

Structuring Machine Learning Projects

Sayak Paul, Deep Learning Associate at [PyImageSearch](#)
DevFest Izmir



November 23, 2019



2019



Agenda

- Machine learning: The life-cycle
- Why care about the structure?
- A **structured** approach to **structuring**
 - A polished directory structure
 - Workspace setup
 - Building mental models of the project's flow
 - Experimentation
 -
- Taking the next steps
- Guiding lights

Life-cycle of a machine learning project

- Problem understanding
- Data collection, wrangling and so on
- Understanding of the data
- Beginning the modeling process
- Evaluate, tune, repeat
- Model deployment, monitoring and so on...

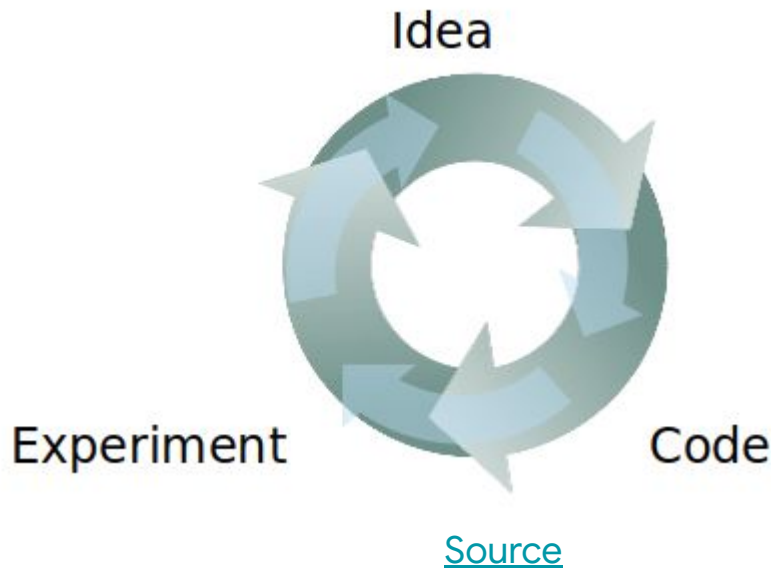
Note that the above steps are very summarized



Structuring machine learning projects: The **need**

Applied machine (and deep) learning is iterative

- Lots of experimentations make it hard to keep track
- Hyperparameter tuning
- Dataset reconstructions
- KPI: Time-Cost tradeoff



Reproducibility crisis



- **Stochasticity** in machine learning models (**neural nets** specifically)
- Code breakdown due to **dependency** mismatches
- Difference in **machine configurations**

Versioning data and codebase

- Training data might change over time (less frequent)

Versioning data and codebase

- Training data might change over time (less frequent)

Pain point!

Versioning data and codebase

- Training data might change over time (less frequent)
- Changes in project's codebase are far more frequent

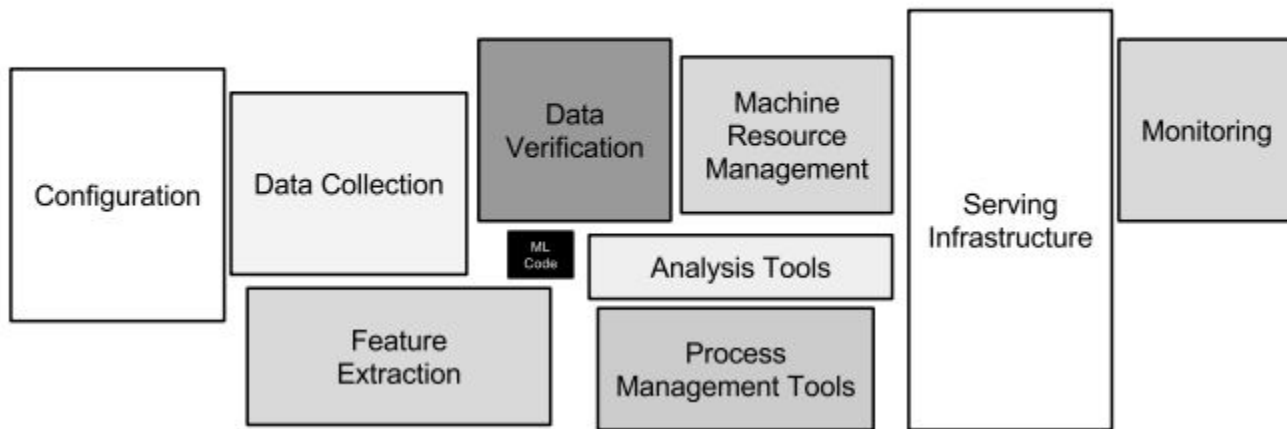
Regularity in checkpointing



Decrease in **validation loss** was steady during the training, but **what if that model is lost?**

A directory structure

- A full-fledged ML project is **not just about models**



[Source](#)

A **high-interest** technical debt

- Lots of ideas, lots of experiments - gives birth to **technical debt**
- Not everyone in the team understands everything
- Lack of *documentation* as it is considered *not-so-cool*

A **high-interest** technical debt

- Lots of ideas, lots of experiments - gives birth to **technical debt**
- Not everyone in the team understands everything
- Lack of *documentation* as it is considered *not-so-cool*

This hurts the overall development of the project!

A **high-interest** technical debt

- Lots of ideas, lots of experiments - gives birth to **technical debt**
- Not everyone in the team understands everything
- Lack of *documentation* as it is considered *not-so-cool*

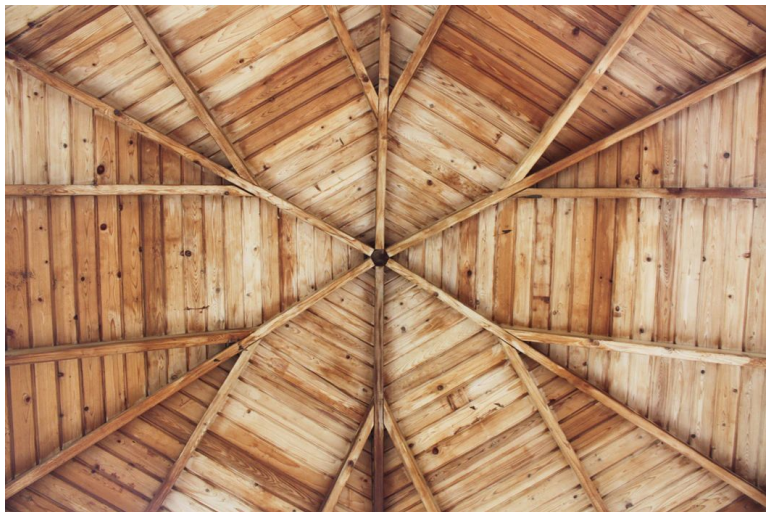
This hurts the overall development of the project!

Machine Learning: The High-Interest Credit Card of Technical Debt

D. Sculley, Gary Holt, Daniel Golovin, Eugene Davydov,
Todd Phillips, Dietmar Ebner, Vinay Chaudhary, Michael Young
{dsculley, gholt, dgg, edavydov}@google.com
{toddphillips, ebner, vchaudhary, mwyong}@google.com
Google, Inc

Quick summary

- Applied machine learning is an iterative process
- Reproducibility crisis
- Versioning data and codebase
- Regularity in checkpointing
- Directory structures
- Technical debt



Structuring Machine Learning projects: A definitive approach

Starting with a polished directory structure

- Data (along with scripts to *download it* and *preprocess it*)
- Experimentations
- Web backend
- Utility scripts
- Model building and model training scripts

A reference directory structure

```

apparel_classifier/           # Package that can be deployed as a self-contained prediction system
  __init__.py

apparel_predictor.py         # Takes a raw image and obtains a prediction

datasets/                   # Code for loading datasets
  __init__.py
  dataset.py                 # Base class for datasets - logic for downloading data (if not available)
  fmnist_dataset.py
  fmnist_essentials.json
  dataset_sequence.py

models/                     # Code for instantiating models, including data preprocessing and loss functions
  __init__.py
  base.py                   # Base class for models
  image_model.py

networks/                   # Code for building neural networks (i.e., 'dumb' input->output mappings) used by models
  __init__.py
  mlp.py

tests/
  support/                  # Raw data used by tests
  test_apparel_predictor.py # Test model on a few key examples

weights/                    # Weights for production model
  Image_Model_FMNIST_Dataset_weights.h5

util.py

training/                   # Code for running training experiments and selecting the best model.
  run_experiment.py         # Parse experiment config and launch training.
  util.py                   # Logic for training a model with a given config
```

Workspace setup

- Use of environment management tools - **pipenv**, **virtualenv** etc
- Incorporating containers - **Docker**, **Kubernetes** etc
- Use of unified machine learning platforms like [FloydHub](#)

Keeping track of the experiments

Machine learning experiments have a lot of components:

- Loss/Accuracy metrics
- Model parameters
- CPU, GPU and disk usage

And so on ...

Keeping track of the experiments

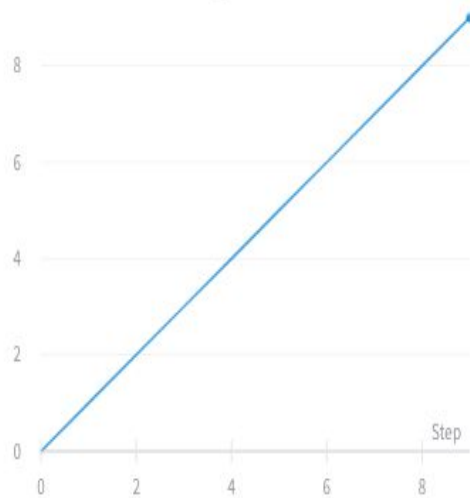
Machine learning experiments have a lot of components:

- Loss/Accuracy metrics
- Model parameters
- CPU, GPU and disk usage

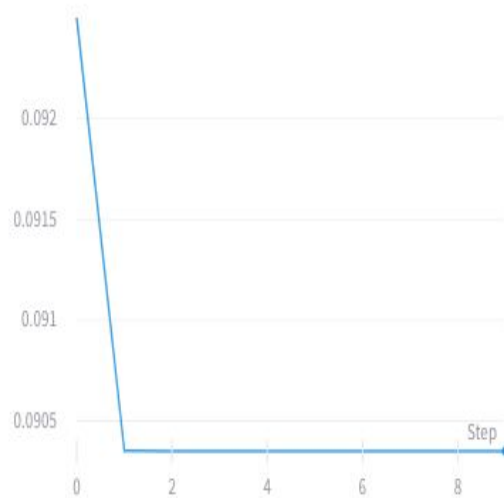
And so on ...

Use [Weights and Biases](#) to automatically keep track of these for you on the **cloud**!

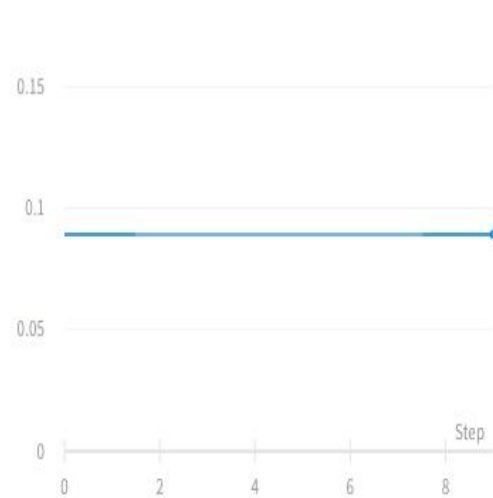
epoch



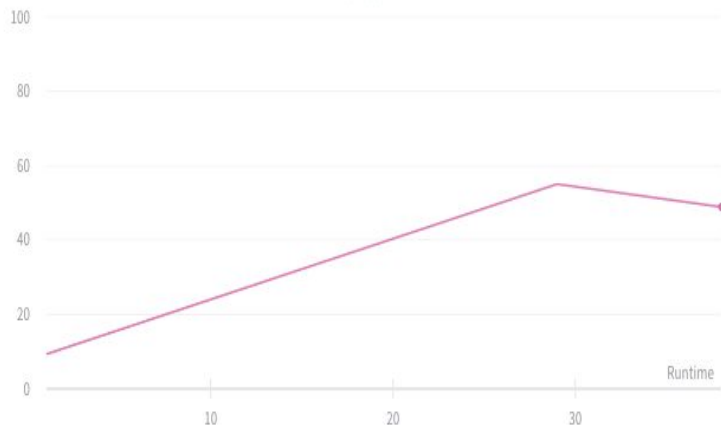
loss



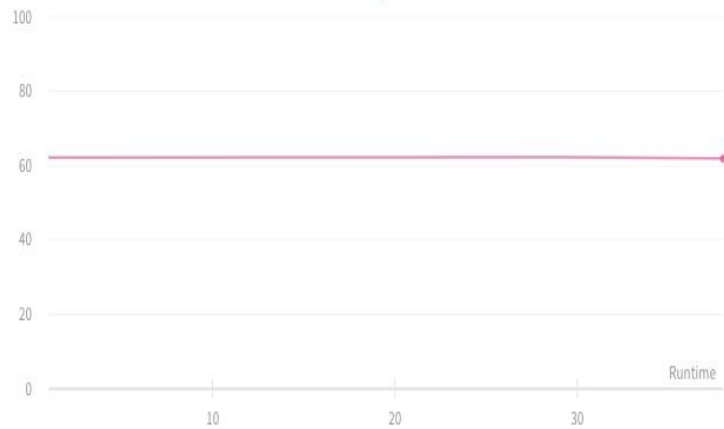
val_loss





























CPU %



Memory %



 Name (28 ...	State	Notes	User	Tags	Created	Runtime	Sweep	epochs	layers	accuracy	epoch	loss	val_accu...	val_loss ▼
  a...	finished	Add notes	saya...		1mo ago	35s	6v8m5w0s	5	32	0.8795	4	0.3328	0.869	0.3756
  d...	finished	Add notes	saya...		1mo ago	9m 36s	-	5	128	0.8918	4	0.2932	0.8693	0.3692
  1...	finished	Add notes	saya...		1mo ago	42s	6v8m5w0s	5	96	0.8895	4	0.3023	0.8691	0.3664
  z...	finished	Add notes	saya...		1mo ago	34s	6v8m5w0s	5	128	0.8906	4	0.2958	0.8729	0.3563
  1...	finished	Add notes	saya...		1mo ago	39s	6v8m5w0s	5	64	0.8869	4	0.3093	0.8747	0.35
  p...	finished	Add notes	saya...		1mo ago	42s	6v8m5w0s	5	256	0.895	4	0.2834	0.8696	0.3488

 Name (28 ...	State	Notes	User	Tags	Created	Runtime	Sweep	epochs	layers	accuracy	epoch	loss	val_accu...	val_loss ▼
  a...	finished	Add notes	saya...		1mo ago	35s	6v8m5w0s	5	32	0.8795	4	0.3328	0.869	0.3756
  d...	finished	Add notes	saya...		1mo ago	9m 36s	-	5	128	0.8918	4	0.2932	0.8693	0.3692
  1...	finished	Add notes	saya...		1mo ago	42s	6v8m5w0s	5	96	0.8895	4	0.3023	0.8691	0.3664
  z...	finished	Add notes	saya...		1mo ago	34s	6v8m5w0s	5	128	0.8906	4	0.2958	0.8729	0.3563
  1...	finished	Add notes	saya...		1mo ago	39s	6v8m5w0s	5	64	0.8869	4	0.3093	0.8747	0.35
  p...	finished	Add notes	saya...		1mo ago	42s	6v8m5w0s	5	256	0.895	4	0.2834	0.8696	0.3488

Try it out here:

http://bit.ly/w_and_b

Building a mental image of the execution flow

- There will be a lot of [inter-connected blocks](#) in a project

Building a mental image of the execution flow

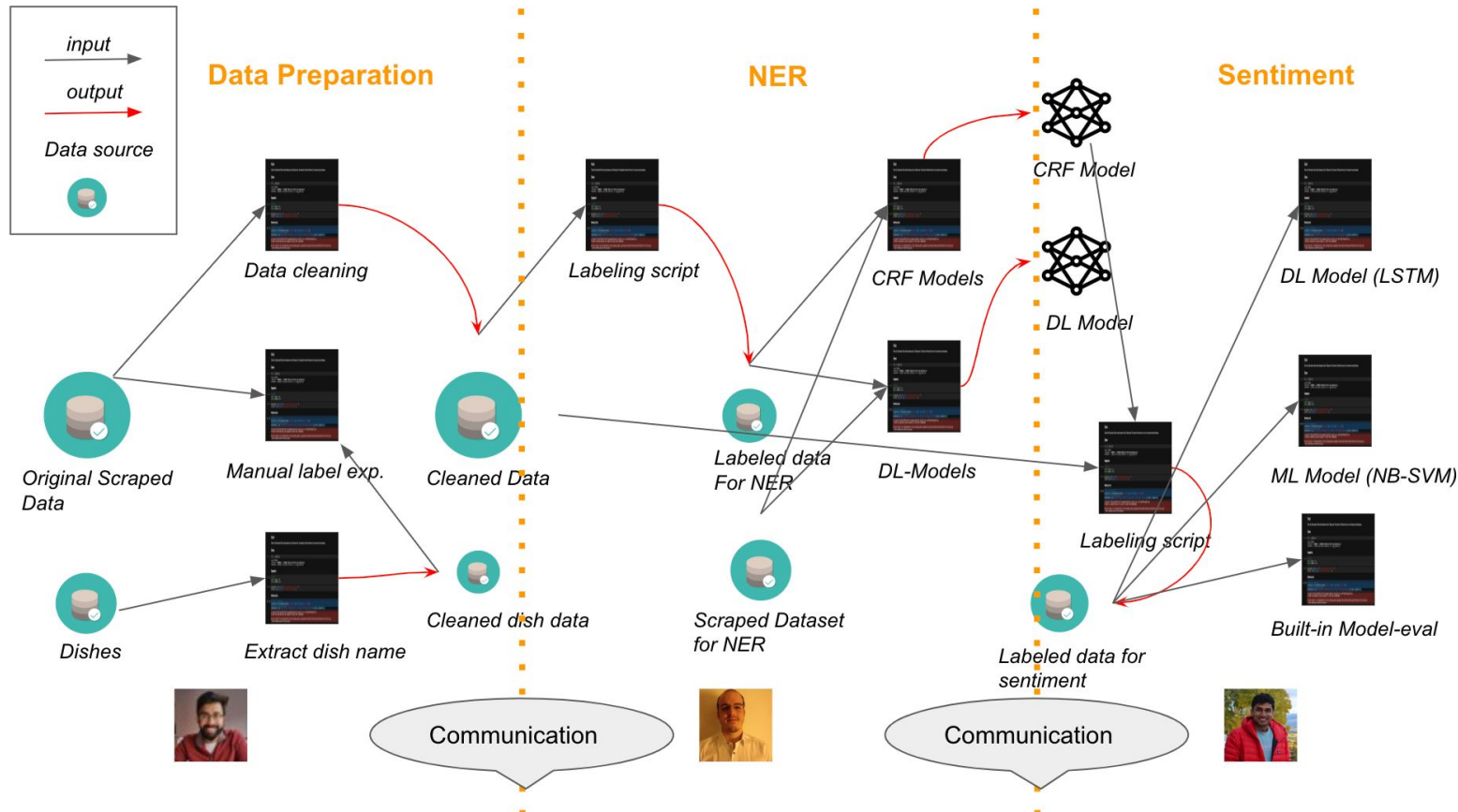
- There will be a lot of **inter-connected blocks** in a project
 - A block can have several sub-blocks too

Building a mental image of the execution flow

- There will be a lot of **inter-connected blocks** in a project
 - A block can have several sub-blocks too

Hence a mental model to keep track of these
lego blocks always helps!

A specimen of a mental model



Version controlling **data** and **codebase** *separately*

Here are some situations:

- Need to replace the **older** images with **newer** ones
- Need to add new images to an already existing training set
- Decision to incorporate **active learning** to select interesting test data points to manually label them and add to the existing training set

Version controlling **data** and **codebase** *separately*

Here are some situations:

- Need to replace the older images with newer ones
- Need to add new images to an already existing training set
- Decision to incorporate **active learning** to select interesting test data points to manually label them and add to the existing training set

Version control of codebase remains traditional!

FloydHub datasets

Create a new dataset

A dataset is a reusable collection of files that can be mounted into your jobs.

Owner

alice

Dataset name

my-cool-dataset

Description (optional)

Visibility

☒  Public

☐  Private

Create dataset

FloydHub datasets are a
good way to version
control your data!

Quick summary

- Maintaining a healthy directory structure
- Setting up the workspace
- Keeping track of your experiments
- Building a mental model of the project flow
- Version control of data and codebase

What did we *not* cover?

- Test cases for machine learning projects
- Model deployment

What did we *not* cover?

- Test cases for machine learning projects
- Model deployment

Maybe next time :)

Slides available here:
<http://bit.ly/izmir-sayak>

References

- [How to plan and execute your ML and DL projects](#)
- [Troubleshooting Deep Neural Networks](#)
- [Production Machine Learning Systems](#)

See you next time



Find me here:

sayak.dev

Thank you very much :)

