

NYT Phrase Rank Lists

From left to right, below is the list of top 30, middle 30, and bottom 30 single word phrases for the NYT dataset with the suggested parameters of *HIGHLIGHT_SINGLE* = .9 and *HIGHLIGHT_MULTI* = .5.

| | |
|-------------|--------------|
| wbc | 0.8512157704 |
| minneapolis | 0.8472041311 |
| dna | 0.8470183261 |
| cska | 0.8468338448 |
| cbs | 0.8423605647 |
| wellington | 0.8421586470 |
| iv | 0.8421261359 |
| fbi | 0.8420124957 |
| dc | 0.8414233784 |
| nba | 0.8411841758 |
| fx | 0.8405590938 |
| istanbul | 0.8399355593 |
| prague | 0.8398781518 |
| amc | 0.8394182441 |
| rbs | 0.8389925043 |
| clarence | 0.8389280748 |
| lima | 0.8385657743 |
| sacramento | 0.8385282580 |
| tripoli | 0.8381023545 |
| dublin | 0.8379963028 |
| lgbt | 0.8379090817 |
| bcs | 0.8374214900 |
| usa | 0.8372757945 |
| albany | 0.8370727411 |
| naples | 0.8368999479 |
| lexington | 0.8368385291 |
| cal | 0.8363055371 |
| milan | 0.8362579035 |
| amsterdam | 0.8361695801 |

Table 1: Top 30

| | |
|--------------|--------------|
| intentional | 0.5870088101 |
| prodced | 0.5870083729 |
| ordinarily | 0.5870083729 |
| confront | 0.5870008670 |
| soprano | 0.5869853623 |
| machinery | 0.5869702128 |
| riot | 0.5869621683 |
| waiting | 0.5869557429 |
| atrocities | 0.5869556507 |
| prevail | 0.5869546957 |
| here's | 0.5869530878 |
| interrupted | 0.5869266819 |
| smokers | 0.5869212730 |
| disruption | 0.5868902329 |
| originally | 0.5868889273 |
| slam | 0.5868875053 |
| poisoning | 0.5868756404 |
| desperately | 0.5868640156 |
| siege | 0.5868636363 |
| widow | 0.5868611921 |
| buy | 0.5868535537 |
| recreational | 0.5868499478 |
| hosting | 0.5868189940 |
| unusually | 0.5867760126 |
| wheelchairs | 0.5867586481 |
| torque | 0.5867586481 |
| strikes | 0.5867349929 |
| ribs | 0.5867307076 |
| endless | 0.5867247648 |

Table 2: Mid 30

| | |
|------------|--------------|
| 172 | 0.2744979322 |
| 06 | 0.2742260498 |
| 332 | 0.2739979322 |
| 299 | 0.2739979322 |
| 285 | 0.2739979322 |
| 139 | 0.2738766169 |
| 126 | 0.2734599502 |
| 214 | 0.2734599502 |
| 03 | 0.2728980737 |
| 320 | 0.2728980737 |
| 325 | 0.2726260498 |
| 157 | 0.2724599502 |
| 02 | 0.2723790519 |
| 375 | 0.2723790519 |
| containing | 0.2721296301 |
| 158 | 0.2719718549 |
| 380 | 0.2719718549 |
| 116 | 0.2710698596 |
| specified | 0.2696099502 |
| 124 | 0.2694850748 |
| 145 | 0.2691612652 |
| 280 | 0.2685555739 |
| 107 | 0.2683433052 |
| 111 | 0.2678667286 |
| 112 | 0.2673443406 |
| 114 | 0.2667911009 |
| sincere | 0.2666137439 |
| ought | 0.2653443406 |
| 102 | 0.2644190375 |

Table 3: Bottom 30

These lists above make sense because the top 30 are all common acronyms and proper nouns, the mid 30 are frequent and meaningful nouns and verbs, and the bottom 30 are mostly random numbers.

From left to right, below is the list of top 30, middle 30, and bottom 30 multi word phrases for the NYT dataset with the suggested parameters of *HIGHLIGHT_SINGLE* = .9 and *HIGHLIGHT_MULTI* = .5. These lists below for multi word phrases are intuitive because the top 30 are names of people and organizations, the middle 30 are common phrases in speech, and the bottom 30 are prepositions and transitional phrases.

| | |
|-------------------|--------------|
| kyrie irving | 0.9859670300 |
| adrian peterson | 0.9858567268 |
| justin bieber | 0.9855607233 |
| aston villa | 0.9852497353 |
| richie incognito | 0.9851547916 |
| marshawn lynch | 0.9848564661 |
| fidel castro | 0.9845802860 |
| napa valley | 0.9845547987 |
| graeme swann | 0.9843076584 |
| nick saban | 0.9842913079 |
| rafael nadal | 0.9841839528 |
| henrik lundqvist | 0.9841120842 |
| justin timberlake | 0.9840859622 |
| gareth bale | 0.9840295465 |
| tiger woods | 0.9839738113 |
| jamaal charles | 0.9839603943 |
| sears holdings | 0.9839533040 |
| dwight howard | 0.9838893572 |
| alan mulally | 0.9838876767 |
| derek stepan | 0.9838231328 |
| walt disney | 0.9837392413 |
| nick foles | 0.9836625843 |
| Michel Djotodia | 0.9836394567 |
| brook lopez | 0.9835856279 |
| hurricane katrina | 0.9835716916 |
| tim duncan | 0.9835537588 |
| ubs ag | 0.9835314120 |
| sebastian vettel | 0.9834178661 |
| janet yellen | 0.9833784342 |
| golan heights | 0.9833098495 |

Table 4: Top 30

| | |
|------------------------|--------------|
| handed down | 0.1278334721 |
| stand up | 0.1278236174 |
| ranges from | 0.1278187086 |
| from public life | 0.1278089004 |
| of comcast corp | 0.1277916275 |
| wars in iraq and | 0.1277562480 |
| doesn't include | 0.1277504663 |
| energy bill | 0.1277337884 |
| as south africa's | 0.1277202886 |
| well aware | 0.1277180673 |
| 35 minutes | 0.1277144797 |
| in tax breaks | 0.1277116959 |
| israeli military said | 0.1276519684 |
| get hurt | 0.1276420641 |
| gave me | 0.1276355675 |
| of political prisoners | 0.1276305578 |
| announced tuesday that | 0.1275690246 |
| in june 2012 | 0.1275658157 |
| dominican republic and | 0.1275312384 |
| slightly better | 0.1275279922 |
| political system | 0.1275074171 |
| to enlist | 0.1275043352 |
| left behind | 0.1274784624 |
| last four years | 0.1274718913 |
| still alive | 0.1274561571 |
| without power | 0.1274171262 |
| reunite with | 0.1273897374 |
| former employees | 0.1273891631 |
| to be | 0.1273639799 |
| 60 years | 0.1273436734 |

Table 5: Mid 30

| | |
|-----------------------|--------------|
| one hand and | 0.0075959662 |
| of months of | 0.0075959662 |
| of pennsylvania and | 0.0075959662 |
| to syria to | 0.0075959662 |
| from home and | 0.0075959662 |
| three games with | 0.0075959662 |
| of meetings with | 0.0075959662 |
| for consumers to | 0.0075959662 |
| a lot to | 0.0075959662 |
| to happen in | 0.0075959662 |
| of murder in | 0.0075959662 |
| in court on | 0.0075959662 |
| for oil and | 0.0075959662 |
| nine years in | 0.0075959662 |
| to sign with | 0.0075959662 |
| to grant a | 0.0075959662 |
| to file for | 0.0075959662 |
| to die in | 0.0075959662 |
| in mind that | 0.0075959662 |
| of violence in the | 0.0075959662 |
| of millions of people | 0.0075959662 |
| for sale in | 0.0075959662 |
| to study in | 0.0075959662 |
| to turn this | 0.0075959662 |
| in court that | 0.0075959662 |
| of billions of | 0.0075959662 |
| to reporters at | 0.0075959662 |
| of money in | 0.0075959662 |
| of state for | 0.0075959662 |
| said she | 0.0070537286 |

Table 6: Bottom 30

Yelp Phrase Rank Lists

From left to right, the lists on the following page are of the top 30, middle 30, and bottom 30 single word phrases for the NYT dataset with the suggested parameters of *HIGHLIGHT_SINGLE* = .9 and *HIGHLIGHT_MULTI* = .5.

These lists make sense because the top 30 are mostly all names of foods and drinks, the middle 30 are common nouns and verbs associated with restaurants and services, and the bottom 30 are mostly prepositions that don't have much meaning.

The lists on page 4 are the top, middle and bottom 30 multi word phrases in the Yelp dataset. Again these lists are intuitive because the top 30 are organization names or names of specific foods, the middle 30 are common foods and speech phrases, and the bottom 30 are more prepositions and transitions.

| | |
|-----------|--------------|
| trail | 0.9283774739 |
| ale | 0.9279517675 |
| sangria | 0.9251320512 |
| yelper | 0.9231119693 |
| hangover | 0.9215918356 |
| eggplant | 0.9210159397 |
| gold | 0.9208643534 |
| library | 0.9205310201 |
| stew | 0.9190822358 |
| church | 0.9189934626 |
| latte | 0.9187752054 |
| cherry | 0.9184212028 |
| espresso | 0.9181488972 |
| yelpers | 0.9176200685 |
| vegas | 0.9171503716 |
| viet | 0.9168684378 |
| brie | 0.9168375948 |
| mole | 0.9162541174 |
| lifetime | 0.9160440823 |
| flatbread | 0.9158635588 |
| parmesan | 0.9155271740 |
| brewery | 0.9152714963 |
| ribeye | 0.9152654524 |
| zin | 0.9152449625 |
| meatball | 0.9151991072 |
| fitness | 0.9151448717 |
| farms | 0.9150685900 |
| coach | 0.9150399127 |
| ranch | 0.9150369352 |

Table 7: Top 30

| | |
|---------------|--------------|
| kinda | 0.7716520296 |
| basement | 0.7716468228 |
| fricken | 0.7716241467 |
| posts | 0.7716208485 |
| caliber | 0.7716175362 |
| diverse | 0.7716133299 |
| dimsum | 0.7716003917 |
| famed | 0.7715747751 |
| cheezy | 0.7715329133 |
| preschool | 0.7715329133 |
| isolated | 0.7715109722 |
| hope | 0.7714825538 |
| traditionally | 0.7714700443 |
| haircut | 0.7714621814 |
| lengthy | 0.7714562774 |
| attendance | 0.7714516128 |
| resteraunt | 0.7714475621 |
| pigeons | 0.7714475621 |
| app's | 0.7714475621 |
| ashes | 0.7714475621 |
| leadership | 0.7714475621 |
| issue | 0.7714269907 |
| 43rd | 0.7714168844 |
| ubiquitous | 0.7713882299 |
| softener | 0.7713855170 |
| chipper | 0.7713338390 |
| wage | 0.7713290252 |
| sites | 0.7713123338 |
| days | 0.7713087739 |

Table 8: Mid 30

| | |
|------------|--------------|
| amongst | 0.3745894890 |
| 06 | 0.3734362821 |
| werent | 0.3732478184 |
| been | 0.3732182178 |
| which | 0.3727403290 |
| containing | 0.3723831115 |
| 09 | 0.3723831115 |
| yourself | 0.3710446346 |
| me | 0.3708902248 |
| awfully | 0.3707345023 |
| 02 | 0.3707345023 |
| becomes | 0.3707275842 |
| can | 0.3704885490 |
| merely | 0.3704784044 |
| could | 0.3700245072 |
| 05 | 0.3674529488 |
| be | 0.3664271359 |
| gives | 0.3662908717 |
| themselves | 0.3659551358 |
| him | 0.3632958684 |
| have | 0.3631863165 |
| ourselves | 0.3618827637 |
| is | 0.3612317711 |
| has | 0.3610493581 |
| them | 0.3608067799 |
| are | 0.3587614793 |
| were | 0.3577612489 |
| was | 0.3570231537 |
| myself | 0.3561846821 |

Table 9: Bottom 30

| | | | |
|--------------------|--------------|-----------------------|--------------|
| jamba juice | 0.9855316025 | the wedge salad | 0.0885694854 |
| grand marnier | 0.9837130519 | sitting on top of | 0.0885663980 |
| hyatt regency | 0.9834038475 | the old days | 0.0885654616 |
| botanical gardens | 0.9819158655 | a bad choice | 0.0885610444 |
| jimmy johns | 0.9815106738 | bent on | 0.0885559050 |
| hobby lobby | 0.9798741782 | especially since | 0.0885509178 |
| palo verde | 0.9795939574 | a shopping center | 0.0885478251 |
| del rey | 0.9795771212 | the mediterranean | 0.0885436086 |
| orange blossom | 0.9792773294 | with arugula | 0.0885353958 |
| peter piper | 0.9791724839 | flow of | 0.0885310512 |
| taliesin west | 0.9789654644 | an incredible | 0.0885294980 |
| cheesecake factory | 0.9788442154 | that point | 0.0885292634 |
| cabernet sauvignon | 0.9786717397 | driving into | 0.0885289634 |
| del mar | 0.9785584259 | the easiest | 0.0885265352 |
| english muffin | 0.9782791894 | starting to get | 0.0885263643 |
| chile relleno | 0.9782021344 | particularly care for | 0.0885217908 |
| el paso | 0.9781792604 | to order drinks | 0.0885201542 |
| pollo asado | 0.9781085274 | to dance | 0.0885161429 |
| daily dose | 0.9780325640 | by weight | 0.0885105557 |
| red robin | 0.9780062009 | a medium | 0.0885099993 |
| squaw peak | 0.9779615017 | diagnosed with | 0.0885098387 |
| union hills | 0.9779233649 | away from | 0.0885085989 |
| upper crust | 0.9777963502 | at ikea | 0.0885037011 |
| delhi palace | 0.9777607192 | a nice outdoor patio | 0.0884998262 |
| carl's jr | 0.9777348202 | does not matter | 0.0884994810 |
| johnny rockets | 0.9776970828 | all night | 0.0884984343 |
| tater tots | 0.9776615624 | for sunday brunch | 0.0884845028 |
| michael mina | 0.9776529286 | a yummy | 0.0884817247 |
| ann taylor | 0.9776391505 | every other | 0.0884773489 |
| dos equis | 0.9775953971 | an enthusiastic | 0.0884728286 |

Table 10: Top 30

| | |
|-----------------------|--------------|
| the wedge salad | 0.0885694854 |
| sitting on top of | 0.0885663980 |
| the old days | 0.0885654616 |
| a bad choice | 0.0885610444 |
| bent on | 0.0885559050 |
| especially since | 0.0885509178 |
| a shopping center | 0.0885478251 |
| the mediterranean | 0.0885436086 |
| with arugula | 0.0885353958 |
| flow of | 0.0885310512 |
| an incredible | 0.0885294980 |
| that point | 0.0885292634 |
| driving into | 0.0885289634 |
| the easiest | 0.0885265352 |
| starting to get | 0.0885263643 |
| particularly care for | 0.0885217908 |
| to order drinks | 0.0885201542 |
| to dance | 0.0885161429 |
| by weight | 0.0885105557 |
| a medium | 0.0885099993 |
| diagnosed with | 0.0885098387 |
| away from | 0.0885085989 |
| at ikea | 0.0885037011 |
| a nice outdoor patio | 0.0884998262 |
| does not matter | 0.0884994810 |
| all night | 0.0884984343 |
| for sunday brunch | 0.0884845028 |
| a yummy | 0.0884817247 |
| every other | 0.0884773489 |
| an enthusiastic | 0.0884728286 |

Table 11: Mid 30

| | |
|-----------------------|--------------|
| to venture to | 0.0067325702 |
| a search for | 0.0067175661 |
| to go again and | 0.0066756286 |
| some research on | 0.0065839591 |
| got seated at | 0.0065818086 |
| no issues with | 0.0065001504 |
| and it's fun to | 0.0064707003 |
| 1 star because | 0.0064616889 |
| and just wanted to | 0.0064463948 |
| the nail on | 0.0064420109 |
| and can't wait to go | 0.0064151448 |
| and very easy to | 0.0063960832 |
| people complain about | 0.0063864553 |
| i'm working on | 0.0061836335 |
| usually sit at | 0.0061836335 |
| got tired of | 0.0061514703 |
| actually care about | 0.0061280780 |
| very reminiscent of | 0.0061191174 |
| so i'll keep | 0.0061012530 |
| the peak of | 0.0060635619 |
| i've worked in | 0.0059302285 |
| a touch more | 0.0059302285 |
| to relax with | 0.0059064550 |
| a block of | 0.0058506258 |
| to call ahead and | 0.0058439550 |
| to kick back | 0.0057881258 |
| the only drawback to | 0.0057675661 |
| a cheeseburger with | 0.0057606217 |
| a scale of | 0.0057047924 |
| any combination of | 0.0049302285 |

Table 12: Bottom 30

Phrasal Segmentation Metrics

The below table shows the number of unique qualified phrases and average number of phrases per sentence for each dataset for *HIGHLIGHT_SINGLE* = .9 and *HIGHLIGHT_MULTI* = .5. These metrics include both single and multi word phrases.

| | NYT | Yelp |
|----------------------|--------|---------|
| Avg per sentence | .94149 | .780041 |
| Total unique phrases | 65552 | 41988 |

Table 13: Phrase metrics

Since NYT is professionally written and edited, I would expect a natural language processor to more easily identify phrases, which is reflected in the data. Yelp reviews contain a lot of slang and bad grammar, which might make the task of identifying phrases more difficult, yielding fewer phrases identified overall.

NYT Phrase Clusters

Below are samples of six clusters obtained from using the k means algorithm with 75 centers on the NYT dataset. I only clustered phrases returned by AutoPhrase and used each phrase’s vector from word2vec to do the clustering. I’ve labeled the concept each cluster is describing in the table captions.

| |
|--------------|
| _government_ |
| _state_ |
| _country_ |
| _public_ |
| _city_ |
| _capital_ |
| _Syria_ |
| _corruption_ |
| _politics_ |
| _army_ |
| _Iraq_ |
| _Muslim_ |
| _Islamic_ |
| _Islamist_ |
| _wave_ |
| _Assad_ |
| _ethnic_ |
| _suicide_ |
| _Libya_ |
| _toll_ |

Table 14: Middle East Politics

| |
|---------------------|
| _revenue_ |
| _range_ |
| _auction_ |
| _billion_euros_ |
| _million_euros_ |
| _million_pounds_ |
| _euros_ |
| _acres_ |
| _billion_pounds_ |
| _lease_ |
| _temperature_ |
| _square_feet_ |
| _ticket_sales_ |
| _trillion_yen_ |
| _miles_per_hour_ |
| _million_shares_ |
| _contemporary_art_ |
| _square_foot_ |
| _million_Americans_ |
| _million_viewers_ |

Table 15: Quantities

| |
|-----------------------|
| _federal_ |
| _health_ |
| _Obama_ |
| _tax_ |
| _website_ |
| _Americans_ |
| _insurance_ |
| _cover_ |
| _bills_ |
| _federal_government_ |
| _health_insurance_ |
| _insurers_ |
| _Medicaid_ |
| _Medicare_ |
| _Affordable_Care_Act_ |
| _minimum_wage_ |
| _HealthCare.gov_ |
| _Obamacare_ |
| _premiums_ |
| _Social_Security_ |

Table 16: Obamacare

| |
|-----------------|
| _U.S._ |
| _United_States_ |
| _American_ |
| _China_ |
| _Chinese_ |
| _region_ |
| _Israel_ |
| _Japan_ |
| _Western_ |
| _Japanese_ |
| _Beijing_ |
| _North_Korea_ |
| _Saudi_ |
| _Iranian_ |
| _South_Korea_ |
| _Pakistan_ |
| _Turkey_ |
| _Middle_East_ |
| _Egypt_ |
| _Arab_ |

Table 17: State Entities

| |
|--------------------------|
| _Republican_ |
| _Senate_ |
| _House_ |
| _Republicans_ |
| _White_House_ |
| _Democrats_ |
| _Democratic_ |
| _Democrat_ |
| _Senator_ |
| _President_Barack_Obama_ |
| _President_Obama_ |
| _congressional_ |
| _legislative_ |
| _Capitol_ |
| _conservatives_ |
| _liberal_ |
| _senator_ |
| _Christie_ |
| _Tea_Party_ |
| _Legislature_ |

Table 18: American Politics

| |
|--------------------|
| _history_ |
| _football_ |
| _manager_ |
| _star_ |
| _college_ |
| _sports_ |
| _English_ |
| _soccer_ |
| _sport_ |
| _basketball_ |
| _professional_ |
| _baseball_ |
| _golf_ |
| _hockey_ |
| _tennis_ |
| _national_team_ |
| _cricket_ |
| _Argentine_ |
| _athlete_ |
| _college_football_ |

Table 19: Sports

Yelp Phrase Clusters

On this and the following page I've used the same process as described for NYT on the Yelp dataset.

| |
|------------------------|
| _food_ |
| _Food_ |
| _pleasantly_surprised_ |
| _dining_experience_ |
| _blown_away_ |
| _Atmosphere_ |
| _poor_service_ |
| _slow_service_ |
| _Taste_ |
| _Customer_service_ |
| _Average_ |
| _mediocre_food_ |
| _bottom_line_ |
| _entire_experience_ |
| _noise_level_ |
| _non-existent_ |
| _pretty_slow_ |
| _wait_times_ |
| _terrible_service_ |
| _Wait_staff_ |

Table 20: Food Service/Experience

| |
|------------------|
| _pizza_ |
| _salad_ |
| _sandwich_ |
| _sandwiches_ |
| _salads_ |
| _pasta_ |
| _hummus_ |
| _Pizza_ |
| _bruschetta_ |
| _pita_ |
| _wrap_ |
| _mushroom_ |
| _soups_ |
| _pepperoni_ |
| _thin_crust_ |
| _gyro_ |
| _grilled_cheese_ |
| _italian_ |
| _falafel_ |
| _panini_ |

Table 21: Italian/Greek Food

| |
|-------------------|
| _ok_ |
| _average_ |
| _OK_ |
| _Ok_ |
| _basic_ |
| _pretty_decent_ |
| _stellar_ |
| _bar_food_ |
| _pretty_tasty_ |
| _pretty_darn_ |
| _pretty_damn_ |
| _sub-par_ |
| _pretty_awesome_ |
| _pretty_standard_ |
| _pretty_bad_ |
| _looked_pretty_ |
| _downhill_ |
| _sub-par_ |
| _below_average_ |
| _pretty_solid_ |

Table 22: Descriptions

| |
|----------------------|
| _Orange_Table_ |
| _La_Grande_Orange_ |
| _Two_Hippies_ |
| _Z_Pizza_ |
| _5th_and_Wine_ |
| _Cheba_Hut_ |
| _Roaring_Fork_ |
| _Local_Breeze_ |
| _Mellow_Mushroom_ |
| _Bagels_ |
| _Chino_ |
| _Buffalo_Wild_Wings_ |
| _Tilted_Kilt_ |
| _Melting_Pot_ |
| _PF_Chang's_ |
| _Al's_ |
| _BBQ_place_ |
| _original_location_ |
| _Jimmy_Johns_ |
| _Ray's_ |

Table 23: Popular Restaurants

| |
|--------------------------|
| _tex-mex_ |
| _Cambodian_ |
| _Taiwanese_ |
| _real_Mexican_food_ |
| _family_-owned |
| _family_-run |
| _plate_lunch_ |
| _barrio_ |
| _highly_rated_ |
| _lunch_option_ |
| _locally-owned_ |
| _mom-and-pop_ |
| _family_owned_ |
| _fast_casual_ |
| _mom_ |
| _pop_ |
| _authentic_Chinese_ |
| _authentic_Japanese_ |
| _Viet_ |
| _Asian_cuisine_ |
| _authentic_Chinese_food_ |

Table 24: Authentic Food

| |
|-----------------------|
| _Hallmark_ |
| _Harkins_theater_ |
| _Bath_and_Body_Works_ |
| _big-box_ |
| _online_shopping_ |
| _Williams_Sonoma_ |
| _Half_Price_ |
| _Stein_Mart_ |
| _Kombucha_ |
| _Brookstone_ |
| _Ted_Baker_ |
| _kayaking_ |
| _randomness_ |
| _Grocery_store_ |
| _In-N-Out's_ |
| _Axis/_Radius_ |
| _Foot_Locker_ |
| _Disney_Store_ |
| _non-fiction_ |
| _Weight_Watchers_ |

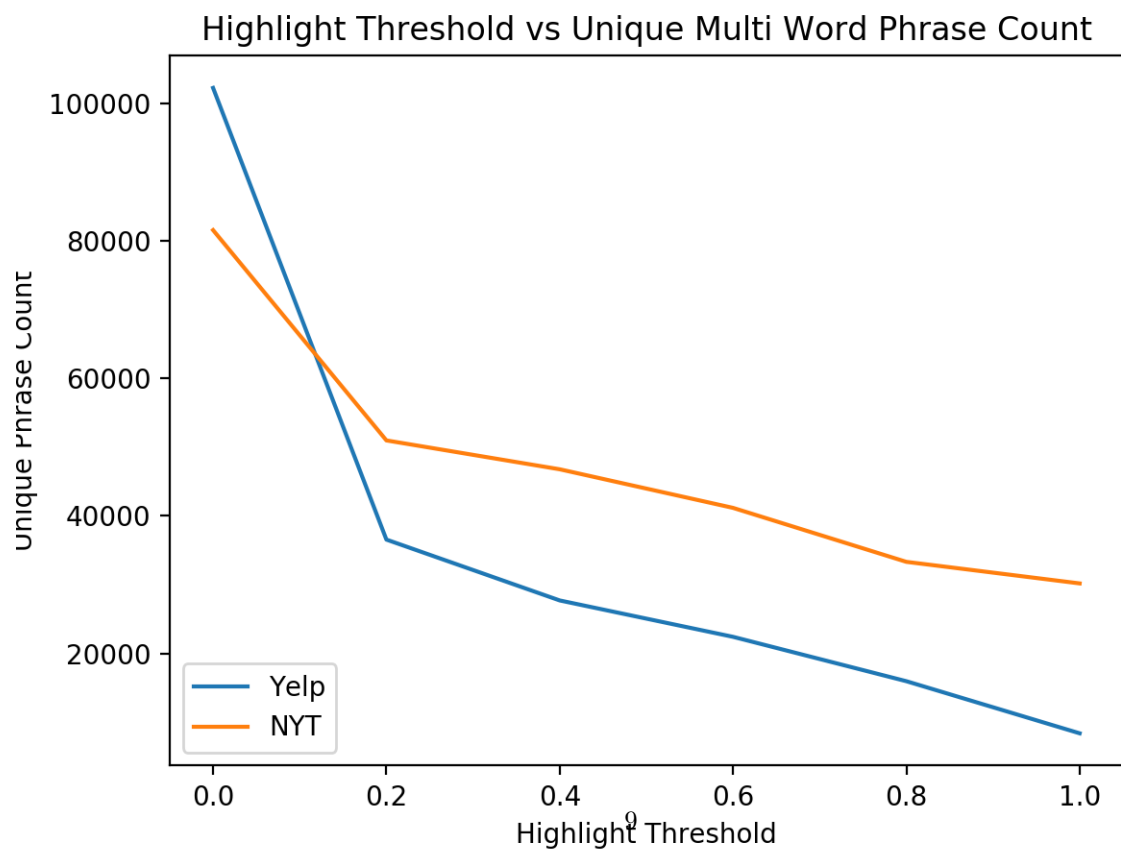
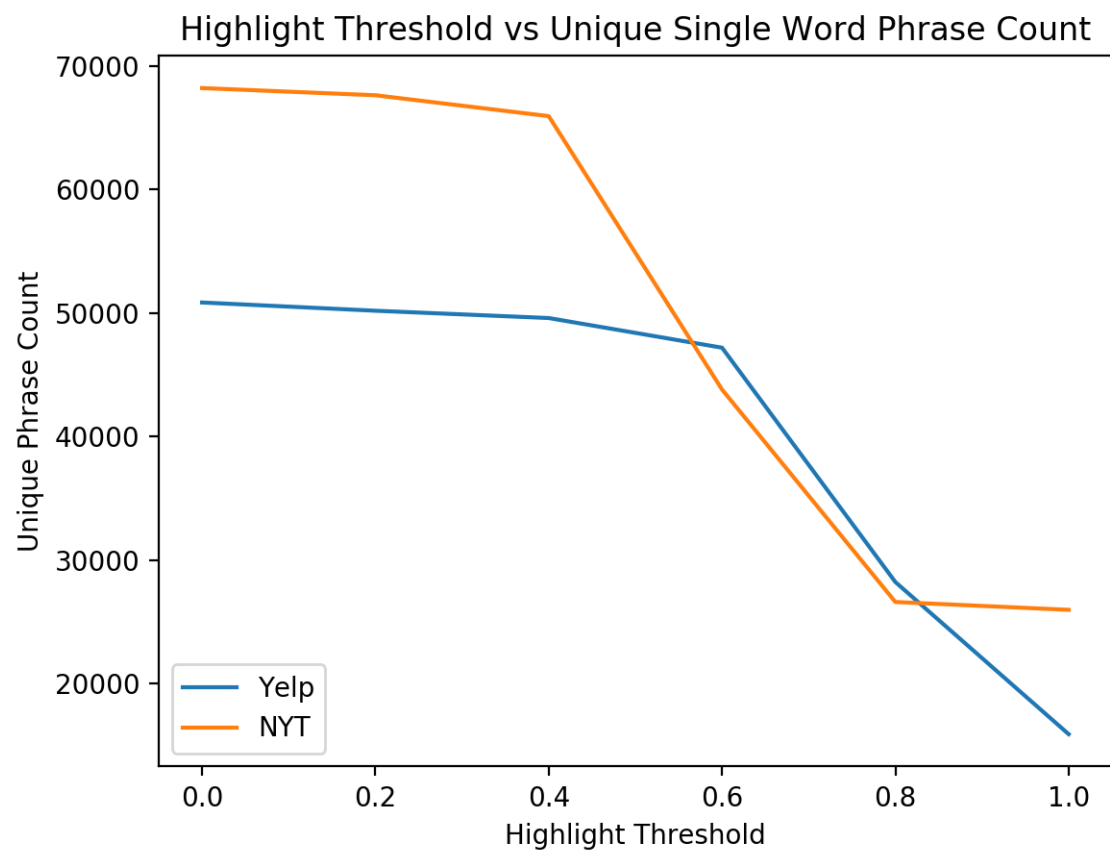
Table 25: Grocery/Department Stores

AutoPhrase Parameters

The images below show plots of the number of unique phrases returned by AutoPhrase with different highlight thresholds. The top figure shows the number of unique single word phrases and the bottom figure shows the number of unique multi word phrases.

There is a very clear inverse relationship between highlight threshold and number of unique phrases. The shapes of the graphs are very similar for both datasets with the NYT dataset yielding more phrases for both single and multi word, except with a very low threshold in the multi word curve. In the single word curve they both plateau for a while before dropping for a threshold in the .5-.6 range. In the multi word curve they both drop very steeply as the threshold rises from 0 to .2 before leveling off a bit.

This shows that AutoPhrase detects more phrases in the NYT dataset, presumably because it is written with correct grammar and less slang than Yelp.



NYT Clustering Parameters

Below I will show samples from 3 different clustering results on the NYT dataset. Each row is from the same clustering iteration and has the same number of centers, the centers are stated in the captions.

The first row of tables shows k=5, second row shows k=20, and the third row shows k=40. You can see the clusters become more coherent as you read through each level. For example, at k=5 the Justice/Politics cluster mixes police, family and politics while a similar cluster at k=20 (U.S Government/Foreign Entities) provides more closely related phrases related to the U.S and foreign policy. Furthermore, there is a k=40 cluster (U.S Government) that narrows this down to only the U.S.

| |
|-----------------|
| _government_ |
| _U.S._ |
| _state_ |
| _United_States_ |
| _country_ |
| _long_ |
| _public_ |
| _city_ |
| _American_ |
| _political_ |

Table 26: k=5: Government/Politics

| |
|-----------|
| _game_ |
| _season_ |
| _play_ |
| _games_ |
| _lead_ |
| _led_ |
| _lost_ |
| _shot_ |
| _free_ |
| _history_ |

Table 27: k=5: Sports/Games

| |
|--------------|
| _police_ |
| _family_ |
| _law_ |
| _court_ |
| _federal_ |
| _death_ |
| _party_ |
| _Republican_ |
| _South_ |
| _President_ |

Table 28: k=5: Justice/Politics

| |
|-----------------|
| _United_States_ |
| _president_ |
| _China_ |
| _European_ |
| _President_ |
| _Iran_ |
| _Europe_ |
| _Russia_ |
| _trade_ |
| _France_ |

Table 29: k=20: U.S Government/Foreign Entities

| |
|------------------|
| _city_ |
| _military_ |
| _security_ |
| _French_ |
| _British_ |
| _violence_ |
| _war_ |
| _corruption_ |
| _United_Nations_ |
| _Afghanistan_ |

Table 30: k=20: Military/War

| |
|-------------------|
| _company_ |
| _business_ |
| _director_ |
| _research_ |
| _firm_ |
| _chief_executive_ |
| _general_ |
| _organization_ |
| _Department_ |
| _Center_ |

Table 31: k=20: Business

| |
|----------------|
| _long_ |
| _high_ |
| _fell_ |
| _average_ |
| _rose_ |
| _short_ |
| _stock_ |
| _lower_ |
| _unemployment_ |
| _reading_ |

Table 32: k=40: Econ-
omy/Stocks

| |
|---------------|
| _Republican_ |
| _Senate_ |
| _House_ |
| _Republicans_ |
| _Congress_ |
| _White_House_ |
| _Texas_ |
| _Democrats_ |
| _Democratic_ |
| _Democrat_ |

Table 33: k=40: U.S Gov-
ernment

| |
|--------------------|
| _England_ |
| _captain_ |
| _Chelsea_ |
| _New_Zealand_ |
| _Premier_League_ |
| _Champions_League_ |
| _Arsenal_ |
| _Liverpool_ |
| _Barcelona_ |
| _squad_ |

Table 34: k=40: Soccer

Yelp Clustering Parameters

Below I will show samples from 3 different clustering results on the Yelp dataset. Each row is from the same clustering iteration and has the same number of centers, the centers are stated in the captions.

The first row of tables shows k=5, second row shows k=20, and the third row shows k=40. You can see the clusters become more coherent as you read through each level. For example, at k=5 there is a Miscellaneous cluster in which the phrases don't seem to be related, while the clusters at k=20 are more conceptually similar. Furthermore, there is a k=40 cluster (Decorations) that is much more specific than any cluster in k=20 or k=5.

| |
|--------------|
| _food_ |
| _restaurant_ |
| _stars_ |
| _sushi_ |
| _happy_hour_ |
| _star_ |
| _Thai_ |
| _average_ |
| _Food_ |
| _Chicken_ |

Table 35: k=5: Restau-
rants/Food

| |
|-----------|
| _order_ |
| _long_ |
| _free_ |
| _water_ |
| _music_ |
| _car_ |
| _money_ |
| _glass_ |
| _manager_ |
| _game_ |

Table 36: k=5: Miscella-
neous

| |
|--------------|
| _nice_ |
| _friendly_ |
| _bit_ |
| _atmosphere_ |
| _high_ |
| _sat_ |
| _Love_ |
| _variety_ |
| _park_ |
| _Nice_ |

Table 37: k=5: Descrip-
tions

| |
|--------------|
| _food_ |
| _restaurant_ |
| _stars_ |
| _sushi_ |
| _star_ |
| _Thai_ |
| _average_ |
| _Food_ |
| _Mexican_ |
| _Italian_ |

Table 38: k=20: Food Types

| |
|--------------|
| _sandwich_ |
| _fries_ |
| _shrimp_ |
| _steak_ |
| _sandwiches_ |
| _appetizer_ |
| _salads_ |
| _pasta_ |
| _egg_ |
| _eggs_ |

Table 39: k=20: Food Items

| |
|-------------|
| _order_ |
| _friendly_ |
| _friends_ |
| _family_ |
| _sat_ |
| _manager_ |
| _party_ |
| _bartender_ |
| _talk_ |
| _chef_ |

Table 40: k=20: Food Service

| |
|--------------------|
| _friendly_ |
| _manager_ |
| _customer_service_ |
| _bartender_ |
| _talk_ |
| _company_ |
| _chef_ |
| _professional_ |
| _management_ |
| _super_friendly_ |

Table 41: k=40: Customer Service

| |
|--------------|
| _sandwich_ |
| _fries_ |
| _bread_ |
| _steak_ |
| _sandwiches_ |
| _bacon_ |
| _salads_ |
| _eggs_ |
| _onion_ |
| _sausage_ |

Table 42: k=40: Food Items

| |
|---------------|
| _art_ |
| _center_ |
| _brand_ |
| _dining_room_ |
| _furniture_ |
| _bar_area_ |
| _private_ |
| _stock_ |
| _bakery_ |
| _design_ |

Table 43: k=40: Decorations