# Text as Data Final Paper

*Jacqueline Gutman, Alex Pine*

*May 10, 2016*

```r
# LDA model for debate: debate_LDA_15
# topic names for debate: debate_LDA_15_names
# data frame for twitter: twitter.df
# dfm for twitter: twitter_dfm
# posterior topic distribution (LDA) = @gamma
# LDA model for twitter: use simple_lda_20, simple_lda_15, simple_lda_10
# LDA posterior for twitter using debate topics: twitter.topics$topics
all(nrow(twitter_dfm) == sum(twitter.df$debate_topic != 0),
    nrow(twitter_dfm)  == nrow(twitter.topics$topics))
```

```
    [1] TRUE
```

```r
table(twitter.df$debate_topic)
```

```
       0     1     2     3     4     5     6     7     8     9    10    11    12    13    14
    2940   525   937   430   533  1387   618   490   617   622   673   406  2340   480   415
      15
     456
```

```r
debate_LDA_15_names
```

```
     [1] "common core"  "mods1"        "foreign pol"  "social sec"
     [5] "mods1"        "immigration"  "economics"    "border"
     [9] "budget"       "Paul Ryan"    "military"     "mods2"
    [13] "gen election" "iran"         "marriage"
```

```r
# pos.neg <- dplyr::select(twitter.df[tweet_indices,], -tweet_created)
# pos.neg <- filter(pos.neg, sentiment != "Neutral")
# pos.neg$candidate[pos.neg$candidate == "OTHER"] <- NA
# pos.neg$subject_matter[pos.neg$subject_matter == "None of the above"] <- NA
# pos.neg <- droplevels(pos.neg)
levels(pos.neg$sentiment)
```

```
    [1] "Negative" "Positive"
```

```r
levels(pos.neg$candidate)
```

```
     [1] "Ben Carson"     "Chris Christie" "Donald Trump"    "Jeb Bush"
     [5] "John Kasich"    "Marco Rubio"    "Mike Huckabee"  "Rand Paul"
     [9] "Scott Walker"   "Ted Cruz"
```

```r
levels(pos.neg$subject_matter)
```

```
 [1] "Abortion"
 [2] "Foreign Policy"
 [3] "FOX News or Moderators"
 [4] "Gun Control"
 [5] "Healthcare (including Medicare)"
 [6] "Immigration"
 [7] "Jobs and Economy"
 [8] "LGBT issues"
 [9] "Racial issues"
[10] "Religion"
[11] "Women's Issues (not abortion though)"
```

```r
dropped.rows <- which(twitter.df[tweet_indices, "sentiment"] == "Neutral")
nrow(pos.neg) + length(dropped.rows) == nrow(twitter.topics$topics)
```

```
[1] TRUE
```

```r
all(dim(simple_lda_15@gamma) == dim(twitter.topics$topics),
    class(simple_lda_15@gamma) == class(twitter.topics$topics))
```

```
[1] TRUE
```

```r
dim(simple_lda_25@gamma[-dropped.rows,])
```

```
[1] 8722   25
```

```r
dim(twitter.topics$topics[-dropped.rows,])
```

```
[1] 8722   15
```

```r
all(abs(rowSums(simple_lda_25@gamma) - 1) < 1e-10)
```

```
[1] TRUE
```

```r
all(abs(rowSums(twitter.topics$topics) - 1) < 1e-10)
```

```
[1] TRUE
```

```r
# build a logistic regression from lda model parameters, additional predictors as parameter
glm_lda_model <- function(lda_model_post, modified_data,
                          predictors = c("candidate", "subject_matter")) {
    x <- lda_model_post[,-2] # need to drop one of the topics, I drop #2
    colnames(x) <- paste("topic", 1:(ncol(x)+1), sep=".")[-2]
    data <- cbind(modified_data, x)
    formula <- paste("sentiment ~ ",
                     paste(c(colnames(x), predictors), collapse = " + "))
```

```
    fit <- glm(as.formula(formula) , data = data, family = "binomial")
    print(summary(fit))
    fit
}

# use forward-backward stepwise procedure with AIC criterion to choose best model from full model
stepwise_twitter <- function(lda_model_post, modified_data,
                    predictors = c("candidate", "subject_matter")) {
  x <- lda_model_post # don't drop any topics
  colnames(x) <- paste("topic", 1:(ncol(x)), sep=".")
  data <- cbind(modified_data, x)
  formula <- paste("sentiment ~ ",
                paste(c(colnames(x), predictors), collapse = " + "))
  fit <- glm(as.formula(formula) , data = data, family = "binomial")
  stepAIC(fit, trace = FALSE) # stops verbose printing
}
```

```
sentiment_twitter_candidate_10 <- glm_lda_model(simple_lda_10@gamma[-dropped.rows,] ,
                            modified_data = pos.neg, predictors = "candidate")
sentiment_twitter_candidate_15 <- glm_lda_model(simple_lda_15@gamma[-dropped.rows,] ,
                            modified_data = pos.neg, predictors = "candidate")
sentiment_twitter_candidate_20 <- glm_lda_model(simple_lda_20@gamma[-dropped.rows,] ,
                            modified_data = pos.neg, predictors = "candidate")
sentiment_twitter_candidate_25 <- glm_lda_model(simple_lda_25@gamma[-dropped.rows,] ,
                            modified_data = pos.neg, predictors = "candidate")
sentiment_twitter_candidate_30 <- glm_lda_model(simple_lda_30@gamma[-dropped.rows,] ,
                            modified_data = pos.neg, predictors = "candidate")
sentiment_twitter_candidate_50 <- glm_lda_model(simple_lda_50@gamma[-dropped.rows,] ,
                            modified_data = pos.neg, predictors = "candidate")
```

```
which.max(c(k10 = simple_lda_10@loglikelihood, k15 = simple_lda_15@loglikelihood,
            k20 = simple_lda_20@loglikelihood, k25 = simple_lda_25@loglikelihood,
            k30 = simple_lda_30@loglikelihood, k50 = simple_lda_50@loglikelihood))
```

```
    k25
      4
```

```
which.min(c(k10 = AIC(sentiment_twitter_candidate_10), k15 = AIC(sentiment_twitter_candidate_15),
          k20 = AIC(sentiment_twitter_candidate_20), k25 = AIC(sentiment_twitter_candidate_25),
          k30 = AIC(sentiment_twitter_candidate_30), k50 = AIC(sentiment_twitter_candidate_50)))
```

```
    k25
      4
```

```
sort(c(k10 = BIC(sentiment_twitter_candidate_10), k15 = BIC(sentiment_twitter_candidate_15),
          k20 = BIC(sentiment_twitter_candidate_20), k25 = BIC(sentiment_twitter_candidate_25),
          k30 = BIC(sentiment_twitter_candidate_30), k50 = BIC(sentiment_twitter_candidate_50)),
     decreasing = TRUE)
```

```
       k50       k30       k15       k25       k20       k10
    5301.986 5183.816 5143.212 5119.215 5089.087 5084.702
```

```
anova(sentiment_twitter_candidate_20, sentiment_twitter_candidate_25, test="Chisq")
```

```
Analysis of Deviance Table

Model 1: sentiment ~ topic.1 + topic.3 + topic.4 + topic.5 + topic.6 +
    topic.7 + topic.8 + topic.9 + topic.10 + topic.11 + topic.12 +
    topic.13 + topic.14 + topic.15 + topic.16 + topic.17 + topic.18 +
    topic.19 + topic.20 + candidate
Model 2: sentiment ~ topic.1 + topic.3 + topic.4 + topic.5 + topic.6 +
    topic.7 + topic.8 + topic.9 + topic.10 + topic.11 + topic.12 +
    topic.13 + topic.14 + topic.15 + topic.16 + topic.17 + topic.18 +
    topic.19 + topic.20 + topic.21 + topic.22 + topic.23 + topic.24 +
    topic.25 + candidate
  Resid. Df Resid. Dev Df Deviance Pr(>Chi)
1      4607      4844.3
2      4602      4832.2  5    12.08   0.0337 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
sentiment_debate_candidate <- glm_lda_model(twitter.topics$topics[-dropped.rows,] ,
                            modified_data = pos.neg, predictors = "candidate")
```

```
Call:
glm(formula = as.formula(formula), family = "binomial", data = data)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-1.6155  -0.7945  -0.6544   1.0755   2.5197

Coefficients:
                       Estimate Std. Error z value Pr(>|z|)
(Intercept)             -0.2373     0.2486  -0.954  0.33988
topic.1                 -0.1039     0.5407  -0.192  0.84760
topic.3                 -0.2809     0.5949  -0.472  0.63680
topic.4                 -0.7597     0.5416  -1.403  0.16069
topic.5                 -0.1932     0.3620  -0.534  0.59353
topic.6                  1.1004     0.4433   2.483  0.01304 *
topic.7                  0.8795     0.5398   1.629  0.10325
topic.8                  0.2195     0.4832   0.454  0.64962
topic.9                  0.5589     0.4449   1.256  0.20901
topic.10                 0.1810     0.4278   0.423  0.67225
topic.11                -0.4272     0.4949  -0.863  0.38800
topic.12                 0.2434     0.3346   0.727  0.46703
topic.13                -0.2289     0.4552  -0.503  0.61505
topic.14                 1.2349     0.4708   2.623  0.00873 **
topic.15                -0.1230     0.4439  -0.277  0.78173
candidateChris Christie -1.7096     0.2362  -7.237 4.57e-13 ***
candidateDonald Trump   -0.9433     0.1441  -6.546 5.91e-11 ***
candidateJeb Bush       -2.7767     0.2191 -12.674  < 2e-16 ***
candidateJohn Kasich     0.4412     0.2087   2.114  0.03450 *
candidateMarco Rubio     0.1148     0.1929   0.595  0.55181
```

```
candidateMike Huckabee   -1.2379      0.2066  -5.993 2.07e-09 ***
candidateRand Paul       -0.8864      0.2200  -4.029 5.59e-05 ***
candidateScott Walker    -1.5672      0.2374  -6.603 4.03e-11 ***
candidateTed Cruz         0.2480      0.1617   1.534  0.12503
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1


(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 5551.6  on 4635  degrees of freedom
Residual deviance: 4967.0  on 4612  degrees of freedom
  (4086 observations deleted due to missingness)
AIC: 5015

Number of Fisher Scoring iterations: 5
```

**AIC**(sentiment_debate_candidate); **BIC**(sentiment_debate_candidate)

```
    [1] 5015.035


    [1] 5169.634
```

**AIC**(sentiment_twitter_candidate_25); **BIC**(sentiment_twitter_candidate_25)

```
    [1] 4900.2


    [1] 5119.215
```

**anova**(sentiment_debate_candidate, sentiment_twitter_candidate_25, test="Chisq")

```
    Analysis of Deviance Table

    Model 1: sentiment ~ topic.1 + topic.3 + topic.4 + topic.5 + topic.6 +
        topic.7 + topic.8 + topic.9 + topic.10 + topic.11 + topic.12 +
        topic.13 + topic.14 + topic.15 + candidate
    Model 2: sentiment ~ topic.1 + topic.3 + topic.4 + topic.5 + topic.6 +
        topic.7 + topic.8 + topic.9 + topic.10 + topic.11 + topic.12 +
        topic.13 + topic.14 + topic.15 + topic.16 + topic.17 + topic.18 +
        topic.19 + topic.20 + topic.21 + topic.22 + topic.23 + topic.24 +
        topic.25 + candidate
      Resid. Df Resid. Dev Df Deviance  Pr(>Chi)
    1      4612     4967.0
    2      4602     4832.2 10   134.84 < 2.2e-16 ***
    ---
    Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

**require**(MASS)

```
    Loading required package: MASS
```

```
step_25_candidate_subject <- stepwise_twitter(simple_lda_25@gamma[-dropped.rows,],
                                pos.neg.sub, predictors = c("candidate", "subject_matter"))
step_25_candidate_subject$anova
```

```
    Stepwise Model Path
    Analysis of Deviance Table

    Initial Model:
    sentiment ~ topic.1 + topic.2 + topic.3 + topic.4 + topic.5 +
        topic.6 + topic.7 + topic.8 + topic.9 + topic.10 + topic.11 +
        topic.12 + topic.13 + topic.14 + topic.15 + topic.16 + topic.17 +
        topic.18 + topic.19 + topic.20 + topic.21 + topic.22 + topic.23 +
        topic.24 + topic.25 + candidate + subject_matter

    Final Model:
    sentiment ~ topic.1 + topic.2 + topic.3 + topic.4 + topic.6 +
        topic.8 + topic.11 + topic.13 + topic.14 + topic.15 + topic.16 +
        topic.17 + topic.19 + topic.22 + topic.23 + topic.24 + candidate +
        subject_matter


                  Step Df     Deviance Resid. Df Resid. Dev       AIC
    1                                      8676    7371.464 7463.464
    2  - topic.25  0 0.00000000      8676    7371.464 7463.464
    3   - topic.5  1 0.04248644      8677    7371.506 7461.506
    4  - topic.10  1 0.06326254      8678    7371.570 7459.570
    5   - topic.9  1 0.06924882      8679    7371.639 7457.639
    6  - topic.21  1 0.16682058      8680    7371.806 7455.806
    7   - topic.7  1 0.49511827      8681    7372.301 7454.301
    8  - topic.18  1 0.72221103      8682    7373.023 7453.023
    9  - topic.20  1 1.08137289      8683    7374.104 7452.104
    10 - topic.12  1 1.47842730      8684    7375.583 7451.583
```

```
summary(step_25_candidate_subject)
```

```
    Call:
    glm(formula = sentiment ~ topic.1 + topic.2 + topic.3 + topic.4 +
        topic.6 + topic.8 + topic.11 + topic.13 + topic.14 + topic.15 +
        topic.16 + topic.17 + topic.19 + topic.22 + topic.23 + topic.24 +
        candidate + subject_matter, family = "binomial", data = data)

    Deviance Residuals:
        Min       1Q   Median       3Q      Max
    -1.9741  -0.6583  -0.4530  -0.2559   2.8157

    Coefficients:
                                          Estimate Std. Error
    (Intercept)                          -5.282085   1.354942
    topic.1                              13.711103   4.049753
    topic.2                             -11.620507   5.948299
    topic.3                              10.878613   4.879918
```

```
topic.4                                                10.579727    3.776517
topic.6                                                11.553613    4.852787
topic.8                                                 8.867997    5.847927
topic.11                                               -8.925792    5.476606
topic.13                                                9.721931    5.511908
topic.14                                               27.278398    4.883344
topic.15                                               30.948666    4.901999
topic.16                                              -13.841999    7.144440
topic.17                                              -11.059203    5.432706
topic.19                                               22.451599    4.801794
topic.22                                               15.048531    3.703971
topic.23                                              -19.988527    6.336070
topic.24                                               24.278079    5.178959
candidateChris Christie                                -1.553775    0.255412
candidateDonald Trump                                  -1.180442    0.138379
candidateJeb Bush                                      -2.574033    0.260114
candidateJohn Kasich                                    0.318412    0.203706
candidateMarco Rubio                                   -0.201864    0.199531
candidateMike Huckabee                                 -0.999883    0.225492
candidateRand Paul                                     -1.052019    0.213535
candidateScott Walker                                  -1.630008    0.232306
candidateTed Cruz                                      -0.123238    0.169722
candidateother                                         -2.110719    0.136885
subject_matterForeign Policy                            0.147148    0.301641
subject_matterFOX News or Moderators                    0.293252    0.245595
subject_matterGun Control                             -13.193683  215.215383
subject_matterHealthcare (including Medicare)           0.487323    0.436864
subject_matterImmigration                               0.912597    0.316465
subject_matterJobs and Economy                         -0.001804    0.318501
subject_matterLGBT issues                               0.148386    0.367211
subject_matterRacial issues                            -0.827960    0.343091
subject_matterReligion                                 -0.711367    0.341760
subject_matterWomen's Issues (not abortion though)     -1.306857    0.398521
subject_matterother                                     0.737697    0.236063
                                                       z value Pr(>|z|)
(Intercept)                                            -3.898 9.68e-05 ***
topic.1                                                 3.386  0.00071 ***
topic.2                                                -1.954  0.05075 .
topic.3                                                 2.229  0.02580 *
topic.4                                                 2.801  0.00509 **
topic.6                                                 2.381  0.01727 *
topic.8                                                 1.516  0.12941
topic.11                                               -1.630  0.10314
topic.13                                                1.764  0.07776 .
topic.14                                                5.586 2.32e-08 ***
topic.15                                                6.313 2.73e-10 ***
topic.16                                               -1.937  0.05269 .
topic.17                                               -2.036  0.04178 *
topic.19                                                4.676 2.93e-06 ***
topic.22                                                4.063 4.85e-05 ***
topic.23                                               -3.155  0.00161 **
topic.24                                                4.688 2.76e-06 ***
candidateChris Christie                                -6.083 1.18e-09 ***
candidateDonald Trump                                  -8.531  < 2e-16 ***
```

```
candidateJeb Bush                                       -9.896  < 2e-16 ***
candidateJohn Kasich                                     1.563  0.11803
candidateMarco Rubio                                    -1.012  0.31169
candidateMike Huckabee                                  -4.434 9.24e-06 ***
candidateRand Paul                                      -4.927 8.36e-07 ***
candidateScott Walker                                   -7.017 2.27e-12 ***
candidateTed Cruz                                       -0.726  0.46777
candidateother                                         -15.420  < 2e-16 ***
subject_matterForeign Policy                             0.488  0.62567
subject_matterFOX News or Moderators                     1.194  0.23246
subject_matterGun Control                               -0.061  0.95112
subject_matterHealthcare (including Medicare)            1.116  0.26464
subject_matterImmigration                                2.884  0.00393 **
subject_matterJobs and Economy                          -0.006  0.99548
subject_matterLGBT issues                                0.404  0.68615
subject_matterRacial issues                             -2.413  0.01581 *
subject_matterReligion                                  -2.081  0.03739 *
subject_matterWomen's Issues (not abortion though)  -3.279  0.00104 **
subject_matterother                                      3.125  0.00178 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1


(Dispersion parameter for binomial family taken to be 1)


    Null deviance: 8774.8  on 8721  degrees of freedom
Residual deviance: 7375.6  on 8684  degrees of freedom
AIC: 7451.6


Number of Fisher Scoring iterations: 14
```

```
step_debate_topics <- stepwise_twitter(twitter.topics$topics[-dropped.rows,],
                    pos.neg.sub, predictors = c("candidate", "subject_matter"))
step_debate_topics$anova
```

```
    Stepwise Model Path
    Analysis of Deviance Table

    Initial Model:
    sentiment ~ topic.1 + topic.2 + topic.3 + topic.4 + topic.5 +
        topic.6 + topic.7 + topic.8 + topic.9 + topic.10 + topic.11 +
        topic.12 + topic.13 + topic.14 + topic.15 + candidate + subject_matter

    Final Model:
    sentiment ~ topic.3 + topic.4 + topic.5 + topic.9 + candidate +
        subject_matter


             Step Df  Deviance Resid. Df Resid. Dev      AIC
    1                              8686   7506.915 7578.915
    2 - topic.15  0 0.0000000     8686   7506.915 7578.915
    3 - topic.11  1 0.0361263     8687   7506.951 7576.951
    4 - topic.13  1 0.1544666     8688   7507.105 7575.105
    5  - topic.2  1 0.1144348     8689   7507.220 7573.220
    6  - topic.1  1 0.1522128     8690   7507.372 7571.372
```

```
7   - topic.10  1 1.0886845      8691   7508.461 7570.461
8   - topic.12  1 1.2258320      8692   7509.687 7569.687
9    - topic.6  1 0.8967099      8693   7510.583 7568.583
10  - topic.14  1 0.9235075      8694   7511.507 7567.507
11   - topic.8  1 0.9577672      8695   7512.465 7566.465
12   - topic.7  1 1.3343270      8696   7513.799 7565.799
```

```
summary(step_debate_topics)
```

```
Call:
glm(formula = sentiment ~ topic.3 + topic.4 + topic.5 + topic.9 +
    candidate + subject_matter, family = "binomial", data = data)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-1.4871  -0.6432  -0.4516  -0.2665   3.0724

Coefficients:
                                                   Estimate Std. Error
(Intercept)                                        -0.38258    0.26277
topic.3                                            -0.84461    0.44540
topic.4                                            -1.38034    0.39837
topic.5                                            -0.54730    0.20716
topic.9                                             0.44076    0.26667
candidateChris Christie                            -1.83143    0.23146
candidateDonald Trump                              -0.95344    0.13538
candidateJeb Bush                                  -2.96046    0.21748
candidateJohn Kasich                                0.31489    0.20195
candidateMarco Rubio                                0.03744    0.19177
candidateMike Huckabee                             -1.27677    0.19832
candidateRand Paul                                 -1.01944    0.20981
candidateScott Walker                              -1.54588    0.22962
candidateTed Cruz                                   0.14478    0.15684
candidateother                                     -1.97076    0.13278
subject_matterForeign Policy                        0.03311    0.29551
subject_matterFOX News or Moderators                0.18841    0.23921
subject_matterGun Control                         -12.40812  131.77517
subject_matterHealthcare (including Medicare)       0.50775    0.43309
subject_matterImmigration                           0.93757    0.29532
subject_matterJobs and Economy                      0.06306    0.31517
subject_matterLGBT issues                           0.08448    0.35991
subject_matterRacial issues                        -0.89991    0.34039
subject_matterReligion                             -0.76217    0.33691
subject_matterWomen's Issues (not abortion though) -1.32550    0.39724
subject_matterother                                 0.69955    0.23190
                                                   z value Pr(>|z|)
(Intercept)                                         -1.456 0.145414
topic.3                                             -1.896 0.057922 .
topic.4                                             -3.465 0.000530 ***
topic.5                                             -2.642 0.008245 **
topic.9                                              1.653 0.098367 .
candidateChris Christie                             -7.912 2.52e-15 ***
candidateDonald Trump                               -7.043 1.89e-12 ***
```

```
candidateJeb Bush                                    -13.613  < 2e-16 ***
candidateJohn Kasich                                   1.559 0.118928
candidateMarco Rubio                                   0.195 0.845204
candidateMike Huckabee                                -6.438 1.21e-10 ***
candidateRand Paul                                    -4.859 1.18e-06 ***
candidateScott Walker                                 -6.732 1.67e-11 ***
candidateTed Cruz                                      0.923 0.355957
candidateother                                       -14.843  < 2e-16 ***
subject_matterForeign Policy                           0.112 0.910784
subject_matterFOX News or Moderators                   0.788 0.430916
subject_matterGun Control                             -0.094 0.924981
subject_matterHealthcare (including Medicare)          1.172 0.241037
subject_matterImmigration                              3.175 0.001500 **
subject_matterJobs and Economy                         0.200 0.841410
subject_matterLGBT issues                              0.235 0.814423
subject_matterRacial issues                           -2.644 0.008199 **
subject_matterReligion                                -2.262 0.023681 *
subject_matterWomen's Issues (not abortion though)  -3.337 0.000848 ***
subject_matterother                                    3.017 0.002557 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1


(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 8774.8  on 8721  degrees of freedom
Residual deviance: 7513.8  on 8696  degrees of freedom
AIC: 7565.8

Number of Fisher Scoring iterations: 13
```

`anova(step_debate_topics, step_25_candidate_subject, test="Chisq")`

```
    Analysis of Deviance Table

    Model 1: sentiment ~ topic.3 + topic.4 + topic.5 + topic.9 + candidate +
        subject_matter
    Model 2: sentiment ~ topic.1 + topic.2 + topic.3 + topic.4 + topic.6 +
        topic.8 + topic.11 + topic.13 + topic.14 + topic.15 + topic.16 +
        topic.17 + topic.19 + topic.22 + topic.23 + topic.24 + candidate +
        subject_matter
      Resid. Df Resid. Dev Df Deviance  Pr(>Chi)
    1      8696     7513.8
    2      8684     7375.6 12   138.22 < 2.2e-16 ***
    ---
    Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```