

# hACKER NEWS UPVOTE PREDICTOR

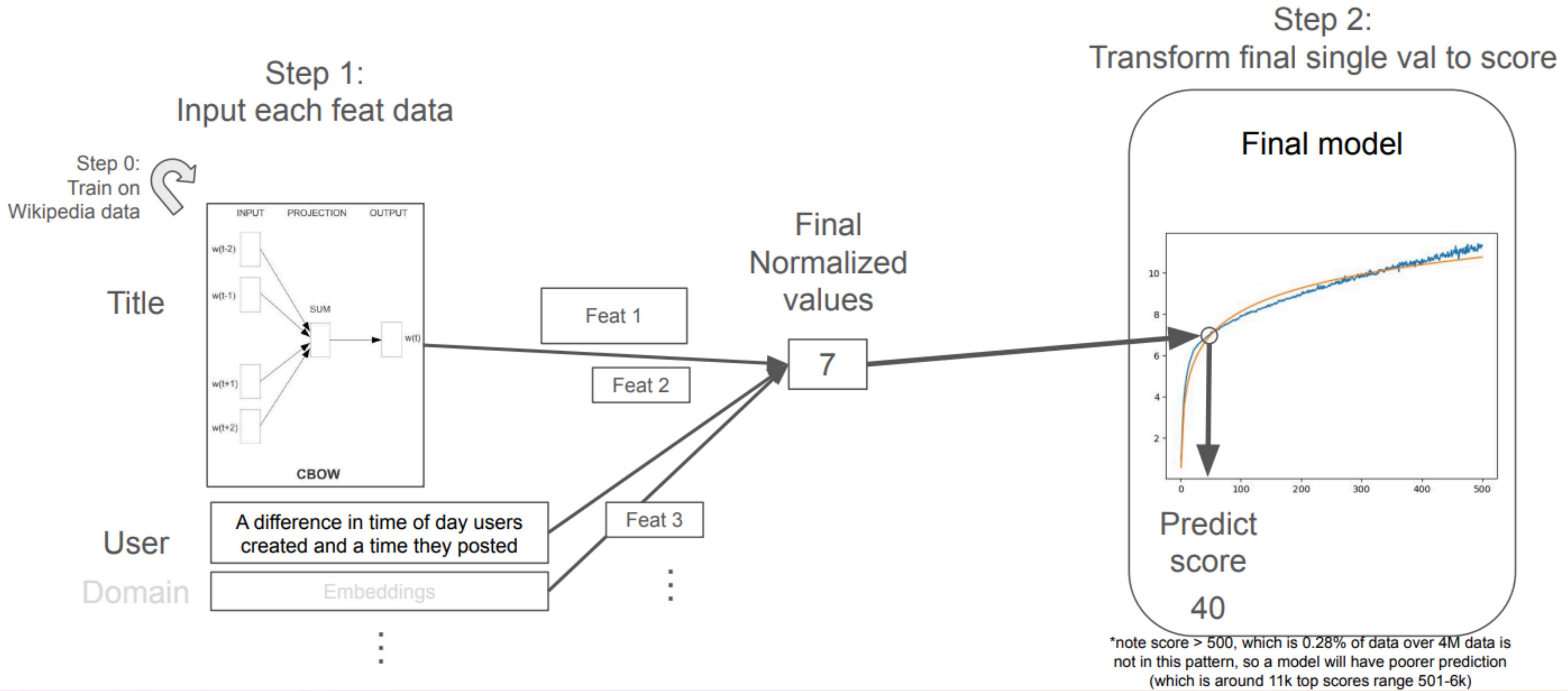
***Gradient Gigglers***

# WORKFLOW

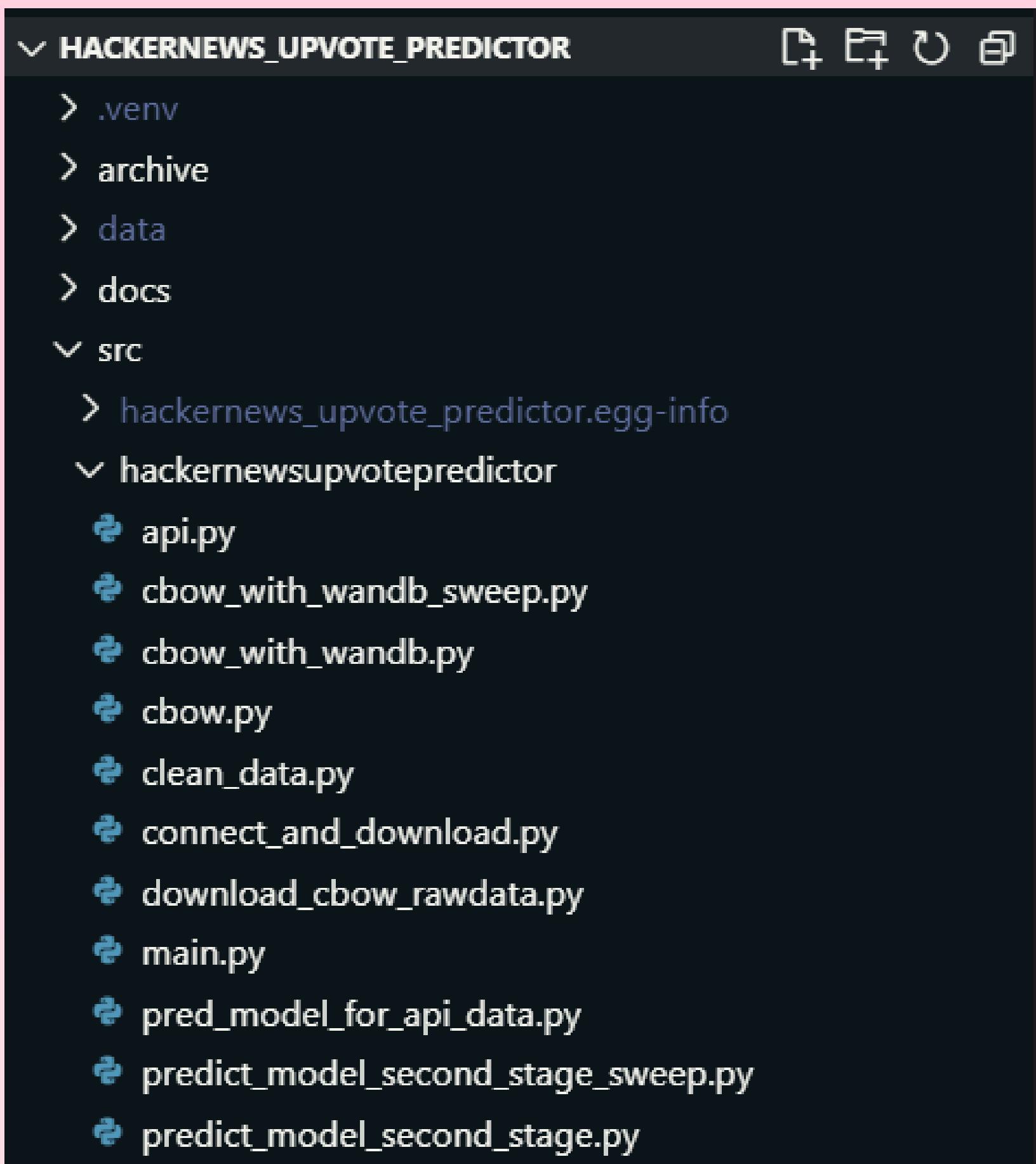


# ARCHITECTURE

## Overall project directions - Approach 3: Target single score



# PROJECT STRUCTURE



The project structure is as follows:

- > temp
- > toy\_models
- > wandb
- ⚙️ .env
- ❖ .gitignore
- ☰ .python-version
- 🔑 LICENSE.md
- ❖ main.py
- ⚙️ pyproject.toml
- ⓘ README.md
- \$ send\_post.sh
- ☰ uv.lock

### **### Sequence of running**

Ensure initial python set up has been done

0. Add a data and temp folder to root
1. Run `connect\_and\_download.py` file - this will connect to the database, and download the hackernews items (joined with user data) into a parquet file
2. Run `download\_cbow\_rawdata.py` that will download the wikipedia text data to data/text8
3. Run `cbow.py` to run the CBOW code
4. Run `clean\_data.py` which will read the above parquet file, and then extract the feature data of how many days the user has existed, title data, as well as the target data (the upvote score)
5. Run `predict\_model\_second\_stage.py` that will train another model on the feature of how many days since the user has been created + title to give final output (score prediction)
  - This will output the loss functions for each run

### **### How to run the api**

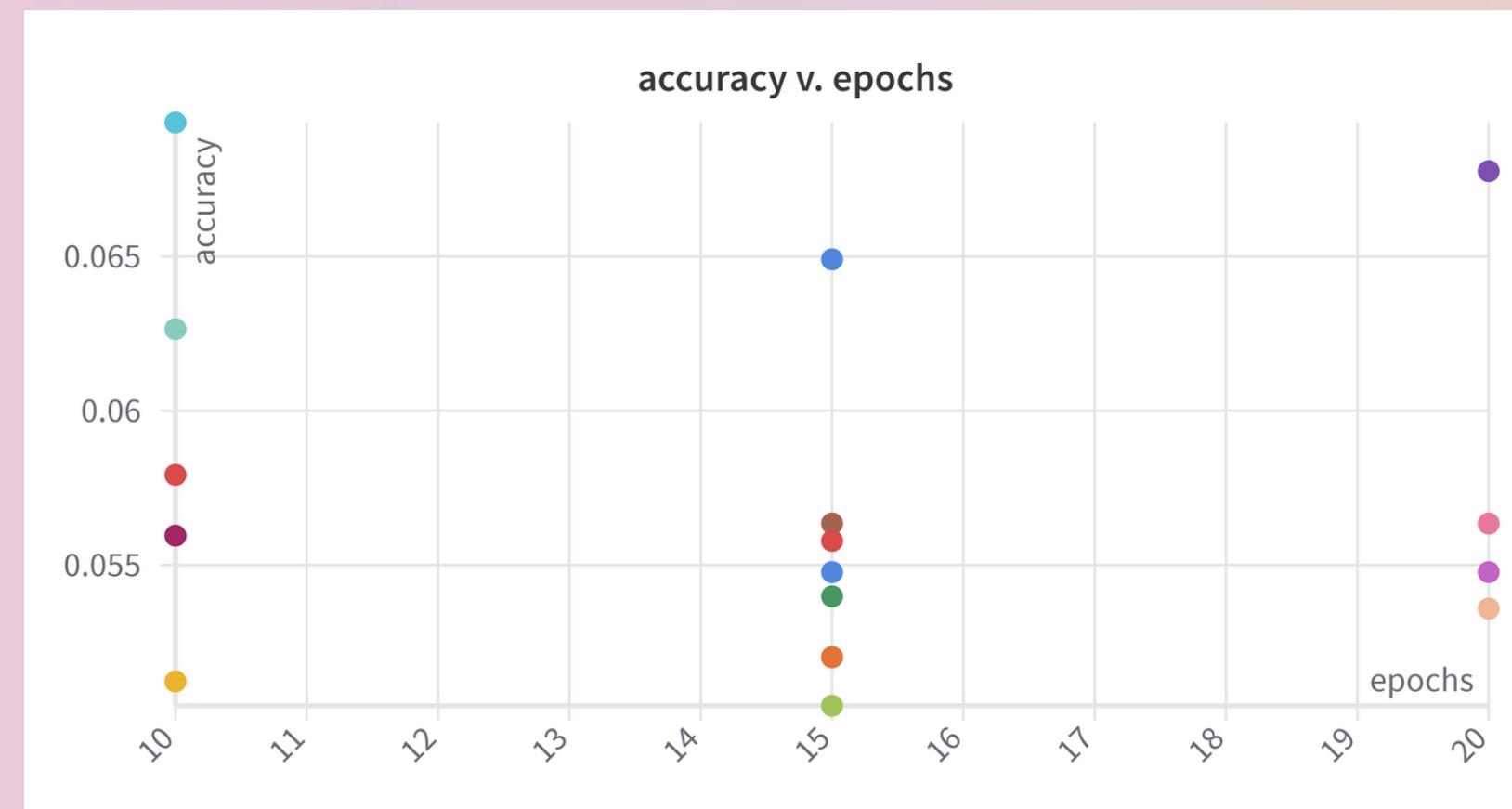
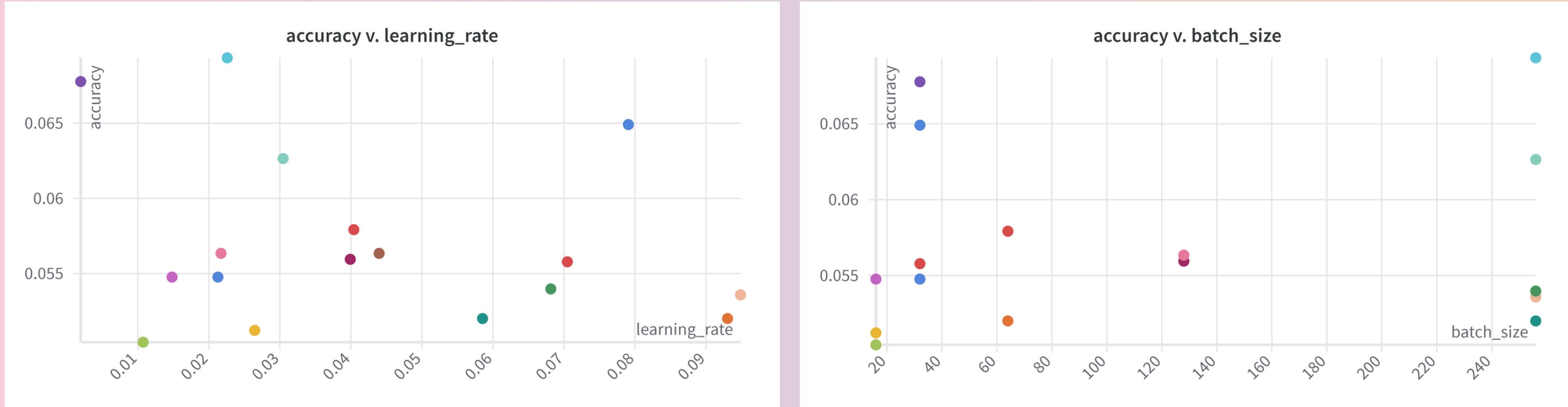
1. Run `uvicorn src.hackernewsupvotepredictor.api:app --reload` to run the API server
  - You should be able to see in <http://localhost:8000/healthcheck> to see it up and running
2. Then edit the send\_post.sh bash script with your custom values
3. Run send\_post.sh to see the predicted upvote from the model



**Time to take action**



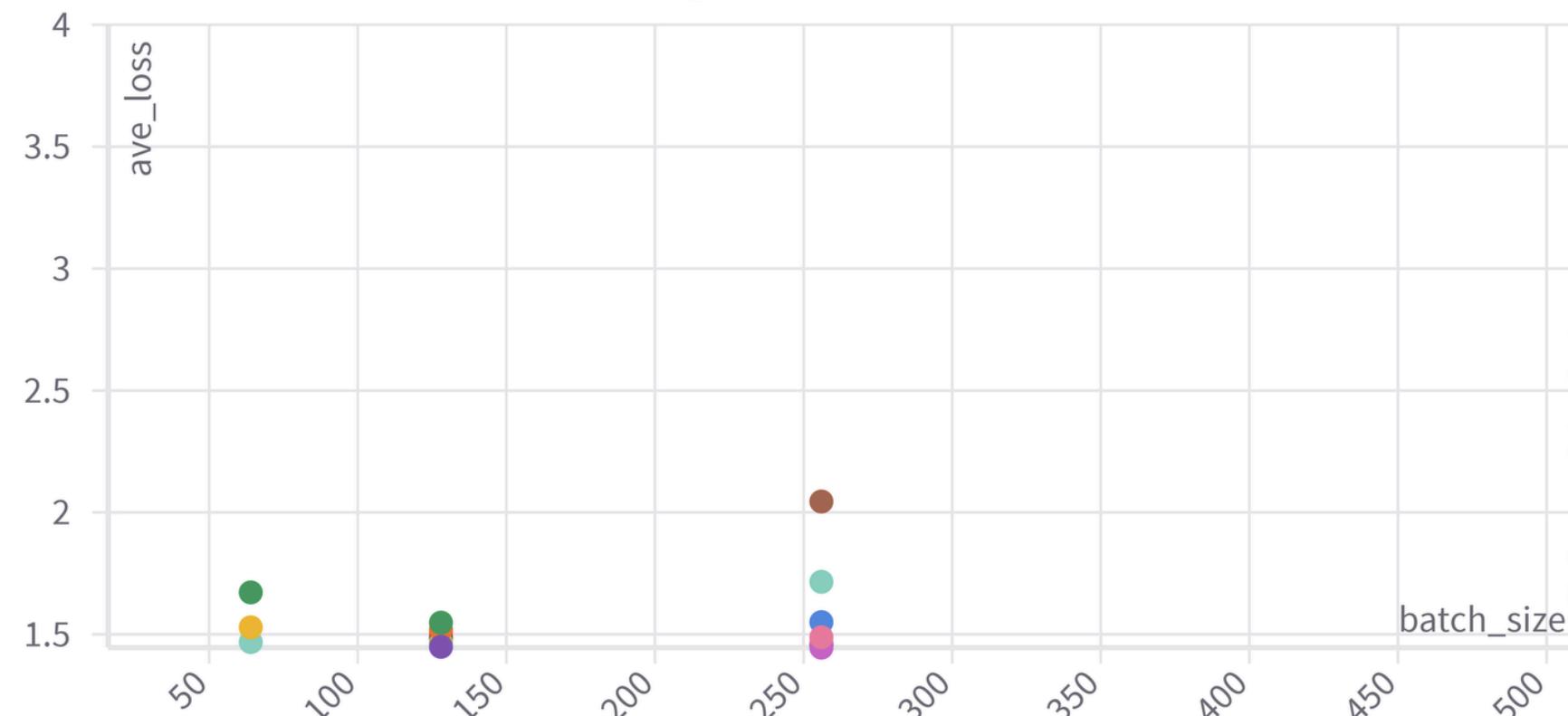
# RESULTS - CBOW



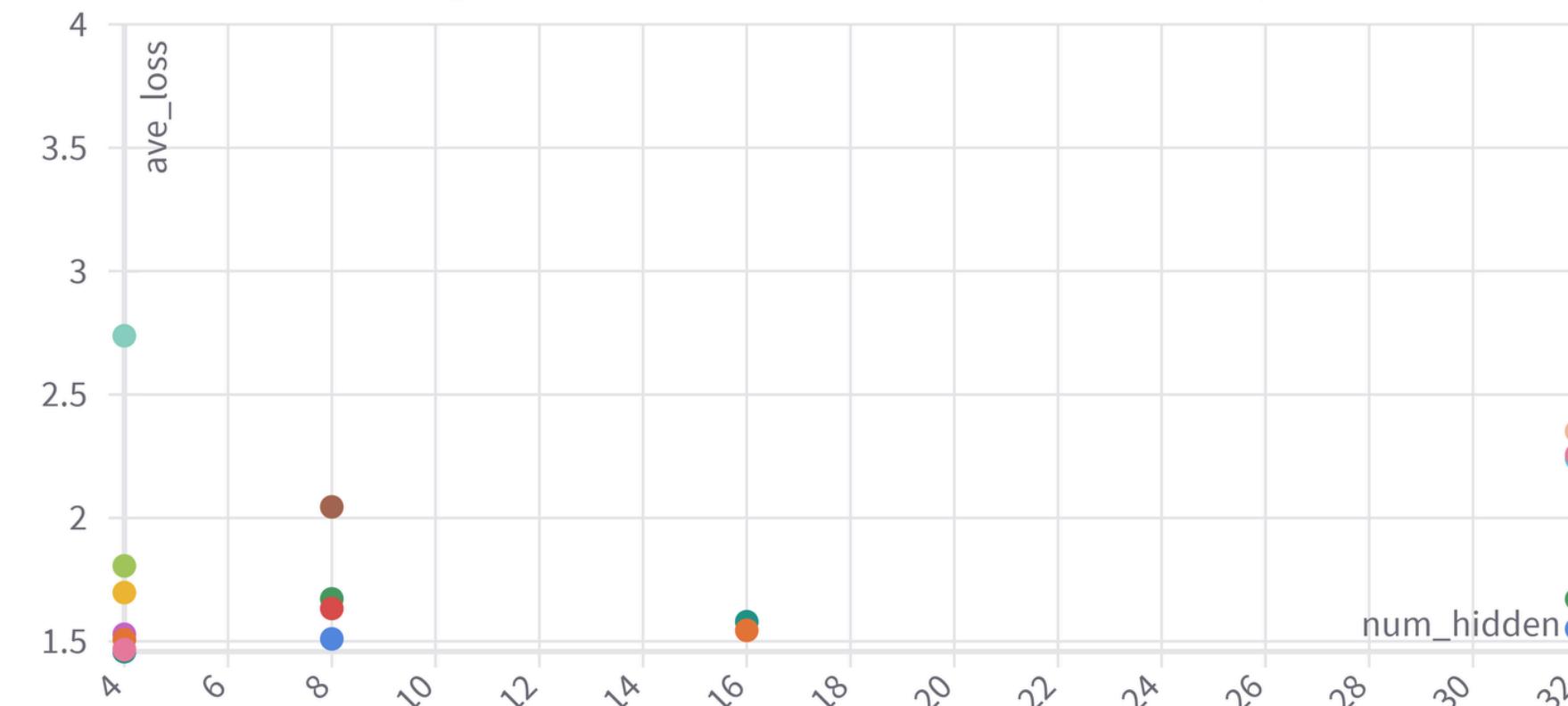
# RESULTS - PREDICTION MODEL

```
Epoch 1, Loss: 730.2390
Epoch 51, Loss: 1.5738
Epoch 101, Loss: 2.0703
Epoch 151, Loss: 1.8469
Epoch 201, Loss: 1.7715
```

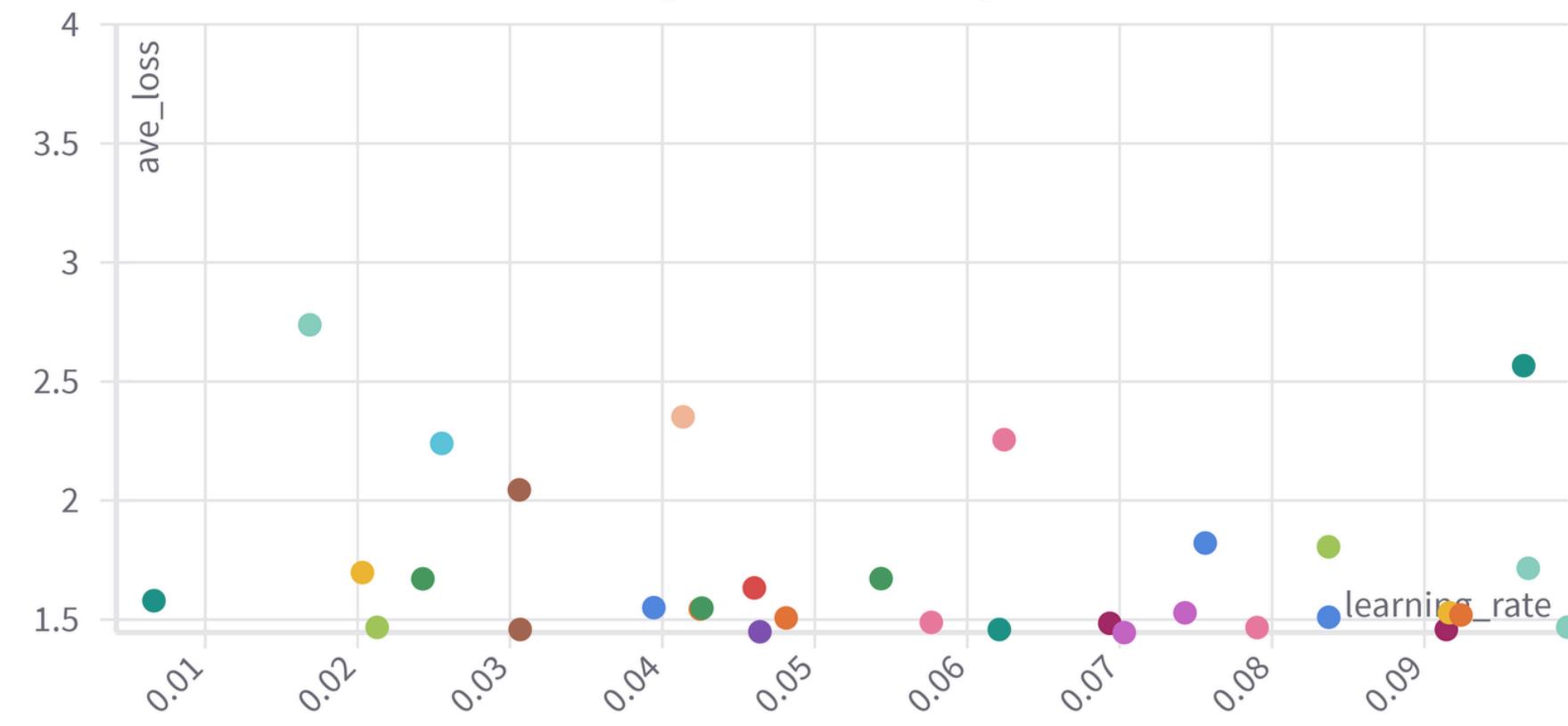
Average Loss vs Batch Size



Average Loss vs Number of Nodes in the Hidden Layer



Average loss vs Learning Rate



When you realize  
that we are closer  
to 5 PM, than we  
are to 10 AM

