Otriangle - Ostriped)
	p (Osolid+Ocircle-Otriangle-Ostroped) + p (Ocircle-Otriangle)
The state of the s	e (Osolid + Ocircle - Otriangle - Ostriped) e (Osolid + Ocircle - Otriangle - Ostriped) + e (Ocircle - Otriangle) + e (Osolid - Ostriped) + 21
	the state of the s
	p (striped cincle)
the second designation with the stand	e (Ocircle - Otriangle)  - (Ocircle + Osolid - Otriangle - Ostriped) + e (Ocircle - Otriangle)
1	e (Ocircle + Osolid - Otriangle - Ostriped) + e (Ocircle - Otriangle)
	+ e (Oso)id - Oskriped) + 1
	p(solid triangle)
	e (Otriang e (Osolid - Ostriped)
	e(Ocivele +Osolid -Otriangle - Ostriped) + e(Ocisele-Otriangle
	+ e (Osolid - Ostriped) + 1.
No. of the last of	

	P(striped triangle)
	= e (Beircle + Osolid - Otriangle - Ostriped) + e (Ociscle - Otri
	t e (Osolid - Ostriped) + 1
	Let Ocincle - Oprionale be « cincle
	Let Ocincle - Opriangle be a cincle Osolid - Ostriped be asolid.
,,	p(solid cincle) e(xcincle + xsolid)
	e («cincle+«solid)+e («cincle)+e («solid)+1
	p(striped circle)
g /	$= \frac{e^{(\alpha_{cincle})}}{e^{(\alpha_{cincle})}}$
	e («circle+asolid) + e (acircle) + e (asolid) +1
	p(solid triangle) e(xsolid)
	= elacizade tasolid) + elacizade) + elasolid) +1



1b.

We know that the derivative of the log likelihood with respect to a parameter is simply the difference of the observed feature count and the expected feature count. Hence feature matching is always possible by making the observed feature count equal to the expected feature counts and calculating the feature weights. Note that this maximizes the log likelihood as well since we are making the derivative zero. Hence this implies that a good move in the log likelihood game implies a good move in the matching game and vice versa.

1c.

In the previous case, if we write down the equations of the probabilities of each of the shapes (solid circle, striped circle, solid triangle and striped triangle) and equate them to the observed count/ total shapes, then we can easily check that there are 4 variables and 5 equations (because of the constraint that sum of probabilities should be 1) and hence the equations can never be exactly satisfied. On the other hand, if we add one more feature, then we have 5 variables and 5 equations and hence we can get the exact shape count as required. Therefore, the added feature in lesson 4 is important.

No, this is not the only additional feature that could have been added. As a matter of fact, any combination of 2 features would have been sufficient for this purpose.

1d.

If no regularization is present, then the feature parameters can take any value to model the training data which can lead to overfitting because of the inherent noise in the data. We can see from the lesson, if no regularization is there then the values of the feature parameters is quite high. Once we use the l1 regularization, it induces sparseness. It can be observed that one of the feature values (hollow) became zero when we solved the optimization problem using l1 regularization. On the other hand, if we use l2 regularization it simply results in small values of the feature parameters for all of them.

1e.

As the number of tokens or training data (N) increases, the log likelihood will be more affected by it than the regularizer term.

$$F(\theta) = (\sum_{i=1}^{N} \log p(y_i | x_i)) - C.R(\theta)$$

Hence the model designed will try to match the empirical data. As more and more training data is added keeping the regularizer term constant, the model will tend to overfit as the regularizer's effect will continue to decrease. On the other hand, if the regularizer term is increased continuously, then it will dominate and most feature parameters will tend to 0. Hence if C is very large compared to the training data all feature parameters will be zero which is the limit and in that case the model will result in equal data for all shapes irrespective of the training data i.e. the training data will not matter anymore.

1f.

In Lesson 10 there are 9 parameters compared to the 6 parameters in lesson 9 and hence lesson 10 involves a more complicated model. For N=5 there is not much difference between the two models but for N=5000 and no regularization or with C=1, it turns out that lesson 10 can almost exactly model the training data whereas the model in lesson 9 fails to model the observed data for 2 of the shapes. Even though lesson 10 is modelling the training data exactly, there exists some smoothness because of the regularizer involved and as C is increased we can find a gradual decrease in the exact modelling of the training data. This also concludes the fact that as the complexity of the model is increased the training error can go to zero but we might over fit.

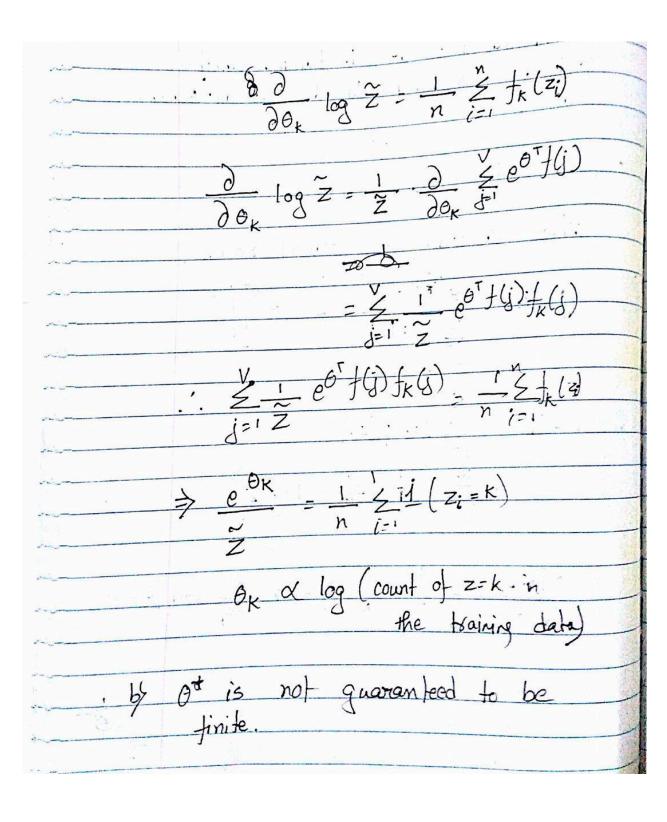
Adding features for contexts that has no relation with the training data outcomes does not change the model at all. The log likelihood or the objective function is completely insensitive to the feature weights of those contexts. This can be observed in lesson 14. Now, if we indeed want to use the context for our prediction, we need to use features that resemble bigram features and conjoin both the context and the outcome. Now the objective function is sensitive to these feature weights. Hence in a real-world scenario, we should use contexts that are related to the outcome or conjoin them to affect the model. Hence the context should be properly selected so that it does not become practically useless.

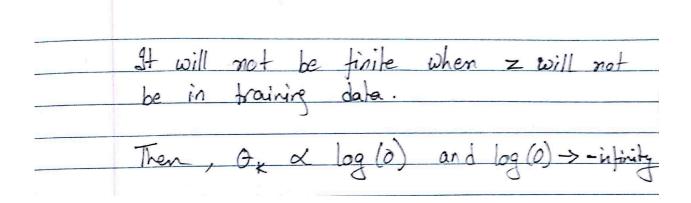
1h.

I am answering this question pertaining to lesson 16.

The model can fit the data moderately well. As it turns out among the bigram features same event and same shape has the highest weight and same fill has the smallest weight. Again, as the training data increases, the model can better fit the data because of lower variance and as the regularization coefficient is increased, the features start tending to zero. Hence, if we use 11 regularization, with C=10, it shows us that the bigram features are quite important compared to some other features like triangle & solid or square & striped which go to zero. With 12 regularization, similar thing happens except more number of features go close to zero other than going exactly to zero. As we increase C, the importance of the training data gradually decreases and the weight of all features gradually tend to zero.

20) Considering there are N training samples and taking the log-likelihood over all training samples,  $L = \prod_{i=1}^{N} P(Z_i)$ where, according to maxent model,  $A^T I(z)$  $p(z_i) = \frac{e^{\theta_i} f(z_i)}{\sum_{i=1}^{n} e^{\theta_i} f(i)} = \frac{e^{\theta_i} f(z_i)}{\sum_{i=1}^{n} e^{\theta_i} f(i)}$  $\log L = \sum_{i=1}^{N} \log p(z_i)$  $= \frac{1}{|z|} \left( \theta_{\overline{z}}^{T} + (z_{i}) - \log \tilde{z} \right)$  $\frac{\partial \log L}{\partial \theta_{K}} = \frac{1}{2} \int_{K} (z_{i}) - n \frac{\partial}{\partial \theta_{K}} \log \widetilde{Z} = 0$ 





3.

In this paper, Ronald Rosenfeld has described their approach in adaptive language modelling. They used a Maximum Entropy(ME) model, a selective unigram cache, a conditional bigram cache, and a conventional static trigram. They have used ARPA's official WSJ corpus for experiments and reported perplexity and word error rate results. To create the selective unigram cache, they stored all words that occurred in the history of the document. The motivation behind unigram cache is very interesting. Once a word occurs in a document, its probability of reoccurring is greatly improved. But this also depends on prior frequency of the word. The occurrence of a common word provides little new information whereas, the occurrence of a rare word provides a lot of information. Also, the occurrence of a more common word deviates less from the expectations of the static model, and hence requires very less modification to it. Similarly, bigram and trigram caches are created by storing consecutive word pairs and word triples respectively. The training data consisted of 38 million words of Wall Street Journal (WSJ) text from 1987-1989. To measure the impact of the amount of training data on language model adaptation, they experimented with systems based on varying amounts of training data. The adaptive language model was based on four component language models: (i)A conventional backoff trigram model (ii)A Maximum Entropy model (iii)A selective unigram cache (iv)A conditional bigram cache. All these four models were combined linearly with varying weights. They chose the specific weights by minimizing the perplexity of unseen data but later concluded that this did not always correspond with minimizing error rate. To evaluate error rate reduction, they used the Nov93 ARPA S1 evaluation set. They figured out that the main disadvantage of the Maximum Entropy framework is the computational requirements of training the ME model. But these are not severe for modest amounts of training data. The approach is thus particularly applicable in domains with a modest amount of LM training data. They have also demonstrated that the ME model significantly improves on the conventional static trigram.

I found the idea in this paper to be quite novel and interesting. Particularly, combining the four different component language models and determining the weights based on perplexity and

error reduction rate was an interesting approach. The paper is a significant contribution in the field of language modelling techniques.

4.

h=0.000001. The code is Q4.py.

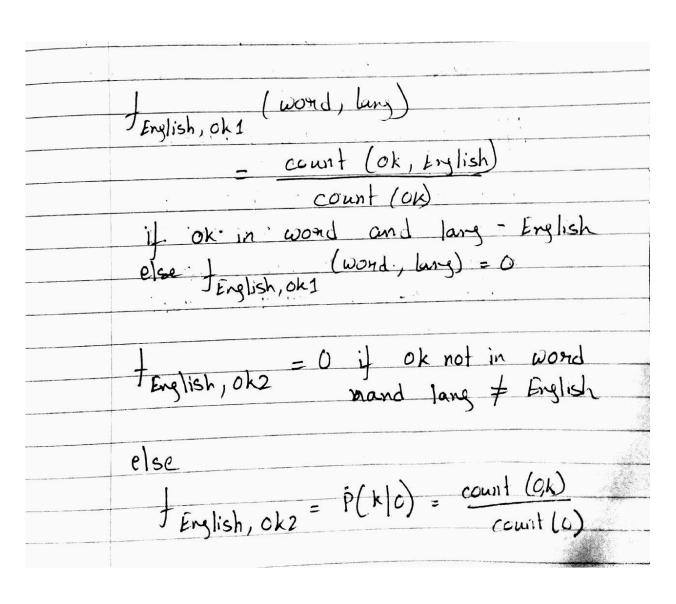
- a. The file produced is  ${\tt unigram\_maxent.csv.}$
- b. The file produced is bigram\_maxent.csv.
- 5. The code is Q5.py

5.	Model 1.	
	p(language   word) = count (word) count (w	nd, language pair) + \( \)
	where V = total number of wi	
	Here > is the only parame  the loglikelih	ler.
	I have optimized this using minimize. I have used the	method 'L-BFGS-B'
	and taken glo = less 1xe	7
	To compute the gradient, I package autograd.	have used the
	After prediction and evaluation micro recall = 0.6922462	1,
	micro precision = $0.6922462$ macro recall = $0.50069115$	
	macro precision = 0.60644593	
		-

	Model 2:
	P(lang   word) = p(word   lang) xp(lang)
	Let the word be okay and larg be English
. · · · · · · · · · · · · · · · · · · ·	p(okay   English) = poor p(& o   <bos), (english)="" td="" x<=""></bos),>
	p (k) o, < English) x p(a k, < English) x
-	P(yla, Linglish)
	p(k/o, (English)) = e + ((English), o, k)
	p(k/0, \English) = e of (\English), 0, k)  \[ \frac{\xeta}{\xeta} e^{\text{of}} \left(\xeta \text{English}\right), 0, \xeta}{\xeta} \]  \[ \xeta e^{\text{of}} \left(\xeta \text{English}\right), 0, \text{B}} \]  \[ \xeta e^{\text{of}} \left(\xeta \text{English}\right), 0, \text{B}} \]  \[ \xeta e^{\text{of}} \left(\xeta \text{English}\right), 0, \text{B}} \]
	character vocabulary consists of all characters in the training data.
	f ((English), o, k) is a 32301 dimensional vector.

Number of languages in the training file X bis character biggram vocabulary size = 32301
lang, biggram
Tlang, bigram (LEnglish), O, k) = count (LEnglish), O, k)  Flang, bigram = O plan  A is a 32301 dimensional vector
Hang, bigram = 0 whom  B is a 32301 dimensional vector  which we need to optimize.
I have optimized the loglikelihood using
I have optimized the collike ishood using scipy optimize minimize. I have used the method = 'L-BFGIS-B' and gtal = 1xe-3
To neduce the time taken for optimization,  I have nandomly selected 200 world-language
To neduce the time taken for optimize of have nandomly selected 200 world-language pairs from the development data and nun the optimization on them.

After prediction and evaluation.
micro siecall = 0.40691927
micro precision = 0.40691927 macro recall = 0.233294536
mauro precision = 0.2285/246186
<u>Model 3</u> :
p(lang   word) = etf (word, lang)  E etf (word, lang)  lang
f (word, larg) is a 64602 dimensional vector.
Number of languages in the training  file X character bigram vocabulary sixess = 64602
There are 2 features for every language biggram pair



Here, 0 is a 64602 dimensional vector which we need to optimize. Here, too I have used the same method to an to optimize and randomly selected 200 word language pairs from der data After prediction and Evaluation; miero recall = 0.570343725019 micro precision = 6.570343725019 macro recall = 0.3012610601479 macro precision = 6.42477854647