



Stellenbosch
UNIVERSITY
IYUNIVESITHI
UNIVERSITEIT

NLP 

[Dr. Hanjo Odendaal]

First things first, lets install the necessary packages that we will need in this lesson as well as create the necessary files to fit everything in a `targets` pipeline:

Step 1: Install Packages

```
install.packages("tidytext", Ncpus = 4)
install.packages("quanteda", Ncpus = 4)
install.packages("topicmodels", Ncpus = 4)
install.packages("ggwordcloud", Ncpus = 4)
```

Step 2: Create Folder for Analysis

```
Texevier::create_template(directory = "~/Downloads"
                          template_name = "NLP")
```

Step 3: Create `Readme.qmd` file

```
---
title: "My readme"
format: html
editor_options:
  chunk_output_type: console
execute:
  echo: false
  eval: false
---
```

```
library(tidyverse)
library(targets)
library(dbbasic)
library(tidytext)
library(quanteda)
library(topicmodels)
library(ggwordcloud)
```

Step 4: Create .Renviron file

```
usethis::edit_r_environ()
```

```
gp_data = datascience  
gp_user = datascience  
gp_pass = f5VPEC8nsU01QKbSxSfv  
gp_host = localhost  
gp_port = 3000
```

Step 6: Restart Rstudio and Test!

```
db_query(  
  "SELECT * FROM  
  pg_catalog.pg_tables  
  WHERE schemaname = 'public';"  
  , db = "psql_datascience")
```



In this session we will be exploring the basics of Natural Language Processing:

Exploratory Analysis:

Concordancing

- Extraction of words from a given text or texts conditional on a context window

Ngrams

- Basic word counts from texts
- Creating word clouds from the text

Modeling:

Sentiment Analysis

- Using dictionary methods

Topic Modeling (Bonus)

- Text > Tokens > Document Frequency Matrix > Topic Model

First things first: Data

Our data is currently stored in a cloud Database. Pull the **sales** data into your session using the `dbbasic` package. Remember to open a `dev.sql` file to test your queries from inside R!:

- Quick look at what the database table contains again:

```
-- !preview conn=db_connect(db = "psql_datascience")  
  
SELECT * FROM gumtree_clean LIMIT 10;
```

type	ad_id	ad_url	location	for_sale_by	dwelling_type	bedrooms	bathrooms	size_sqm	parking	price	available_from	for_rent_by	furnished	smoking	pet_friendly	description
sales	1b4d044419256d4999586eae6baca8c	https://www.gumtree.co.za/a-houses-flats-for-sale/b...	bantry_bay_atlantic_seaboard	agency	apartment	3	3	546	garage	39995000	NA	NA	NA	NA	NA	This incredible property in a reno
sales	991f741a8f3619089ba10cae7d5b4072	https://www.gumtree.co.za/a-houses-flats-for-sale/b...	bantry_bay_atlantic_seaboard	agency	apartment	1	2	211	garage	11500000	NA	NA	NA	NA	NA	RE/MAX Living operates in terms
sales	f9a57d1d6310bc561da6d86ec3fb6f53	https://www.gumtree.co.za/a-houses-flats-for-sale/b...	bantry_bay_atlantic_seaboard	agency	apartment	2	2	118	covered	16995000	NA	NA	NA	NA	NA	Asking Price: R 16 995 000Nestle
sales	8b3c83c89488ee15eb5a441360ef44d	https://www.gumtree.co.za/a-houses-flats-for-sale/b...	bantry_bay_atlantic_seaboard	agency	house	3	3	589	garage	19500000	NA	NA	NA	NA	NA	RE/MAX Living operates in terms
sales	0a0f80664100990b7ebc40169081fd79	https://www.gumtree.co.za/a-houses-flats-for-sale/b...	bantry_bay_atlantic_seaboard	agency	house	1	1	138	covered	14500000	NA	NA	NA	NA	NA	RE/MAX Living operates in terms
sales	812555990c3cede1ccf24b07fc85800a	https://www.gumtree.co.za/a-houses-flats-for-sale/b...	bantry_bay_atlantic_seaboard	agency	apartment	2	2	118	garage	14600000	NA	NA	NA	NA	NA	RE/MAX Living operates in terms
sales	c2c7c97bd2cc996cee5b0821177a259	https://www.gumtree.co.za/a-houses-flats-for-sale/b...	bantry_bay_atlantic_seaboard	agency	apartment	4	3	301	covered	12995000	NA	NA	NA	NA	NA	Exclusive Mandate Asking Price:
sales	cc71a7903344e3d7debca37f13a48d01	https://www.gumtree.co.za/a-houses-flats-for-sale/b...	bantry_bay_atlantic_seaboard	agency	apartment	2	2	118	covered	16995000	NA	NA	NA	NA	NA	Aurum Luxury Residences introdu
sales	844b53eb84712852ce4241ddca268f2c	https://www.gumtree.co.za/a-houses-flats-for-sale/b...	bantry_bay_atlantic_seaboard	agency	apartment	1	2	69	NA	4650000	NA	NA	NA	NA	NA	RE/MAX Living operates in terms
sales	2dd9514800082e2e3ce91f973151a338	https://www.gumtree.co.za/a-houses-flats-for-sale/b...	bantry_bay_atlantic_seaboard	agency	apartment	1	1	NA	NA	4650000	NA	NA	NA	NA	NA	Nestled along the picturesque Be

- What are the different variables within the `type` column?

```
SELECT
  type,
  COUNT(*)
FROM gumtree_clean
GROUP BY
  type
LIMIT 10;
```

- Filter the dataset to only bring in **sales** into R environment:

```
gumtree ← db_query(" SELECT * FROM gumtree_clean WHERE type = 'sales'",
                   db = "psql_datascience")

gumtree_texts ← gumtree %>% select(ad_id, description) %>%
  mutate(description = tolower(description))
```

Case Study: 'Ocean'

Concordancing

As mentioned, "concordancing" is where we would like to see what how words are used in context, or "keywords-in-contexts" analysis. Let us use the `kwic` function from `quanteda` to see how the word **ocean** is used in the text.

```
ocean_kwic ← kwic(
  # define text
  gumtree_texts$description,
  # define search pattern
  pattern = "ocean",
  # define context window size
  window = 5) %>%
  as_tibble

# [main]>ocean_kwic
# A tibble: 1,772 × 7
#   docname from to pre keyword post pattern
#   <chr> <int> <int> <chr> <chr> <chr> <fct>
# 1 text1 16 16 has views overlooking the open ocean , robben island and on ocean
# 2 text3 116 116 stunning vistas of the azure ocean that greet you from every ocean
# 3 text5 134 134 unobstructed views of the atlantic ocean from every angle along with ocean
# 4 text7 85 85 panoramic views of the atlantic ocean and lion's head , emphasising ocean
# 5 text7 281 281 views of the majestic atlantic ocean . designed with entertainment in ocean
# # i 1,762 more rows
# # i Use `print(n = ...)` to see more rows
```

We can also extract exact phrases, by using the function `phrase()`:

```
ocean_kwic ← kwic(
  # define text
  gumtree_texts$description,
  # define search pattern
  pattern = phrase("blue ocean"),
  # define context window size
  window = 5) %>%
  as_tibble

# [main]>ocean_kwic
# A tibble: 6 × 7
#   docname   from   to pre                keyword   post                pattern
#   <chr>     <int> <int> <chr>                <chr>     <chr>                <fct>
# 1 text5849    78    79 the view of the wide  blue ocean stretching to the horizon .  blue ocean
# 2 text6106   124   125 with sliding windows overlooking the blue ocean ! with its stunning ocean  blue ocean
# 3 text6160   100   101 up to the most bright blue ocean views . neutral colour palette  blue ocean
# 4 text12564   67    68 by the beauty of the blue ocean shimmering through glass windows and blue ocean
# 5 text12605  255   256 onto table mountain with the blue ocean at its feet , complete  blue ocean
# 6 text14439   42    43 sands beach with its beautiful blue ocean right on your doorstep ?  blue ocean
```

Concordancing

Exercise time!

- How many instances are there of "swimming pool"?

05:00

Ngrams forms the basis of most text analysis. It is the fundamentals of *tokenization* or breaking up texts into words or sequences of words. Lets take our ocean example a little bit further and analyse the top words (after *stopwords*) within the contexts of the 'ocean' in property ads.

We are going to use a really nice function from `tidytext` called `unnest_tokens` for this.

```
ocean_kwic ← kwic(gumtree_texts$description, pattern = "ocean",  
                 window = 5) %>% as_tibble
```

- Combine the `pre` and `post` columns into one:

```
ocean_pre_post ← ocean_kwic %>% unite("text", c(pre, post)) %>% select(docname, text)
```

```
# A tibble: 1,772 × 2  
#   docname text  
#   <chr>   <chr>  
# 1 text1   has views overlooking the open_, robben island and on  
# 2 text3   stunning vistas of the azure_that greet you from every  
# 3 text5   unobstructed views of the atlantic_from every angle along with  
# 4 text7   panoramic views of the atlantic_and lion's head , emphasising
```


Now that we have our text ready, lets create `unigrams` or single word tokens from the text.

```
ocean_tokens ← ocean_pre_post %>% unnest_tokens(input = text, output = word, n = 1)
```

```
# A tibble: 14,787 × 2  
#   docname word  
#   <chr>   <chr>  
# 1 text1   has  
# 2 text1   views  
# 3 text1   overlooking  
# 4 text1   the
```

⚠ We need to now get rid of `stopwords` or words like 'the', 'a' or 'on' as these do not add contextualization.

```
ocean_tokens ← ocean_tokens %>%  
  mutate(word = gsub("_", "", word)) %>% anti_join(stop_words, by = join_by(word))
```

```
# A tibble: 8,890 × 2  
#   docname word  
#   <chr>   <chr>  
# 1 text1   views  
# 2 text1   overlooking  
# 3 text1   open  
# 4 text1   robben  
# 5 text1   island
```

Plotting the word clouds

Once the text is in a nice tidy format, we can now do a lot with it... first lets plot the clouds to see what are the words closely associated around the 'ocean':

```
ocean_tokens %>%  
  count(word, name = "obs", sort = TRUE) %>%  
  sample_frac(weight = obs, size = 0.1) %>%  
  ggplot(., aes(label = word, size = obs,  
                color = obs)) +  
  geom_text_wordcloud() +  
  scale_color_gradient(low = "#189bcc",  
                       high = "#960018") +  
  scale_size_area(max_size = 20) +  
  theme_minimal()
```



Plotting Bi-Grams (Bonus)

```
ocean_pre_post %>%
  unnest_tokens(input = text, output = word,
               token = "ngrams", n = 2) %>%
  mutate(word = gsub("_", "", word)) %>%
  separate(word, c("word1", "word2"), sep = " ") %>%
  filter(!word1 %in% stop_words$word) %>%
  filter(!word2 %in% stop_words$word) %>%
  unite(bigram, word1, word2, sep = " ") %>%
  count(bigram, name = "obs", sort = TRUE) %>%
  sample_frac(weight = obs, size = 0.1) %>%
  ggplot(., aes(label = bigram, size = obs,
               color = obs)) +
  geom_text_wordcloud() +
  scale_color_gradient(low = "#189bcc",
                      high = "#960018") +
  scale_size_area(max_size = 20) +
  theme_minimal()
```



Modeling

Sentiment analysis in action

Provided that you now have the basics of text analysis and how to get from text > tokens, we can now apply some basic modeling techniques.

Sentiment analysis is a very nice kick-off point as it ranges in complexity from basic dictionary techniques (what we will be using) to intricate deep learning models. The following is an illustrative example of how sentiment is calculated using a news article and a dictionary approach:

Global capital flows into emerging markets including SA are drying up. Andrew McNulty looks at how the changing game will affect investors. For most of 2010 and 2009 SA and other emerging markets were popping champagne corks but money was flowing into emerging markets. Global capital flows helped lift prices of these countries equities and other assets. Their currencies and underpinned growth but the capital flows have dwindled or turned since last October. About relative values using inflation in these economies occurred in North Africa and the Middle East and surging energy prices have turned investors' emerging markets that include SA. JSE equities residence and the rand has fallen 16 dollar after the Reserve Bank's monetary policy committee meeting last Thursday. Governor Gill Marcus said SA had experienced net sales of bonds and equities by own residents in the year to date totalling R102bn. As the world economic recovery momentum and bonders' inventory views on the JSE and global markets have changed. Marcus says the global recovery appears to have remained on track but the European debt crisis oil price and events in Japan may moderate the pace of recovery in the near term. Global inflows have increased particularly in emerging market economies some of which have raised short-term interest rates. The dwindling of foreign capital inflows has not led to steep falls in emerging equity markets but these markets are no longer rising and many have fallen. The developed markets over the past six months. The MSCI Emerging Markets Group economist Joe Lauer says capital flows not necessarily keep flowing into emerging markets at a steady rate. The inflows depend on market circumstances he says. I see the recovery as a tactical or cyclical response to changes in relative values. The long-term trend should continue. However developed markets may be expected to see growth resumes. For a while the recovery from the crisis has been uneven. They could be emerging market equities or bonds. Returns that they could expect from developed markets where growth remains solid. Other high markets such as the JSE were favourably re-rated on recovering domestic growth. Emerging company earnings and rising commodity prices. Rising emerging market currencies. Returns for foreign institutional investors indicate that the game has changed. In the past continued high growth and developed economies have returned to growth. Company earnings have risen and equity markets in developed economies may now offer better value. Does the change in capital flows affect the growth outlook for emerging markets particularly SA and the investor? The reversal of capital flows and the impact of other developments such as the recent spike in the oil price. Capital flows into emerging markets after the financial crisis mainly because of their high growth and structural. Higher returns in many countries favourable economic fundamentals such as low inflation, supplies of labour, in current accounts large foreign exchange reserves and an expanding middle class should support long-term growth. The real returns are about their short-term outlook and relative value which no longer looks attractive. Developed economies if capital outflows continue that could be a growth outlook and increase. For some countries capital flows to developing economies slipped in 2010 but have been in upward trend since early in the past decade after the recovery from the 1997-1998 Asian crisis. The Institute of International Finance IIF says net private capital inflows into emerging markets rose from less than US\$200bn or below 1% of GDP in 2002 to 900bn or about 7% of GDP in 2010. An IIF report in January forecast the figure would rise to 960bn this year and to just over 1trillion in 2021. Predicted growth in portfolio equity inflows in the rest of an emerging equity market. Increased and credit measures curbed credit growth. It identified two main factors that markets in emerging economies could make such as the drying up of avoiding monetary or fiscal policies to stall currency appreciation making. Must cycle more likely. And developed economies could start raising rates sooner than expected reducing the relative attractiveness of emerging market yields. That may be slowing the Eurozone Centre. Expanded earlier this month it may raise rates at its next meeting. Nonetheless the IIF forecasts continuing capital inflows mainly on the prospect of near-term growth. Developing nations particularly in America and Asia have other attractions including increased public infrastructure investments. In the past decade Asian countries have seen sharp reversal of capital inflows as well as the Asian crisis. Almost 18 years ago International Monetary Fund IMF research shows net inflows to the Asian countries were about 6% of their GDP in 1993 and 5% in 1996. In 1997 net outflows were 2% of GDP. In 1998 they rose to 2% of GDP. IIF notes that growth over years of fast growth. In Asian countries including Korea Indonesia and Thailand particularly. IIF was wary of a capital flight. Much of the capital inflows were lending. After a huge expansion of credit had led to a boom. Tighter regulation of banking and financial sectors. IIF says to build institutions in the sector. IIF asserts to start through. However were slightly elevated. There was a common theme in Indonesia Korea and Thailand. Despite the lack of supervision they had a history of heavy government involvement in credit allocation. These were reflected in lending practices and inflated asset prices. There were also financial regulations. Some countries had large current account or trade surpluses but maintained pegged exchange rates for too long. Foreign lenders and investors contributed to asset bubbles as the IIF to evaluate. After Thailand's debt crisis under speculative pressure in the first half of 1997 that year. The government had to raise the currency as well as a managed float. The Thai government's exchange rate. Not among other countries in the region. IIF says that the currency should be allowed to float. IIF also had large foreign exchange reserves. IIF could a similar pattern be repeated. Asian economies have again experienced almost a decade of growth. China and other Asian countries have avoided measures aimed at curbing the rise in property prices. Most of the high-growth developing countries report or managed to maintain. IIF also had large foreign exchange reserves. In October 2010 Global Financial Crisis Report the IMF said rating agencies had downgraded 27 developing country sovereign debt issues 23 times since early 2000. In 2010 there were upgrades of emerging market sovereign debt. The IMF predicts the trend to continue. A potential source of risk is the carry trade that occurs when investors borrow money at low rates and invest it in assets in countries where rates are higher. Australia and SA have been favoured destinations for Japanese investors see story on page 22. Rates start rising in Europe. The US or Japan it could make the carry trade less attractive. However countries such as Brazil India and Australia have raised short-term rates many times in the past year. SA's rates remain at 30-year lows though at the bottom of the cycle. SA is becoming a yield play relative to some other emerging markets. Says Frank for Economic Research. BEI senior economist Hugo Penner. The BEI said last month financial markets had been overly aggressive in pricing in a local rate increase in SA before the end of this year. The BEI forecasts a first 50 basis point increase in November followed by more increases totaling 100 basis points in the first quarter of 2012. IIF forecasts capital flows will be on a rise. The IIF says inflation in emerging markets is expected to increase 47% in 2011 and 50% in 2012. However a combination of continued capital outflows high oil prices and a higher rand if that occurs may change that view. In its economic update earlier this year the IMF forecast emerging and developing economies would grow at 4.5% in 2011 and 2012. IIF is developing the forecast for 2012. IIF forecasts a 2% increase in the forecast. The BEI forecasts SA will grow at 5% this year and 3% in 2012. Many analysts believe the long-term outlook is compelling. Global Standard Chartered chief economist Gerard Lyons says the world may have entered a new supercycle driven by industrialisation and urbanisation of emerging markets and global trade. China's growth rate could average 9% over the next two decades. In 2010 and India may grow faster. China's income per capita still is only about that of the US in 2010 leaving room for catch-up. Emerging markets now account for almost 40% of global consumption and more than two-thirds of global growth. However, the carry trade policies that led to asset price bubbles in the past two decades. IIF to avoid asset price bubbles it would have to raise interest rates on these economies that would deal a blow to the global recovery. Its little wonder that foreign investors have turned their backs on SA and value while these countries adjust to the effects of overvalued and inflated

From the illustration, we can pick out that the article contains 93 negative words and 42 positive words:

$$A_{it} = \frac{\text{Positive Words} - \text{Negative Words}}{\text{Positive Words} + \text{Negative Words}}$$
$$-0.37 = (42 - 93) / (42 + 93)$$

This news article is thus deemed to be mostly negative. Beware, there are directional biases within sentiment analysis and its a good idea to normalise the scores for analysis.

Sentiment analysis in action

Now that we have then basic idea of sentiment analysis using a dictionary approach, lets answer the following question:

- Are the ads for houses more positive than apartments?
- Is there a correlation between price and sentiment?

To conduct this analysis we have to perform the following steps:

- 1) Transform our text into tokens
- 2) Remove stopwords
- 3) Join in a sentiment dictionary
- 4) Group by advert type and analyse

Sentiment analysis: Houses vs Apartments

We start off by turning our text into tokens using `unnest_tokens`:

```
gumtree %>% select(dwelling_type, description) %>%  
  unnest_tokens(word, description) %>%  
  anti_join(stop_words, by = "word")
```

Next, lets use the `bing` sentiment dictionary to determine the relative sentiment per advert:

```
tidytext::get_sentiments(c("bing", "afinn", "loughran", "nrc")[1])  
# A tibble: 10 × 2  
#   word      sentiment  
#   <chr>    <chr>  
# 1 trouble-free positive  
# 2 displaced negative  
# 3 imposers negative  
# 4 achievable positive  
# 5 temptation negative  
# 6 cliché negative
```


Sentiment analysis: Houses vs Apartments

Create the sentiment score per ad:

```
gumtree_sentiment <- gumtree %>%
  select(ad_id, dwelling_type, description) %>%
  unnest_tokens(word, description) %>%
  anti_join(stop_words, by = "word") %>%
  left_join(get_sentiments("bing"), by = "word") %>%
  drop_na() %>%
  count(ad_id, dwelling_type, sentiment, name = "obs") %>%
  pivot_wider(names_from = "sentiment", values_from = "obs",
              values_fill = 0) %>%
  mutate(sentiment = (positive - negative)/(positive + negative))

# [main]>gumtree_sentiment
# A tibble: 14,562 × 5
#   ad_id                dwelling_type positive negative sentiment
#   <chr>                <chr>          <int>    <int>    <dbl>
# 1 00026b744459e17f11d5d66a9634f159 apartment         3         0         1
# 2 000513d19f7999faf08868e98d1a5dde house             12         0         1
# 3 0007b1cbdc9ef27ad3d663aaa5a11240 house              9         2         0.636
# 4 000a2cd6656b24763bd1fc416ea01b00 house              9         1         0.8
# 5 000b5a78eeb86987a0f8f4fba68fd568 house              3         1         0.5
```


Sentiment analysis: Houses vs Apartments

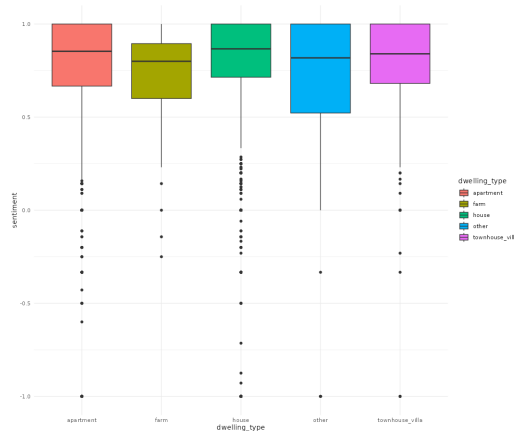
Using the sentiment data frame we are now able to test whether adverts for houses are more positive than apartments. To do this, we turn to stats 

- Boxplots and Density plots for visual
- Wilcoxon test for same continuous distribution (non-parametric version of a t-test)

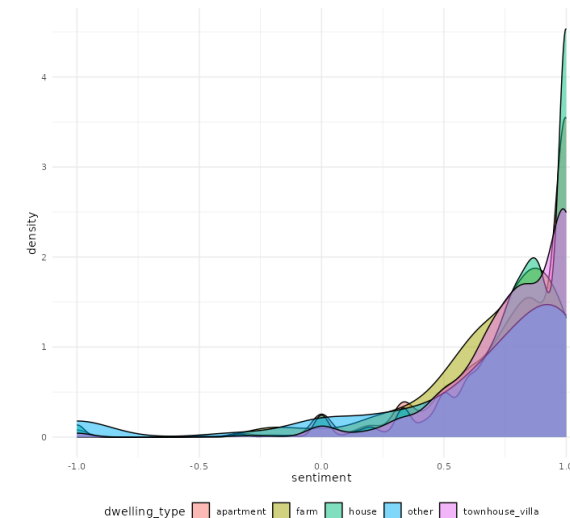
Sentiment analysis: Houses vs Apartments

Lets start with visual inspections:

```
gumtree_sentiment %>%  
  ggplot(., aes(dwelling_type, sentiment,  
               fill = dwelling_type)) +  
  geom_boxplot() +  
  theme_minimal()
```



```
gumtree_sentiment %>%  
  ggplot(., aes(sentiment, fill = dwelling_type)) +  
  geom_density(alpha = 0.5) +  
  theme_minimal() +  
  theme(legend.position = "bottom")
```



Sentiment analysis: Houses vs Apartments

It is good that we performed the visual inspection as one would have seen two things:

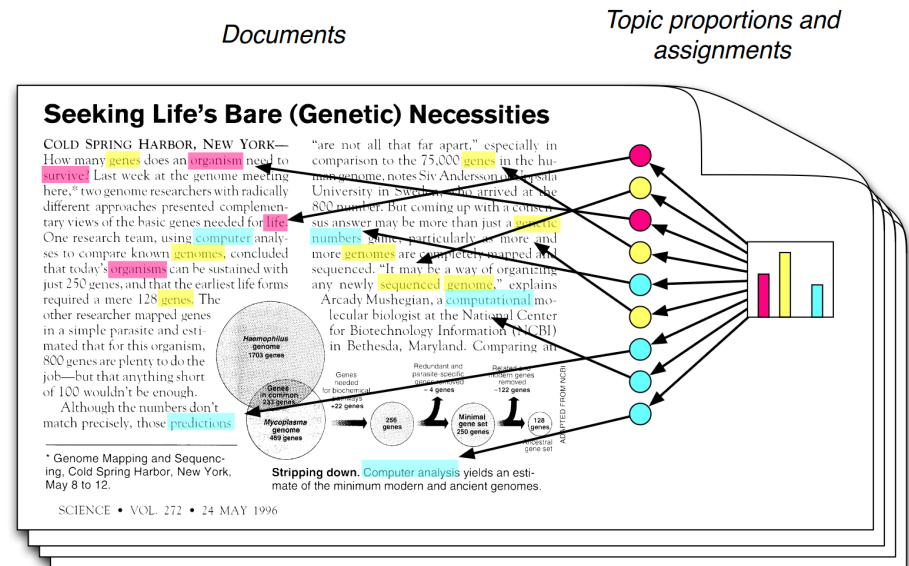
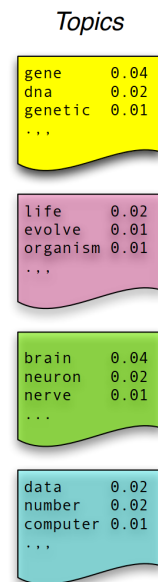
- The data is *not* normally distributed and as thus, we cannot use a *parametric* t-test
- There are other categories which we have to filter out

$$H_0 : S_{apartment} \geq S_{house}$$

```
gumtree_sentiment_list <- gumtree_sentiment %>%  
  filter(dwelling_type %in% c("apartment", "house"))  
  
wilcox.test(sentiment ~ dwelling_type, data = gumtree_sentiment_list,  
            alternative = "less")  
  
#      Wilcoxon rank sum test with continuity correction  
#  
# data:  sentiment by dwelling_type  
# W = 19358132, p-value = 0.0006765  
# alternative hypothesis: true location shift is less than 0
```

Topic modeling is one of the core tools within Natural Language Processing (NLP). The goal of using topic modeling, is to assist the analyst in order to better segment large pieces of text into various clusters or "topics". A single piece of text will be a mixture of various topics with (hopefully) one of the topics being a dominant feature.

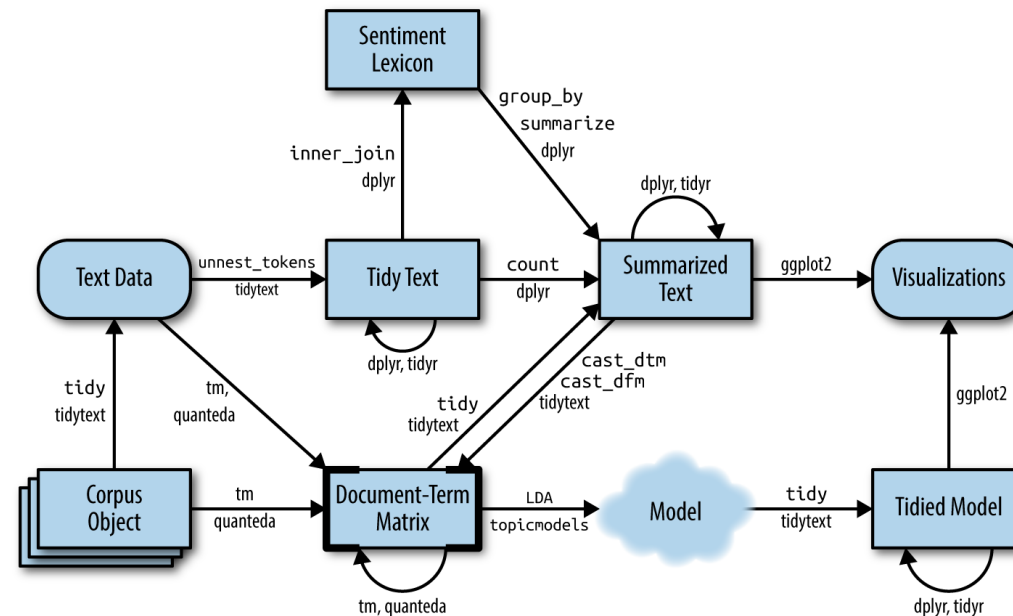
- The analyst has to make a subjective choice on the number of cluster
- Every document is a mixture of topics
- Every topic is a mixture of words



Blei, D.M., 2012. Probabilistic topic models. Communications of the ACM, 55(4), pp.77-84.

Topic Modeling

A flowchart of a text analysis that incorporates topic modeling. The topicmodels package takes a Document-Term Matrix as input and produces a model that can be tidied by tidytext, such that it can be manipulated and visualized with dplyr and ggplot2.



See <https://www.tidytextmining.com/topicmodeling>

Case Study: Farm descriptions

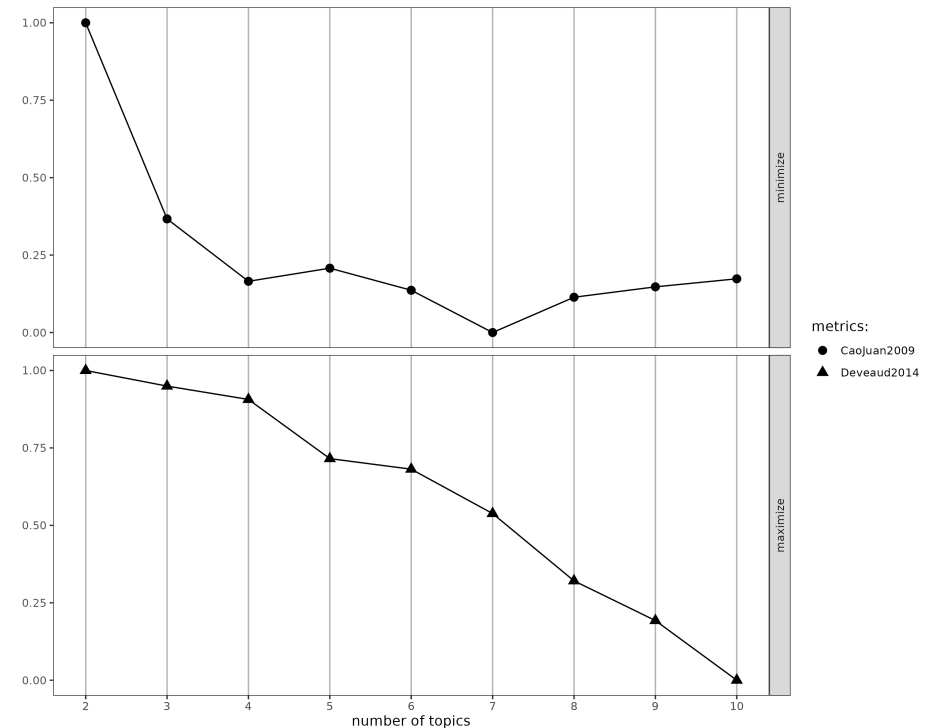
Topic models need to have some kind of design matrix: **DFM** - Document Frequency Matrix or **DTM** - Document Term Matrix. Luckily for us we already know how to get the count of terms per document!

```
gumtree_dtm ← gumtree_clean %>%  
  filter(dwelling_type = "farm") %>%  
  select(ad_id, description) %>%  
  unnest_tokens(word, description) %>%  
  anti_join(stop_words, by = join_by(word)) %>%  
  filter(!grepl("[0-9]+", word)) %>%  
  count(ad_id, word) %>%  
  cast_dtm(ad_id, word, n)  
  
# <<DocumentTermMatrix (documents: 76, terms: 3306)>>  
# Non-/sparse entries: 9306/241950  
# Sparsity           : 96%  
# Maximal term length: 34  
# Weighting          : term frequency (tf)
```

How many topics?

In order to find out what K, or number of topics should be, is a bit more of an "art" than a pure science. One can start by looking at some statistics, but they are not absolutes and you will have to use your own judgement when conducting research:

```
library(ldatuning)
result ← FindTopicsNumber(
  gumtree_dtm,
  topics = seq(from = 2, to = 10, by = 1),
  metrics = c("CaoJuan2009", "Deveaud2014"),
  method = "Gibbs",
  control = list(seed = 77),
  verbose = TRUE
)
FindTopicsNumber_plot(result)
```



We can then use the `LDA()` function to create a four-topic model. This is also mostly driven by theoretical hypothesis. My believe would be we should see: small holdings, citrus farms, wine farms and game farms... lets see if I am correct:

```
gumtree_lda ← LDA(gumtree_dtm, k = 4, control = list(seed = 1234))
gumtree_lda

# [main]>gumtree_lda
# A LDA_VEM topic model with 4 topics.
```

Now lets analyse the output:

- What words are within the topics?
- Prevalence of each topic in the corpus?

```
topics_beta ← tidy(gumtree_lda, matrix = "beta")
topics_gamma ← tidy(gumtree_lda, matrix = "gamma")
```

Understanding the topics

To understand the topics better we can analyse what words are most prevalent in a topic. This is called the `beta` matrix:

```
topics_beta

# [main]>topics_beta
# # A tibble: 13,224 × 3
#   topic term      beta
#   <int> <chr>    <dbl>
# 1     1 1 access 2.16e- 3
# 2     2 2 access 1.45e- 3
# 3     3 3 access 9.67e- 4
# 4     4 4 access 3.26e- 3
# 5     5 1 approx 5.78e-284
```

Lets use `slice_max` to find the topic 10 words per topic:

```
top_terms ← topics_beta %>%
  group_by(topic) %>%
  slice_max(beta, n = 5) %>%
  ungroup() %>%
  arrange(topic, -beta)
```

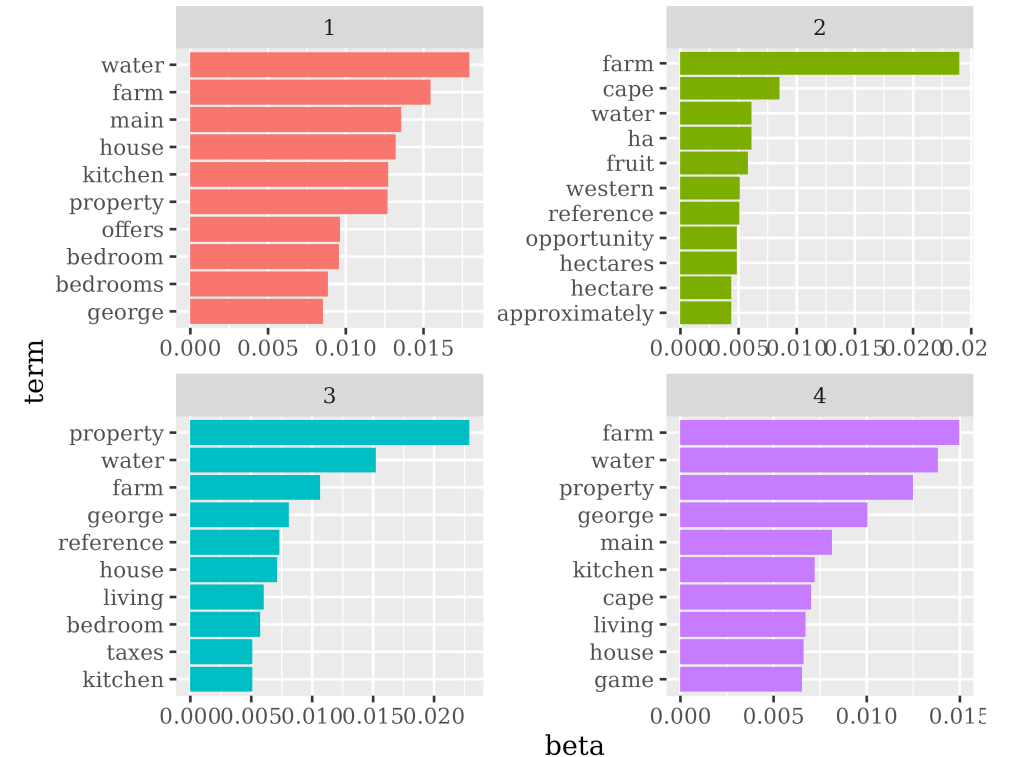
Understanding the topics

To understand the topics better we can analyse what words are most prevalent in a topic. This is called the `beta` matrix:

```
library(ggplot2)

top_terms <- topics_beta %>%
  group_by(topic) %>%
  slice_max(beta, n = 10) %>%
  ungroup() %>%
  arrange(topic, -beta)

top_terms %>%
  mutate(term = reorder_within(term, beta, topic))
  ggplot(aes(beta, term, fill = factor(topic))) +
  geom_col(show.legend = FALSE) +
  facet_wrap(~ topic, scales = "free") +
  scale_y_reordered()
```



Topic Prevalence

To analyse how prevalent a given topic is in the corpus we use the `gamma` matrix or "topic probability per document". This tells us if a certain topic dominates or not:

```
top_terms_group <- top_terms %>%  
  group_by(topic) %>%  
  slice_max(beta, n = 10) %>%  
  summarise(top_words = paste0(term, collapse = ","))  
  
topics_gamma %>%  
  group_by(topic) %>%  
  summarise(mean_gamma = mean(gamma)) %>%  
  left_join(top_terms_group) %>%  
  mutate(topic = glue("topic ({round(mean_gamma, 3)*100}%")) %>%  
  ggplot(., aes(reorder(topic, mean_gamma), mean_gamma,  
               label = top_words)) +  
  geom_col(fill = "#189bcc") +  
  geom_label() +  
  ylim(0, 0.45) +  
  labs(x = "Topic", y = "Gamma") +  
  coord_flip() +  
  theme_minimal()
```

