

Reproducible Pipeline for Large-scale Data Analysis

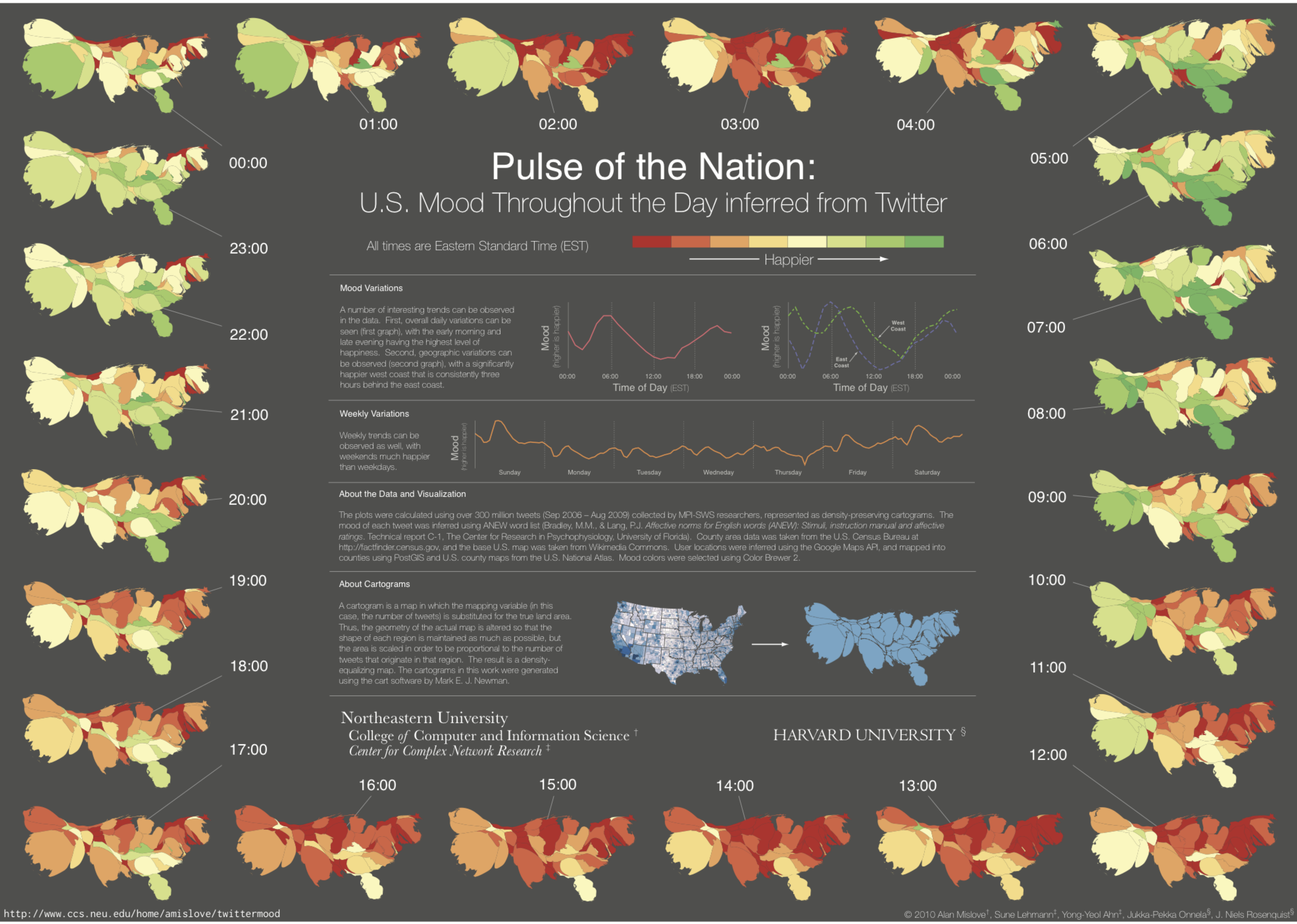
Lecture for ISEA Training Program

Wei Ai <aiwei@umd.edu>

University of Maryland

So, you did
some data
science...

And you have
a beautiful
trained model /
visualization /
dashboard to
show for it





But Behind the Theme?

Reproducibility and Reliability Matters?

- Other people will thank you.
 - Collaborate more easily with you on this analysis
 - Learn from your analysis about the process and the domain
 - Feel confident in the conclusions at which the analysis arrives

- You will thank you.
 - Which notebook should I run first?
 - Where to go if I need to add one more column?
 - What if the questions are asked 4 years later?

```
cohort0.ipynb  
cohort_new.ipynb  
data_preprocessing.ipynb  
dataset_preperation.ipynb  
...
```

-

Cookiecutter Data Science

- A open-source project structure developed by DrivenData.
 - A template directory structure
 - Some opinions of good data science practice
 - And some suggested workflows.

LICENSE	<-	Open-source license if one is chosen
Makefile	<-	Makefile with convenience commands like `make data` or `make train`
README.md	<-	The top-level README for developers using this project.
data		
external	<-	Data from third party sources.
interim	<-	Intermediate data that has been transformed.
processed	<-	The final, canonical data sets for modeling.
raw	<-	The original, immutable data dump.
docs	<-	A default mkdocs project; see www.mkdocs.org for details
models	<-	Trained and serialized models, model predictions, or model summaries
notebooks	<-	Jupyter notebooks. Naming convention is a number (for ordering), the creator's initials, and a short `-` delimited description, e.g. `1.0-jqp-initial-data-exploration`.
pyproject.toml	<-	Project configuration file with package metadata for <code>{{ cookiecutter.module_name }}</code> and configuration for tools like black
references	<-	Data dictionaries, manuals, and all other explanatory materials.
reports	<-	Generated analysis as HTML, PDF, LaTeX, etc.
figures	<-	Generated graphics and figures to be used in reporting
requirements.txt	<-	The requirements file for reproducing the analysis environment, e.g. generated with `pip freeze > requirements.txt`
setup.cfg	<-	Configuration file for flake8
{{ cookiecutter.module_name }}	<-	Source code for use in this project.
__init__.py	<-	Makes <code>{{ cookiecutter.module_name }}</code> a Python module
config.py	<-	Store useful variables and configuration
dataset.py	<-	Scripts to download or generate data
features.py	<-	Code to create features for modeling
modeling		
__init__.py		
predict.py	<-	Code to run model inference with trained models
train.py	<-	Code to train models
plots.py	5	<- Code to create visualizations

Managing a Data Science Project

- Manage **data**
 - Data Pipeline as a DAG
- Manage **code**
 - Notebook vs script
- Manage **project**
 - Version control
- Manage **collaboration** (internal and external)
 - Design doc
 - Issue & Kanban

Rethink about Plumbing — Data Pipeline should be a Directed Acyclic Graph (DAG)

- Each step of your analysis is a node in a *directed graph with no loops*.
 - You can run through the graph forwards to recreate any analysis output, or you can trace backwards from an output to examine the combination of code and data that created it.
- **Raw data must be treated as immutable.**
 - it's okay to read and copy raw data to manipulate it into new outputs, but never okay to change it in place.
- The goal: anyone should be able to re-run your analysis using only your code and raw data and produce the same final products.

The data directory structure

— LICENSE	<— Open-source license if one is chosen
— Makefile	<— Makefile with convenience commands like `make data` or `make train`
— README.md	<— The top-level README for developers using this project.
— data	
— external	<— Data from third party sources.
— interim	<— Intermediate data that has been transformed.
— processed	<— The final, canonical data sets for modeling.
— raw	<— The original, immutable data dump.
— docs	<— A default mkdocs project; see www.mkdocs.org for details
— models	<— Trained and serialized models, model predictions, or model summaries
— notebooks	<— Jupyter notebooks. Naming convention is a number (for ordering), the creator's initials, and a short `—` delimited description, e.g. `1.0-jqp-initial-data-exploration`.

- Most data should not be kept in version control
 - Use .gitignore to specify files not tracked by git.
 - Small and stable data files are fine (such as zipcode-city directory).

Discussion: Data Management

- Where do you store your data?
- How do you manage processed data or interim data?
- How do you share processed data with your teammates?

— LICENSE	<— Open-source license if one is chosen
— Makefile	<— Makefile with convenience commands like `make data` or `make train`
— README.md	<— The top-level README for developers using this project.
— data	
— external	<— Data from third party sources.
— interim	<— Intermediate data that has been transformed.
— processed	<— The final, canonical data sets for modeling.
— raw	<— The original, immutable data dump.
— docs	<— A default mkdocs project; see www.mkdocs.org for details
— models	<— Trained and serialized models, model predictions, or model summaries
— notebooks	<— Jupyter notebooks. Naming convention is a number (for ordering), the creator's initials, and a short `—` delimited description, e.g. `1.0-jqp-initial-data-exploration`.

Notebooks are Good for Demonstration and Experimentation

- Small-scale data wrangling (manipulating dataset into desired structure).
- Running sanity check (checking summary statistics),
- Developing models with unit tests.
- Visualizing results and generating reports.

But Notebooks are bad for repetition (reproduction).

- The flexibility of running cells in arbitrary order creates additional dependency to manage
- Avoid writing in your documentation:
 - “In order to reproduce this figure, run cell A first, skip cell B, then run cell C, and then rerun cell A ...”
- Notebooks are hard for version control
- One more reason to hate them: They don't work well with AI coding assistants.

If You Open a Jupyter Notebook with a Text Editor

```
opensmile.ipynb > ...
Open in Notebook Editor
1 {
2 "cells": [
3 {
4 "cell_type": "code",
5 "execution_count": 11,
6 "metadata": {},
7 "outputs": [],
8 "source": [
9 "# To run this notebook, please install librosa, opensmile, in addition to standard packages such
10 "\n",
11 "import librosa\n",
12 "import opensmile\n",
13 "from IPython.display import Audio, display\n",
14 "\n",
15 "import matplotlib.pyplot as plt\n",
16 "import librosa.display"
17 ]
18 },
19 {
20 "cell_type": "code",
21 "execution_count": 12,
22 "metadata": {},
23 "outputs": [],
24 "source": [
25 "# Filename identifiers of the RAVDESS dataset\n",
26 "# See: https://zenodo.org/records/1188976\n",
27 "\n",
28 "# Modality (01 = full-AV, 02 = video-only, 03 = audio-only).\n",
29 "# Vocal channel (01 = speech, 02 = song).\n",
30 "# Emotion (01 = neutral, 02 = calm, 03 = happy, 04 = sad, 05 = angry, 06 = fearful, 07 = disgust,\n31 "# Emotional intensity (01 = normal, 02 = strong). NOTE: There is no strong intensity for the 'nei\n32 "# Statement (01 = \"Kids are talking by the door\", 02 = \"Dogs are sitting by the door\").\n",
33 "# Repetition (01 = 1st repetition, 02 = 2nd repetition).\n",
34 "# Actor (01 to 24. Odd numbered actors are male, even numbered actors are female).\n",
35 "\n",
36 "male_calm_normal_kids1 = \"Audio_Speech_Actors_01-24/Actor_01/03-01-02-01-01-01-01.wav\"\n",
37 "male_happy_strong_kids1 = \"Audio_Speech_Actors_01-24/Actor_01/03-01-03-02-01-01-01.wav\"\n",
38 "female_calm_normal_kids1 = \"Audio_Speech_Actors_01-24/Actor_02/03-01-02-01-01-01-02.wav\"\n",
39 "female_happy_strong_kids1 = \"Audio_Speech_Actors_01-24/Actor_02/03-01-03-02-01-01-02.wav\"
40 ]
```

```
6d27t9XNA0CD2Hs0gNcICAg09PaRI0†qHdjZ6ACAGQhNALzGkSNHHN/PnD LTRowYITt27HAcK1eunL4AgBkYngPgNa
+f0dXmvr1q16h3f1n0oxDzzwgJw4ccKCswbgKQhNAHzeRx99JJUrV5a1a9fqAPXY4/JnXfeKW3btpXExETp2rWrDk
Xr15eEhAQJCwvTIWrgwIH6mBrm03nypPz666/6/hMmTNCB6ZVXXpG4uDj9/bRp02Tp0qWyc+d0q08HgEwoaQLg85o2
+yz8nJycvTXs2fPSq9evWTs2LF5nqtatWqmtxeAZyI0AcBlrrvu0pk1a5bExsZKcDAfkwAuoqYJAC4zaNAg0XXqLnx
FWAjaVzAQAArog/mQAAAAwgNAEABhAaAIAADCA0AQAGAAoQKAAMAAQhMAAIABhCYAAAADCE0AAAAGEJoAAAAMIDQ
"text/plain": [
"<Figure size 600x300 with 1 Axes>"
]
},
"metadata": {},
"output_type": "display_data"
},
{
"name": "stdout",
"output_type": "stream",
"text": [
"Extracting LLD features using OpenSMILE feature set: eGeMAPSv02\n",
"Extracted 25 LLD features with 366 frames\n"
]
},
{
"data": {
"image/png":
"iVBORw0KGgoAAAANSUHEUgAAAK0AAAGGCAYAAABmPbWYAAAAOXRFWHRTb2Z0d2FyZQBNYXRwbG90bGliIHZlcnNpb
R4n02dB3hT5dvGH7pp6QAKHXS957CyKKDAVBGYoLcaCf/sWJiAo0xAGKooADRBWQv2XuVXuoptKWM7j3yXc+b
67rNMnJyTnvezL03WdWUSgUCgIAAAAAA0XiVv7TAAAAAACAgWgCAAAAADACiYAAAAAAC0AaAIAAAAAAMAKIJgAAAA
AHAQIJJoAsCNLly4VF2V9y/Tp0+09PFmTmJgoRM+JEyesfiwWmobex40bN1r9+AAA2+Bho+MAAMrh/
ffffp3r16mmta9mypo3G4yyiaebMmcJ61rZtW6sfz9vbm3744Ycy69u0aWP1YwMabANEWA0wJAhQ6hjx472HgaoBB4
+u1CoWCPvjgA4qIiCBfX1/q168fxcTElDnG3bt36dVXX6VWrVpRtWrVKCAgQAi8kydPam23c+d0cew//xTWHfq1Kl
kodOnSgqlWrUo0aNeiRRx6hhIQErW369u0rrHBnz54V4+bx8/E++eQTrf06tRJ30djSa4yPkfMf//9Rw8//DBFRUU
+FaTgwcP0r333kuBgYHIXPbp00d8hvSN//Lly+IzFxQUJLbnc5uTk1MnDwfHh/vj+fau3dv2rx5s3hu4sSJFBwcTI
+++K0TTfffdJ5Zjx46JC1ZBQYHWvq5cuUJr1qwRAoPdhdv3qTvvvt0XFhZAISHh2ttP3v2bHFB53gsvsB+/fXX50n
+5cqV9Prrrwtxx6KuWbNmWuXJ+3/66aepV69e4nXdu3cXtytWrBAX+yLTPldNmjXp0KFD4ljXr18Xz5mL7nvIc2dxI
cXgoyFjyY8Rn7f+H3mzwy7I1kIf/zxx+ptWDTz+8zj5PfAy8tLCDM+Fn/Gxo8fTz//DNt2rSJ7r//fvXrkp0TxTY8
55ZcV7u7uirS0NPE4JSVF4eXlPrg6dKjWdm+++aZ4veY+8/LyFMFXvRHzYuLU3h7eyvef/999bod03aI17Zs2VJRU
+fn5itDQUMWoUaPU6w4fPiy24/OiS050TP1s2fPFmONj49Xr+PzbszPJJ8zfe8hj1VzP3w+NDF23nx0a9eurWjbtq
```


Refactor codes from .ipynb to .py

- Consciously and continuously extracting **common building blocks**
 - From notebooks (.ipynb) to source files (.py)
- Prevents duplication of code across multiple notebooks
- Separates the multiple concerns into logical units.
 - data layer (ETL code) -> modeling and experimentation -> output layer
- Make sure each notebook is always correct if “run all”
- Use naming conventions to indicate “owner” and give a sense of the order the analysis
 - e.g. <step>-<initials>-<description>.ipynb
 - 3.1-aiwei-extract-audio-features.ipynb

Discussion: Do You Love/Hate Jupyter Notebooks?

- How long do you keep your notebook running?
- Have you tried Google CoLab? AWS SageMaker Studio? JupyterLab?
- Have you tried running Jupyter Notebook on a remote server?

Always use Version Control

- The most widely-acknowledged norm of software engineering.
- It is a good idea to initiate a project folder with a template
 - e.g. the Cookiecutter Data Science
- It makes it natural for you to adopt other good habits:
 - Setting up an environment (reduce the dependency error, increase consistency)
 - Enabling code reviews.
 - Making it easy to document decisions and design choices.
- It is also how AI coding assistant works (by creating a “diff”)

Keep the Following Away from Version Control

- Large data file,
- Secrets
 - OpenAI API keys, AWS Credentials, SSH keys, Database Passwords, ...
 - Store them in a configuration file. (.env, .config, etc.)
- Other configuration file
 - Personal IDE setting (keyboard mapping)
 - Machine-dependent variables
 - Username, path-to-project-folder, etc.
- Use a package to load them automatically.

```
import os
from dotenv import load_dotenv, find_dotenv

# find .env automatically by walking up directories until it's found
dotenv_path = find_dotenv()

# load up the entries as environment variables
load_dotenv(dotenv_path)

database_url = os.environ.get("DATABASE_URL")
other_variable = os.environ.get("OTHER_VARIABLE")
```


More on Jupyter Notebook

- Some IDE (e.g. VS Code, Cursor, PyCharm), with proper plugin, supports pseudo-notebooks.
- Use **#%%** to indicate new cells.

```
tooluse.py > calculator
Run Cell | Run Below | Debug Cell
1  #%%
2  # Block 1: Import necessary libraries
3  import json
4  from typing import Dict, List, Optional, Union
5  import openai
6  from dotenv import load_dotenv
7  import os
8
Run Cell | Run Above | Debug Cell
9  #%%
10 # Load environment variables
11 load_dotenv()
12
Run Cell | Run Above | Debug Cell
13 #%%
14 # Block 2: Define the calculator tool
15 def calculator(operation: str, numbers: List[float]) -> float:
16     """
17     A simple calculator that performs basic arithmetic operator
18     """
19     Args:
20     operation: One of '+', '-', '*', '/'
21     numbers: List of numbers to operate on
22
23     Returns:
```

```
opensmile.ipynb > ...
Open in Notebook Editor
1  {
2  "cells": [
3  {
4  "cell_type": "code",
5  "execution_count": 11,
6  "metadata": {},
7  "outputs": [],
8  "source": [
9  "# To run this notebook, please install librosa, opensmile, in addition to standard packages such
10 "\n",
11 "import librosa\n",
12 "import opensmile\n",
13 "from IPython.display import Audio, display\n",
14 "\n",
15 "import matplotlib.pyplot as plt\n",
16 "import librosa.display"
17 ]
18 },
19 {
20 "cell_type": "code",
21 "execution_count": 12,
22 "metadata": {},
23 "outputs": [],
24 "source": [
25 "# Filename identifiers of the RAVDESS dataset\n",
26 "# See: https://zenodo.org/records/1188976\n",
27 "\n",
28 "# Modality (01 = full-AV, 02 = video-only, 03 = audio-only).\n",
29 "# Vocal channel (01 = speech, 02 = song).\n",
30 "# Emotion (01 = neutral, 02 = calm, 03 = happy, 04 = sad, 05 = angry, 06 = fearful, 07 = disgust,\n31 "# Emotional intensity (01 = normal, 02 = strong). NOTE: There is no strong intensity for the 'nei\n32 "# Statement (01 = \"Kids are talking by the door\", 02 = \"Dogs are sitting by the door\").\n",
33 "# Repetition (01 = 1st repetition, 02 = 2nd repetition).\n",
34 "# Actor (01 to 24. Odd numbered actors are male, even numbered actors are female).\n",
35 "\n",
36 "male_calm_normal_kids1 = \"Audio_Speech_Actors_01-24/Actor_01/03-01-02-01-01-01.wav\"\n",
37 "male_happy_strong_kids1 = \"Audio_Speech_Actors_01-24/Actor_01/03-01-03-02-01-01-01.wav\"\n",
38 "female_calm_normal_kids1 = \"Audio_Speech_Actors_01-24/Actor_02/03-01-02-01-01-01-02.wav\"\n",
39 "female_happy_strong_kids1 = \"Audio_Speech_Actors_01-24/Actor_02/03-01-03-02-01-01-02.wav\"\n",
40 ]
```

Discussion: Version Control

- Have you heard about git?
- Have you used git?
- What do you think is the most effective way to learn git?
-

Project Management — More Than Version Control

- Issues
- Project
- I personally keep two documents:
 - A “Design Doc”
 - A “Lab Log”

Start with a simple pipeline

- Data ETL,
 - Feature extraction,
 - Preliminaries - summary statistics, visualizing distribution, etc.
 - Model fitting (a linear model is OK!)
 - Model evaluation
-
- Allow you to think about problem formulation and evaluation criteria early on.

Other notes

- Test
 - Unit Test, Integration Test
- Log
- Sanity check
 - Rule-of-thumb for validating findings from exploratory data analysis:
 - Findings that match intuition
 - Findings that are serendipitous but can be explained
 - Findings that warrant further investigation or intervention

```
2022-11-01 14:10:02 INFO starting run on branch master (HEAD @ 37c02ab)
2022-11-01 14:10:02 DEBUG set random seed of 42
2022-11-01 14:10:02 DEBUG reading config file
2022-11-01 14:10:02 DEBUG ... settings: {"n_threads": 4, <...snip...>}
2022-11-01 14:10:03 INFO reading in and merging data files
2022-11-01 14:10:39 INFO finished loading data: 2,835,824 rows
2022-11-01 14:10:40 WARNING ... dropped 126,664 rows where ID was duplicated (4.47%)
2022-11-01 14:10:41 WARNING ... dropped 2,706 rows where column `total` was null (0.09%)
2022-11-01 14:10:41 INFO creating train/test split
2022-11-01 14:10:41 INFO ... train: 0.5 [1,353,227 rows]
2022-11-01 14:10:41 INFO ... test: 0.5 [1,353,227 rows]
2022-11-01 14:10:42 INFO starting grid search cross validation ...
2022-11-01 14:21:01 DEBUG ... 30/120
2022-11-01 14:31:46 DEBUG ... 60/120
2022-11-01 14:42:31 DEBUG ... 90/120
2022-11-01 14:53:03 DEBUG ... 120/120
2022-11-01 14:53:03 INFO finished cross validation, writing best parameters
2022-11-01 14:53:03 DEBUG ... runs/2022-11-01_14-10-02/parameters/xgboost.yml
2022-11-01 14:53:03 DEBUG ... runs/2022-11-01_14-10-02/parameters/random_forest.yml
2022-11-01 14:53:03 DEBUG ... runs/2022-11-01_14-10-02/parameters/adaboost.yml
2022-11-01 14:53:03 INFO training voting classifier on ensemble of 3 best models...
2022-11-01 14:53:19 INFO making predictions
2022-11-01 14:53:22 INFO writing results to runs/2022-11-01_14-10-02/results.yml
2022-11-01 14:53:22 DEBUG ... results: precision=0.9624 recall=0.9388 f1=0.9514
2022-11-01 14:53:23 INFO writing predictions to runs/2022-11-01_14-10-02/predictions.csv
```

Missing semester: did you ever

- Ssh to a machine, but figuring out that you need to type in password every time?
- Run a job on a remote machine, but afraid that the job will get lost if you disconnect
- If so, check out **The Missing Semester of Your CS Education**
 - <https://missing.csail.mit.edu/>

Open Discussion

- How comfortable do you find using AI assistant in coding?
- Are there “best” practices for using AI assistants?