



In [37]:	<pre>#new vectorizer + subreddit ppl1= Pipeline([</pre>
In [38]:	<pre> 1)), 1)), ('logreg', LogisticRegression (C=10000.0, max_iter=100)) 1) ppll.fit(train_data, train_labels_enc) pred1 = ppll.predict(validation_data) pl, r1,accl, f11 = evaluation_summary_q2q3("Body + Subreddit", pred1, validation_labels_enc, flag=1) Body + Subreddit ' has Acc=0.739 P=0.506 R=0.552 F1=0.524 SINGLE FEATURE: Majority_type #new vectorizer + majority_type #new vectorizer + subreddit</pre>
In [39]:	<pre> 1)),</pre>
	<pre>#new vectorizer + subreddit ppl3= Pipeline([</pre>
In [39]:	('logreg', LogisticRegression (C=100000.0, max_iter=100))]) ppl3.fit(train_data,train_labels_enc) pred3 = ppl3.predict(validation_data) p3, r3,acc3, f13 = evaluation_summary_q2q3("Body + Title",pred3,validation_labels_enc,flag=1) 'Body + Title' has Acc=0.731 P=0.429 R=0.494 F1=0.454 Create new features with Vader
In [41]:	<pre>vader = SentimentIntensityAnalyzer() vs_train = [] for i in range(len(train_data['body'])): scores = vader.polarity_scores(train_data['body'][i]) vs_train.append(scores['compound']) vs_valid = [] for i in range(len(validation_data['body'])): scores = vader.polarity_scores(validation_data['body'][i])</pre>
	<pre>vs_valid.append(scores['compound']) vs_test = [] valid in range(len(test_data['body'])): scores = vader.polarity_scores(test_data['body'][i]) vs_test.append(scores['compound']) train_data_new = train_data.copy() valid_data_new = validation_data.copy() test_data_new = test_data.copy() train_data_new['vader'] = np.zeros_like(len(vs_train)) valid_data_new['vader'] = np.zeros_like(len(vs_valid))</pre>
In [42]:	<pre>test_data_new['vader'] = np.zeros_like(len(vs_test)) vs_train_df = pd.DataFrame(vs_train) vs_valid_df = pd.DataFrame(vs_valid) vs_test_df = pd.DataFrame(vs_test) #features have values in range [-1,1], thus i split the range in 5 parts (For all sets) train_data_new['vader'][(vs_train_df[0]<=(-0.6))] =0 train_data_new['vader'][((vs_train_df[0]>(-0.6)) & (vs_train_df[0]<=(-0.2)))] = 0 train_data_new['vader'][((vs_train_df[0]>(-0.2)) & (vs_train_df[0]<=(0.2)))] = 0 train_data_new['vader'][((vs_train_df[0]>(0.2)) & (vs_train_df[0]<=(0.6)))] = 0 train_data_new['vader'][((vs_train_df[0]>(0.2)) & (vs_train_df[0]<=(0.6)))] = 0 train_data_new['vader'][(vs_train_df[0]>(0.2)) & (vs_train_df[0]<=(0.6)))] = 0</pre>
	<pre>test_data_new['vader'][(vs_test_df[0]<=(-0.6))] =0 test_data_new['vader'][((vs_test_df[0]>(-0.6)) & (vs_test_df[0]<=(-0.2)))] = 0 test_data_new['vader'][((vs_test_df[0]>(-0.2)) & (vs_test_df[0]<=(0.2)))] = 0 test_data_new['vader'][((vs_test_df[0]>(0.2)) & (vs_test_df[0]<=(0.6)))] = 0 test_data_new['vader'][(vs_test_df[0]>(0.2)) & (vs_test_df[0]<=(0.6)))] = 0 valid_data_new['vader'][(vs_valid_df[0]<=(-0.6))] = 0 valid_data_new['vader'][((vs_valid_df[0]>(-0.2)) & (vs_valid_df[0]<=(-0.2)))] = 0 valid_data_new['vader'][((vs_valid_df[0]>(-0.2)) & (vs_valid_df[0]<=(0.2)))] = 0 valid_data_new['vader'][((vs_valid_df[0]>(0.2)) & (vs_valid_df[0]<=(0.2)))] = 0 valid_data_new['vader'][((vs_valid_df[0]>(0.2)) & (vs_valid_df[0]<=(0.2)))] = 0 valid_data_new['vader'][(vs_valid_df[0]>(0.2)) & (vs_valid_df[0]<=(0.2)))] = 0</pre>
In [43]:	<pre>ppl4= Pipeline([</pre>
	<pre>('selector', NumericalItemSelector(key='vader')),</pre>
In [44]:	<pre>#measure the len of each post, get statistics of it to split it (repeated for all sets) train_body_length = train_data_new['body'].apply(len) train_data_new['body_length'] = train_body_length train_descr= train_data_new['body_length'].describe() valid_body_length = valid_data_new['body'].apply(len) valid_data_new['body_length'] = valid_body_length valid_descr = valid_data_new['body_length'].describe() test_body_length = test_data_new['body'].apply(len) test_data_new['body_length'] = test_body_length</pre>
Out[44]:	test_descr = test_data_new['body_length'].describe() columns = ["Train Body Length", "Test Body Length", "Valid Body Length"] descr = pd.concat([train_descr ,valid_descr,test_descr],axis=1) descr.columns = columns descr Train Body Length
In [45]:	25% 42.250000 44.000000 103.000000 50% 101.000000 103.000000 101.000000 75% 227.000000 222.000000 230.000000 max 11747.000000 10605.000000 5031.000000 #split all sets based on individual the statistics showabove at 4 parts & assign categorize them (0-1-2-3) train_data_new['body_length'][train_data_new.body_length<=10] = 0 train_data_new['body_length'][((train_data_new.body_length>10)) & (train_data_new.body_length>10))] = 1
	<pre>train_data_new['body_length'][((train_data_new.body_length>101) & (train_data_new.body_length<= 227))] = 2 train_data_new['body_length'][train_data_new.body_length>227] = 3 train_data_new['body_length'] = ((np.array(train_data_new['body_length'])).reshape(-1,1)) valid_data_new['body_length'][valid_data_new.body_length<= 24] = 0 valid_data_new['body_length'][((valid_data_new.body_length>44) & (valid_data_new.body_length'= 203))] = 1 valid_data_new['body_length'][((valid_data_new.body_length>103) & (valid_data_new.body_length'= 223))] = 2 valid_data_new['body_length'][valid_data_new.body_length>223] = 3</pre>
	<pre>valid_data_new['body_length'] = (np.array(valid_data_new['body_length'])).reshape(-1,1) test_data_new['body_length'][test_data_new.body_length<='3] = 0 test_data_new['body_length'][((test_data_new.body_length>3) & (test_data_new.body_length<='01))] = 1 test_data_new['body_length'][((test_data_new.body_length>101) & (test_data_new.body_length<='010)] = 2 test_data_new['body_length'][test_data_new.body_length>30] = 3 test_data_new['body_length'] = (np.array(test_data_new['body_length'])).reshape(-1,1)</pre> SINGLE FEATURE: Body_length
In [46]:	<pre>#new vectorizer + body_length ppl5= Pipeline([</pre>
In [47]:	<pre>('logreg', LogisticRegression(C=10000.0, max_iter=100)) ppl5.fit(train_data_new, train_labels_enc) pred5 = ppl5.predict(valid_data_new) p5, r5,acc5, f15 = evaluation_summary_q2q3("Body + Body_length", pred5, validation_labels_enc, flag=1) 'Body + Body_length' has Acc=0.751 F=0.542 R=0.613 F1=0.572 #vec+ is_first_post ppl6= Pipeline([</pre>
	<pre>transformer_list=[('body', Pipeline([</pre>
In [48]:	<pre>pred6 = ppl6.predict(validation_data) p6, r6,acc6, f16 = evaluation_summary_q2q3("Body + is_first_post",pred6,validation_labels_enc,flag=1) 'Body + is_first_post' has Acc=0.747 P=0.537 R=0.588 F1=0.557 #vec+ author ppl7= Pipeline([</pre>
	<pre>('CV', CountVectorizer(tokenizer = tokenize_normalize,max_features='0000)),</pre>
In [49]:	<pre>'Body + author' has Acc=0.752 P=0.516 R=0.582 F1=0.542 pp=np.array((p_vo,p1,p2,p3,p4,p5,p6,p7)) rr=np.array((r_vo,r1,r2,r3,r4,r5,r6,p7)) accacc=np.array((acc_vo,acc1,acc2,acc3,acc4,acc5,acc6,acc7)) f1f1=np.array((f1_vo,f11,f12,f13,f14,f15,f16,f17)) columns=['Precision(macro)','Recall(macro)','F1(macro)','Accuracy'] index = ['Baseline: body','Add: subreddit', 'Add: majority_type','Add: title',</pre>
Out[49]:	print('\tcombination of Body with one Extra Feature , And New Vectorizer') COMBINATION OF BODY WITH ONE EXTRA FEATURE , AND New Vectorizer Precision(macro) Recall(macro) F1(macro) Accuracy Baseline: body 0.540 0.589 0.561 0.744 Add: subreddit 0.506 0.552 0.524 0.739 Add: majority_type 0.522 0.559 0.537 0.746 Add: title 0.429 0.494 0.454 0.731 Add:body_length 0.521 0.563 0.539 0.758 Add:vader 0.542 0.613 0.572 0.751 Add: is_first_post 0.537 0.588 0.557 0.747
	Add: Is_Tirst_post
	<pre>('selector', ItemSelector(key='body')), ('CV', CountVectorizer(tokenizer = tokenize_normalize,max_features=18888)),</pre>
In [51]:	<pre>('logreg', LogisticRegression (C=10000.0, max_iter=100))]) pipe1.fit(train_data_new, train_labels_enc) predp1 = pipe1.predict(valid_data_new) pln, rln,accln, fln = evaluation_summary_q2q3("Body + body_length+ Subreddit",predp1,validation_labels_enc,flag=1) 'Body + body_length+ Subreddit' has Acc=0.744 P=0.532 R=0.567 F1=0.547 #new vectorizer + body_length + majority_type pipe2= Pipeline([</pre>
	<pre>('body', Pipeline([</pre>
In [52]:	<pre>('logreg', LogisticRegression(C=.0000.0, max_iter=.00))]) pipe2.fit(train_data_new, train_labels_enc) predp2 = pipe2.predict(valid_data_new) p2n, r2n,acc2n, f2n = evaluation_summary_q2q3("Body + body_length+ Majority_type",predp2,validation_labels_enc,flag=1) 'Body + body_length+ Majority_type' has Acc=0.746 P=0.526 R=0.572 F1=0.546 #new vectorizer + vader + body_length pipe3= Pipeline([</pre>
	<pre>transformer_list=[('body', Pipeline([</pre>
In [53]:	<pre>pipe4= Pipeline([</pre>
	<pre>('body', Pipeline([</pre>
In [54]:	<pre>pipe4.fit(train_data_new, train_labels_enc) predp4 = pipe4.predict(valid_data_new) p4n, r4n,acc4n, f4n = evaluation_summary_q2q3("Body + body_length+ is_first_post",predp4,validation_labels_enc,flag=1) 'Body + body_length+ is_first_post' has Acc=0.752 P=0.552 R=0.607 F1=0.576 #new vectorizer + vader + author pipe5= Pipeline([</pre>
	<pre>('body', Pipeline([</pre>
In [55]:	<pre>('logreg', LogisticRegression (C=10000.0, max_iter=100))]) pipe5.fit(train_data_new, train_labels_enc) predp5 = pipe5.predict(valid_data_new) p5n, r5n,acc5n, f5n = evaluation_summary_q2q3("Body + body_length+ author",predp5,validation_labels_enc,flag=_) 'Body + body_length+ author' has Acc=0.750 P=0.546 R=0.577 F1=0.559 ppl=np.array((p_vo,pln,p2n,p3n,p4n,p5n)) rrl=np.array((r_vo,rln,r2n,r3n,r4n,r5n)) accaccl=np.array((acc_vo,accln,acc2n,acc3n,acc4n,acc5n)) flf1l=np.array((fl vo,fln,f2n,f3n,f4n,f5n))</pre>
Out[55]:	columns=['Precision(macro)','Recall(macro)','Fl(macro)','Accuracy'] index = ['Baseline: body','Add: body_length & subreddit', 'Add: body_length & majority_type',
	Add: body_length & subreddit 0.532 0.567 0.547 0.744 Add: body_length & majority_type 0.526 0.572 0.546 0.746 Add:vader & body_length 0.514 0.546 0.528 0.760 Add: body_length & is_first_post 0.552 0.607 0.576 0.752 Add:body_length & author 0.546 0.577 0.559 0.750 Final Model, evaluation on Test set & confusion matrix Occam's razor "rule", the simplest model was choosen among the 2 that peformed better than baseline, as their performance difference was not that huge to choose the more complex model
In [56]:	<pre>best_model1 = Pipeline([('union', FeatureUnion(</pre>
In [57]:	<pre>('logreg', LogisticRegression (C=10000, max_iter=100))]) best_model1.fit(train_data_new, train_labels_enc) preds_best_model1 = best_model1.predict(valid_data_new) p_bm1, r_bm1,acc_bm1, f1_bm1 = evaluation_summary_q2q3("Body+body_length, & new vectorisation: ",preds_best_model1,validation_labels_enc,flag=1)</pre> 'Body+body_length, & new vectorisation: ' has Acc=0.751 P=0.542 R=0.613 F1=0.572
In [58]:	<pre>'Body+body_length, & new vectorisation: ' has Acc=0.757 P=0.532 R=0.571 F1=0.549 p_final=np.array((prec_opt,p_bm_test1)) r_final=np.array((rec_opt,r_bm_test1)) acc_final=np.array((acc_opt,acc_bm_test1)) f_final=np.array((f1_opt,f1_bm_test1)) columns=['Precision(macro)','Recall(macro)','F1(macro)','Accuracy'] index = ['Baseline-Optimised on Q2',</pre>
Out[58]: In [59]:	<pre>print("\t\t\tTEST SET COMPARISON") df_final.style.highlight_max(color='red') TEST SET COMPARISON Precision(macro) Recall(macro) F1(macro) Accuracy Baseline-Optimised on Q2</pre>
Out[59]:	negative neutral positive very negative very positive negative 88 172 14 6 2 neutral 102 2116 270 9 17 positive 6 282 791 1 22 very negative 6 12 1 13 0 very positive 0 11 42 0 33