1.

a)

Training Set:
{'ZHO': 593, 'JPN': 557, 'KOR': 557, 'TEL': 533, 'ITA': 516, 'TUR': 504, 'ARA': 494, 'FRA': 473, 'SPA': 450, 'HIN': 352, 'DEU': 337}

| 'ZHO' | 'JPN' | 'KOR' | 'TEL' | 'ITA' | 'TUR' | 'ARA' | 'FRA' | 'SPA' | 'HIN' | 'DEU' |
|---|---|---|---|---|---|---|---|---|---|---|
| 593 | 557 | 557 | 533 | 516 | 504 | 494 | 473 | 450 | 352 | 337 |

Dev Set:
{'ZHO': 69, 'TEL': 62, 'KOR': 60, 'JPN': 60, 'TUR': 57, 'FRA': 53, 'ITA': 53, 'SPA': 52, 'ARA': 51, 'HIN': 47, 'DEU': 34}

| 'ZHO' | 'TEL' | 'KOR' | 'JPN' | 'TUR' | 'FRA' | 'ITA' | 'SPA' | 'ARA' | 'HIN' | 'DEU' |
|---|---|---|---|---|---|---|---|---|---|---|
| 69 | 62 | 60 | 60 | 57 | 53 | 53 | 52 | 51 | 47 | 34 |

b)
The majority class baseline accuracy on the dev set is 69/598=11.538%

2.
a) Yes. The training accuracy reaches 1 after 13 iterations.
b) According to the dev accuracy, it seems 13th iterations is the best. Although training accuracy reaches maximum value (1.0), the dev accuracy is also the highest (0.679). The test accuracy is 0.710

3.

| Feature Set | Test Set Accuracy | Number of Iterations to Converge | Number of Iteration Selected (Max Dev Accuracy) |
|---|---|---|---|
| Baseline (Unigram) | 0.71 | 13 | 13 |
| Baseline+Lemmatization +Normalization+Bigram (All Featrues) | 0.755 | 9 | 7 |
| All features except for Lemmatization | 0.753 | 11 | 6 |
| All features except for Normalization | 0.76 | 7 | 6 |
| All features except for Bigram | 0.657 | 17 | 15 |

# 4.

The most accurate model is the model with all features except for normalization.

a)

Confusion Matrix:

| | | Predicted Class | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | ARA | DEU | FRA | HIN | ITA | JPN | KOR | SPA | TEL | TUR | ZHO |
| Actual Class | ARA | 39 | 6 | 1 | 1 | 2 | 1 | 4 | 2 | 1 | 1 | 2 |
| | DEU | 0 | 37 | 0 | 0 | 2 | 1 | 0 | 0 | 0 | 0 | 1 |
| | FRA | 1 | 4 | 41 | 1 | 3 | 0 | 1 | 0 | 0 | 0 | 0 |
| | HIN | 1 | 1 | 0 | 19 | 0 | 0 | 0 | 0 | 6 | 2 | 1 |
| | ITA | 0 | 2 | 2 | 0 | 49 | 0 | 1 | 0 | 0 | 0 | 0 |
| | JPN | 0 | 1 | 1 | 1 | 0 | 41 | 10 | 0 | 1 | 0 | 7 |
| | KOR | 1 | 2 | 1 | 1 | 0 | 3 | 51 | 0 | 1 | 0 | 1 |
| | SPA | 1 | 1 | 3 | 2 | 5 | 0 | 3 | 42 | 1 | 2 | 1 |
| | TEL | 0 | 0 | 0 | 9 | 0 | 0 | 0 | 1 | 53 | 0 | 1 |
| | TUR | 1 | 5 | 1 | 0 | 1 | 0 | 3 | 1 | 2 | 40 | 1 |
| | ZHO | 1 | 1 | 0 | 1 | 0 | 4 | 7 | 1 | 2 | 1 | 47 |
| | Precision | 0.867 | 0.617 | 0.820 | 0.543 | 0.790 | 0.820 | 0.638 | 0.894 | 0.791 | 0.870 | 0.758 |
| | Recall | 0.650 | 0.902 | 0.804 | 0.633 | 0.907 | 0.661 | 0.836 | 0.689 | 0.828 | 0.727 | 0.723 |
| | F1 | 0.743 | 0.733 | 0.812 | 0.585 | 0.845 | 0.732 | 0.723 | 0.778 | 0.809 | 0.792 | 0.740 |

b)

Language: ARA
bias: -2
10 most common features and their weights:
[('alot', 30), ('statment', 27), ('many|reason', 26), ('any', 20), ('self', 20), ('alot|of', 20), ('his', 18), ('and|they', 18), ('thier', 18), ('fun', 18)]
10 least common features and their weights:
[('of|the', -24), ('difficult', -19), ('go', -18), ('.|Because', -18), ('to|study', -17), ('only', -17), ('everything', -17), ('Because', -17), ('decide', -16), ('whether', -16)]
Language: DEU
bias: -13
10 most common features and their weights:
[(',|that', 23), ('often', 23), ('special', 20), ('more|important', 19), ('beeing', 19), ('statement', 18), (',|because', 17), ('important|to', 17), ('big', 16), ('a|high', 15)]
10 least common features and their weights:
[(',|and', -30), ('"s', -23), ('i', -21), (';', -17), ('two', -17), ('we', -17), ('we|can', -16), ('be|that', -15), ('to|the', -15), ('our|life', -15)]
Language: FRA
bias: -9
10 most common features and their weights:
[('...', 34), ('Indeed', 29), ('.|Indeed', 26), (',|we', 23), ('be|to', 23), ('totally', 20), ('even|if', 19), ('Indeed|,', 19), ('why', 18), ('stay', 18)]
10 least common features and their weights:
[('the|statement', -23), ('the|people', -21), (',|but', -20), ('reason', -19), ('probably', -18), ('there', -18), ('get', -18), ('from', -18), ('see|the', -17), ('this|be', -17)]
Language: HIN
bias: -8
10 most common features and their weights:
[('then', 24), ('various', 20), ('which', 19), ('Now', 19), ('go|for', 19), ('behind', 18), ('.|Now', 17), ('old|age', 17), ('help', 16), ('today', 16)]
10 least common features and their weights:
[('Because', -18), ('time|.', -17), ('that|be', -16), (',|but', -16), ('for|a', -16), ('For|example', -16), (':', -16), ('may', -16), ('also|the', -15), (',|I', -15)]
Language: ITA
bias: -7
10 most common features and their weights:
[(':', 27), ('think|that', 27), (',|for', 24), ('probably', 20), ('a|specific', 20), ('happen', 18), ('that|in', 18), ('study', 17), ('at|the', 17), ('that|a', 17)]
10 least common features and their weights:
[('get', -22), ('may', -20), ('do|not', -18), (',|I', -17), ('hard', -17), ('technology', -16), ('country', -15), ('those', -15), ('.|Because', -15), ('which', -15)]
Language: JPN
bias: -4
10 most common features and their weights:
[('Japan', 39), ('<s>|I', 33), ('in|Japan', 33), ('Japan|,', 27), ('reason|.', 23), ('Japanese', 20), (',|because', 18), ('how', 18), ('two|reason', 18), ('fact|,', 18)]
10 least common features and their weights:
[('he', -24), ('a|good', -23), ('an', -22), ('be|a', -20), ('every', -20), ('present', -20), ('a|very', -19), ('Korea', -19), ('buy', -19), (',|that', -19)]
Language: KOR

bias: -6
10 most common features and their weights:
[('Korea', 33), ('their|own', 22), ('in|Korea', 22), ('.|However', 19), ('can|be', 18), ('have|their', 18), ('fail', 18), ('other|people', 17), ('enjoy|their', 17), ('people|want', 17)]
10 least common features and their weights:
[('take', -26), ('may', -24), ('this', -22), ('always', -21), ('how', -21), ('u', -20), ('take|the', -19), ('the|people', -19), ('reach', -18), ('Japan', -18)]
Language: SPA
bias: -9
10 most common features and their weights:
[(',|etc', 27), ('of|the', 22), ('be|go', 20), (',|be', 20), ('moment', 20), ('have|be', 19), ('learn', 18), ('that|you', 17), ('maybe', 17), ('Many', 16)]
10 least common features and their weights:
[('which', -28), ('from', -28), ('on', -21), ('people|.', -20), ('people|to', -19), ('be|also', -18), ('i', -18), ('after', -18), ('he', -17), ('will', -17)]
Language: TEL
bias: -4
10 most common features and their weights:
[('statement', 22), ('may', 21), ('strongly', 21), ('with|out', 19), ('and|also', 17), ('the|above', 17), ('present', 16), ('by', 16), ('some', 16), ('when', 16)]
10 least common features and their weights:
[('if', -23), ('just', -22), ('I|think', -22), ('what', -19), ('often', -18), (',|a', -18), (':', -18), ('keep', -17), ('to|be', -17), ('you', -17)]
Language: TUR
bias: -6
10 most common features and their weights:
[('Turkey', 25), ('Because', 22), ('As|a', 22), ('condition', 22), ('.|Because', 21), ('idea', 20), ('not|be', 18), ('anything', 18), ('succesful', 18), ('start|to', 17)]
10 least common features and their weights:
[(',|and', -34), (',|but', -22), ('could', -21), ('out', -20), ('a|lot', -20), ('in|the', -20), ('might', -19), ('feel', -19), ('First|,', -19), ('because|they', -18)]
Language: ZHO
bias: -15
10 most common features and their weights:
[('still', 32), ('tell', 26), ('just', 23), ('people|to', 23), ('Take', 23), ('.|Take', 23), (',|the', 21), ('not|only', 21), ('may', 20), ('always', 20)]
10 least common features and their weights:
[('have|to', -33), ('an', -24), ('his', -20), ('but', -20), ('know|about', -19), ('change', -19), ('explain', -18), ('clear', -18), ('`', -18), ('try|to', -18)]

c)

Precision, Recall and F1 for each language are as shown in part a) confusion matrix

d)

The language with the highest precision is **SPA**, the lowest precision is **HIN**.

The language with the highest recall is **DEU**, the lowest recall is **HIN**.

The language with the highest F1 is **ITA**, the lowest F1 is **HIN**.

According to F1 score:

Overall, the model performs best on language **ITA**, **FRA** and **TEL**.

The model performs worst on language **HIN**, **KOR** and **JAP**, and especially on **HIN**.

It seems this model cannot distinguish between **HIN** and **TEL**. It is very likely these two languages are natively similar so that their native speakers have similar pattern when writing in English. I briefly searched on the Internet, and it seems both **Hindi** and **Telegu** are used in India.

Besides, it is also interesting to see that all three languages with worst performance are used in **Asian countries** (India, Korea and

Japan).