execute ./download

For each category:
1. query category

    wiki_module.query_category

    dict

2. write to db

    database_module.insert_category

3. add page to list of page ids

    FAIL

Parse command line arguments

args

page_ids

category.yml

load yaml file if necessary

category_titles

For each page:
1. query page

    wiki_module.query_page

    dict
    html
    text
    id
    summary

2. encode page

    encoding_module.encode_page

    dict
    id:vector

3. write to db

    database_module.insert_page

    FAIL

# Pipeline 1

## download

minimum viable implementation

execute ./search

Parse command line arguments

args

encode search query

encoding_module.encode_query

dict
query:vector

select encoded pages

database_module.select_page_vectors

DataFrame    page_vectors

Identify 5 nearest vectors

list of page_ids

fetch page text

database_module.select_pages

DataFrame    pages

Pipeline 2

search

minimum viable implementation

execute ./train

select encoded pages

`database_module.select_page_vectors`

select corresponding categories

`database_module.select_category_ids`

X
(page_vectors)

page_ids

y
(category_ids)

Make data dictionary

dict of DataFrames
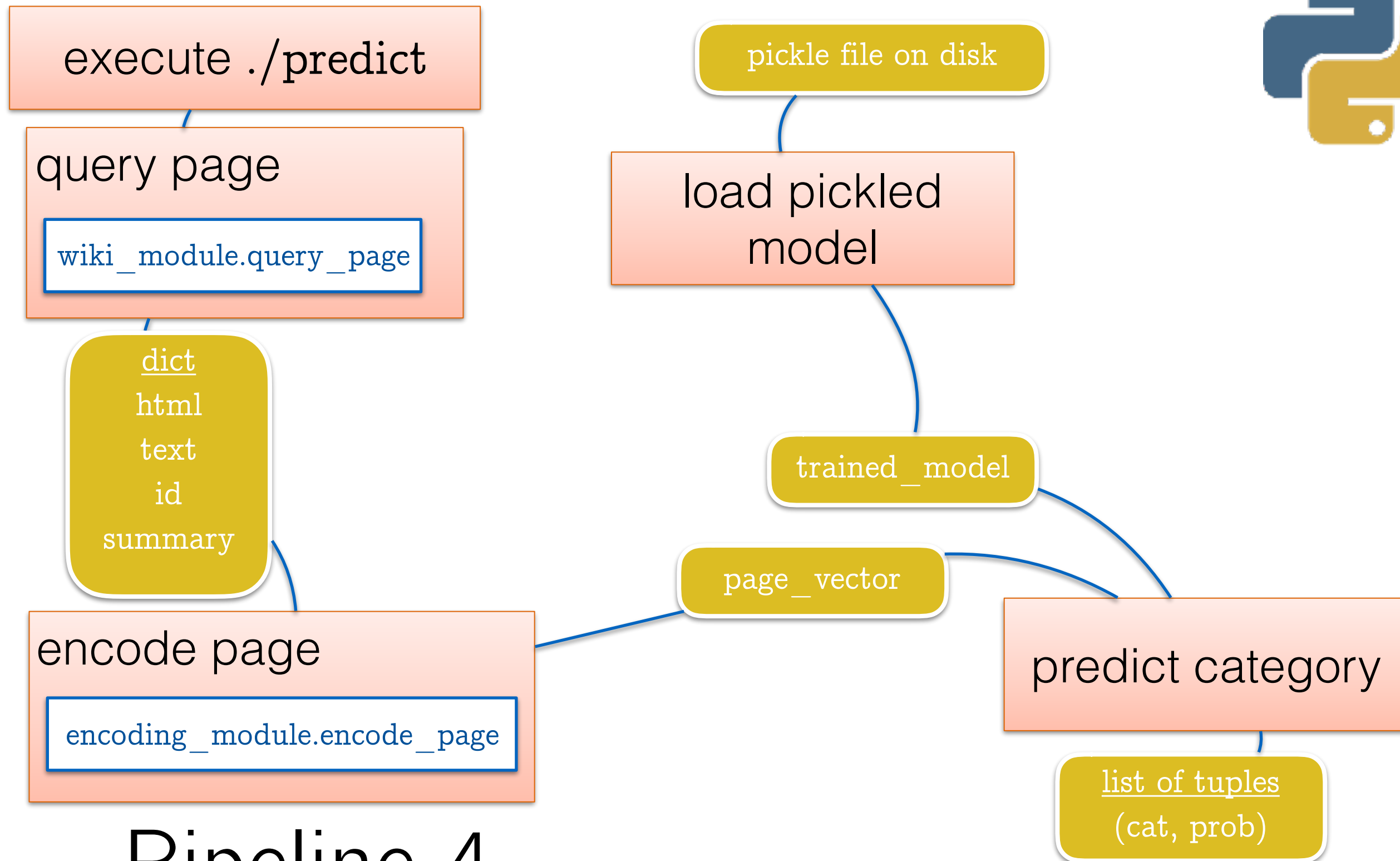
trained_model

Fit, Score, Tune Model

Pickle model

pickle file on disk

# Pipeline 3

train

minimum viable implementation

execute ./predict

query page

`wiki_module.query_page`

**dict**
html
text
id
summary

encode page

`encoding_module.encode_page`

pickle file on disk

load pickled
model

trained_model

page_vector

predict category

**list of tuples**
(cat, prob)

# Pipeline 4

predict

minimum viable implementation