

End to End Gen AI pipeline

① Data Acquisition:

- Available Data (CSV, txt, PDF, Docs, XLSX)
- Other Data (DB, Internet, API, Scraping)
- No Data → create your own Data
→ LLM

Note: need - 10,000 , collection - 5000

Data Augmentation

① Replace with Synonyms

- I am a Data Scientist
- I am a AI Engineer

② Bi-gram flip

I am Bappy ←
Bappy is my name ←

③ Back Translations

① add additional data/noise

(I am a Data scientist, I love
to work here)

② Data preparation:

① cleanup: HTML, emoji, spelling, correction

② Basic preprocessing

✓ ③ Advanced preprocessing →

✓ # Basic preprocessing

→ Tokenization



Optional preprocessing

✓ ① stop word removal

✓ ② stemming — less used

✓ ③ lemmatization — none used

✓ ④ punctuation removal (2, !, ., & @)

✓ ⑤ lower case

✓ ⑥ Language detection —

My name is Bary

word = ["my", "name", "is", "Bary"]

my name is Bary. I am a good boy

sen = ["my name is Bary", "I am a good boy"]

play, played, playing, plays

play

→ sports

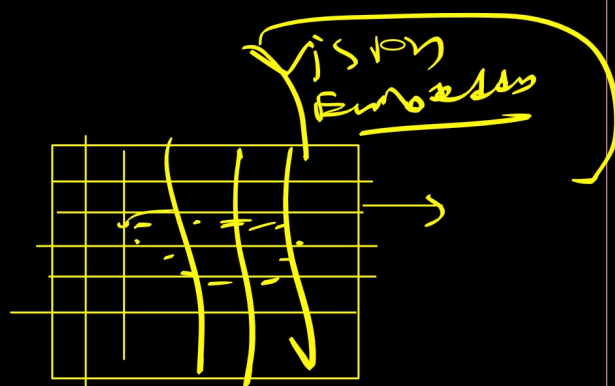
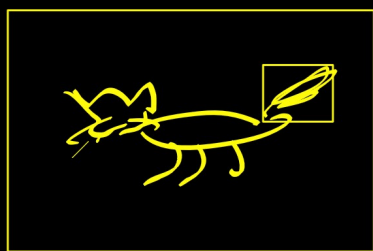
→
bary is a data scientist

bary is a good boy → ASCII

③ Feature Engineering

→ Text vectorization

→ TFIDF
→ Bag of word (BoW)
→ word2vec
→ one hot encoding
→ Transformer Embedding →



④ Modeling

→ ① Open Source

→ ② paid model → API (OpenAI)



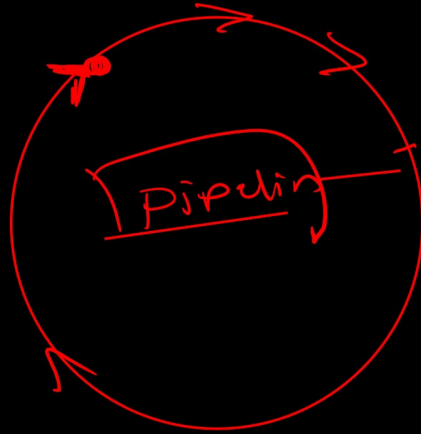
⑤ Model evaluation

① Intrinsic Eval →

② Extrinsic Eval →

from

Metric, P, R, F1, L



⑥ Deployment → Host

AWS, Azure, GCP
↑ ↑ ↑

common terms ←

- ① corpus (Entire text)
- ② Vocabulary (Unique words)
- ③ Documents (Single line)
- ④ words (single word)

Image augmentation generates random images based on existing training data to improve the generalization ability of models. In order to obtain definitive results during prediction, we usually only apply image augmentation to training examples, and do not use image augmentation with random operations during prediction.