

## Inferring Gender From the Context of Tweets: A Region Specific Example

Uses a subset of users from Nigeria, contains both a machine and non-machine learning approach

Determines gender from tweet content alone.

### **Data Collection and Preparation**

Data was collected using Twitter's Search API with Geocode.

45 Nigerian cities with population > 100000 with 40 miles radius

Name gender pairs were obtained by searching name on Facebook, which has a gender field. This provided ground truth for twitter users gender

Any Non English names were translated

Left with 11155 users with 18.5 million tweets

Emoticons and symbols were translated to their meaning

### **Feature Selection and Machine Learning Experiments**

Features were generated for unigrams, hashtags and psychometric properties.

Linguistic Inquiry and Word Count was performed on each tweet. This maps certain high frequency words to psychometric properties.

Unigrams features included all words but stop words

Hashtags contained the 75% most common ones

SVMLight was used with various combinations of these feature vectors

SVMLight is an SVM classifier with a linear kernel

80/20 training to test ratio

<b>Feature Class</b>	<b>Prec.</b>	<b>Rec.</b>	<b>F</b>	<b>Accu.</b>
<b>Hashtag</b>	70.82	46.99	56.48	63.81
<b>Hashtag, LIWC</b>	71.51	47.17	56.83	64.18
<b>Hashtag, LIWC, Unigram</b>	<b>82.53</b>	<b>77.64</b>	<b>80.00</b>	<b>80.60</b>
<b>Hashtag, Unigram</b>	82.06	77.46	79.69	80.26
<b>LIWC</b>	70.76	75.07	72.85	72.02
<b>LIWC, Unigram</b>	82.22	77.79	79.93	80.48
<b>Unigram</b>	<b>82.50</b>	<b>77.47</b>	<b>79.90</b>	<b>80.51</b>

*Table 2 – Classification Results*