

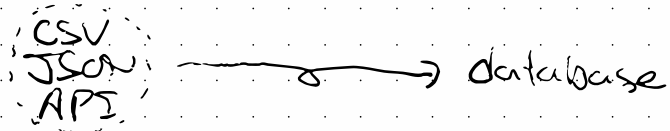
# Data Engineering

## Introduction to data engineering

- ↳ Gather data from different sources
  - ↳ Process & combine
  - ↳ Store clean data for use
- Data  
Pipelines  
↓
- ↳ typically separate from DS
  - ↳ But, DS might need to know some in smaller companies.

## ETL Pipelines

ETL → Extract Transform Load



→ Extract

- ↳ Practical:

→ Transform

- ↳ clean & combine data
- ↳ encodings
- ↳ remove duplicates
- ↳ Dummy variables / Normalize
- ↳ remove outliers
- ↳ etc.

→ Data formats

- ↳ CSV ← comma sep. values  
imported w/ pandas
- ↳ JSON ← key: value pair file
- ↳ XML ← the same
- ↳ API
- ↳ SQL ← language to communicate w/  
databases

→ While cleaning data you might  
need to combine data or parse it  
to the correct data type

→ you may also need to delete rows  
or create appropriate values where  
possible.

→ Handling outliers

- ↳ Cause, who knows? (Programmatic or human error)
- ↳ Statistical tests to remove
  - ↳ Use pca or euclidean distance to check
- ↳ corr plot can even help

→ Scaling data

- ↳ Gives each feature a fair shot

→ Feature engineering

- ↳ creating new features from old ones  
(helps w/ underfitting)  
↳ increases variance of model

(Load)

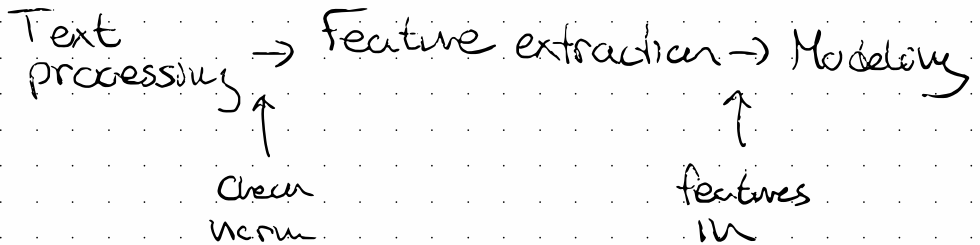
↖ Store data appropriately

Recap

- Data engineering - prepare data
- done w/ ETL

# NLP pipelines

↙ May not be linear



processing → example: website: remove HTML tags  
→ creates plain text.

example library "Beautiful Soup"

- ↳ returns page as soup
- ↳ transforms to plain text w/ "get\_text"
- ↳ possible to find tags w/ "find\_all"



plain text!

Text processing

→ Normalization

- ↳ remove variations of the same words
- ↳ remove punctuation (?)

→ Tokenization

↳ take 1 string & split into multiple strings

(NLTK - used to perform some of this cleaning)

→ Stopword removal

↳ remove words that are not important to text analysis

→ Part of speech tagging

↳ Takes string & classifies based on being noun/adverb/etc.

→ Named entity recognition

↳ can recognise person/org/GPE/etc.

→ Stemming & lemmatization

↳ takes word & removes further variation

branching  
" ed } branch  
" es }



## TF-IDF

- ↳ BOW is naive (doesn't attach importances to words)
- ↳ add 1 row to BOW w/  $\Sigma$  of word freq.
  - ↳ divide each col by number in this row for each col.

one hot enc.

- ↳ Similar to BOW only it's a word/word matrix (?)
- ↳ Doesn't work well for large documents

word embeddings

- ↳ cluster words based on similarity in 3d space.

Other

word2vec  
Glove

↳ good for deep learning

## → Modeling

- ↳ use features to make predictions for seen & unseen data.
- ↳ create a model to predict on any target you want.

→ Machine Learning Pipelines

## Automatische ML workflows

example: classify corporate messages

Step 1: clean / tokenize

Step 2: for a model

Step 3: wrap everything into a pipeline

ie pipeline performs necessary steps before fitting/predicting

↑ Cost of key, value pairs.

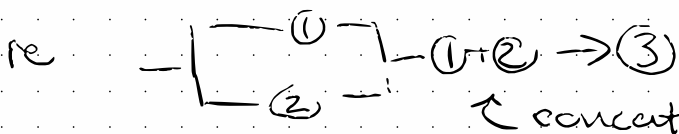
← All of these steps must be transformers (except final)

## Adv of pipelines

- Simple
- Makes workflow easier to handle/maintain

## Feature univars

- ↳ Allows for the creation of multiple features at once in a single pipeline





## Custom transformers

↳ in the case that your pipeline requires complicated transformations, you can also make custom ones.

↳ create classes that extend sklearn.base.BaseEstimator & TransformerMixin

↳ need to create fit & transform methods

↳ or use sklearn.preprocessing.FunctionTransformer which turns a function into a transformer.

## Gridsearch

✓ pseudo code

parameters = dict of params

svc = SVC()

clf = GS(svc, parameters)

clf.fit(X, y)

⚠ take care to transform data only after the split has taken place