

1. turn data into dataframe
2. turn dataframe into corpus
3. turn into dfm, dtm, matrix dependent on model choice

remove  
urls,  
emojis  
etc.

# PREPROCESSING

= input text in right format

manipulate text features

stem or  
lemmatize

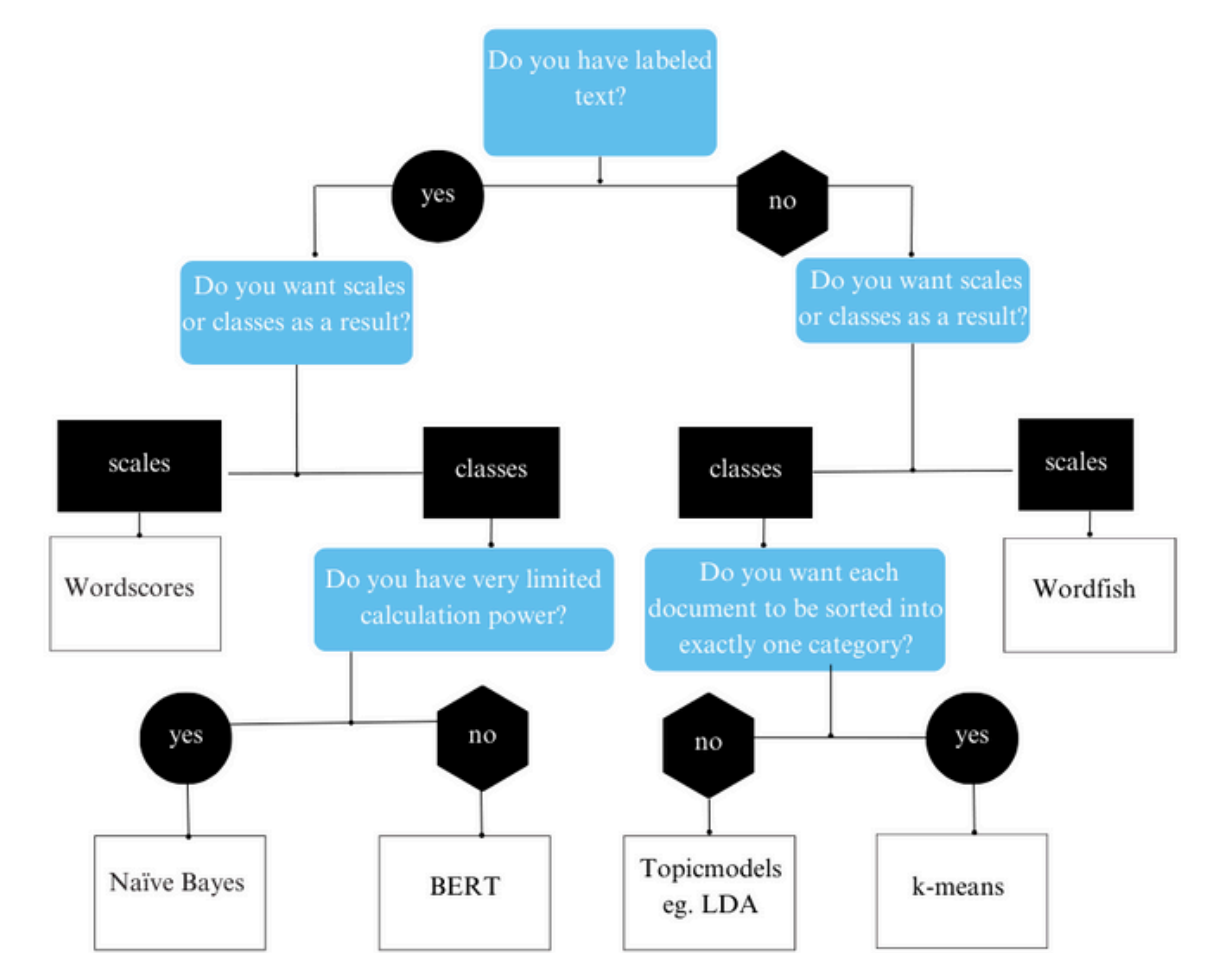
remove  
punctuation,  
numbers,  
symbols,  
stopwords

What format  
do I need?  
What features  
do I keep?

What is my  
use case,  
data format,  
resources?

research methods  
used for your goal  
data type?

# MODEL CHOICE



# MACHINE LEARNING IN QTA

iterations how  
often the  
algorithm should  
run

examples

choosing the  
number of  
topics/clusters  
"k"

algorithm e.g.  
sorting  
algorithm

# HYPERPARAMETERS

= parameters you have to  
determine

determine through  
literature  
evaluation  
methods

How can I  
optimise my  
model?

set.seed() construct validity

reliability  
-delivers the same  
result in repeated  
measurements

reliability  
of hand-  
coded text

face-, convergent-, validation

validity  
-whether a measure  
measures what it  
should

# EVALUATION & VISUALISATION

`ggplot2`  
framework

statistical-,  
semantic-, predictive