# It's The Data, Stupid!

March 24, 2022
Peter Gao

# Peter Gao

Cofounder and CEO at *Aquarium*

Early employee at Cruise, led CV team

Deep Learning research at Berkeley

**Alex Gude**
@alex_gude

Here is a real use case from work for model improvement and the steps taken to get there:

- Baseline: 53%
- Logistic: 58%
- Deep learning: 61%
- **Fixing your data: 77%**

Some good ol' fashion "understanding your data" is worth it's weight in hyperparameter tuning!

12:48 PM · Apr 24, 2019

♡ 1.3K    💬 377 people are Tweeting about this

# How To Improve Your ML System

- Improve your model code

- Improve your training dataset

- <u>Do these faster and more frequently</u>
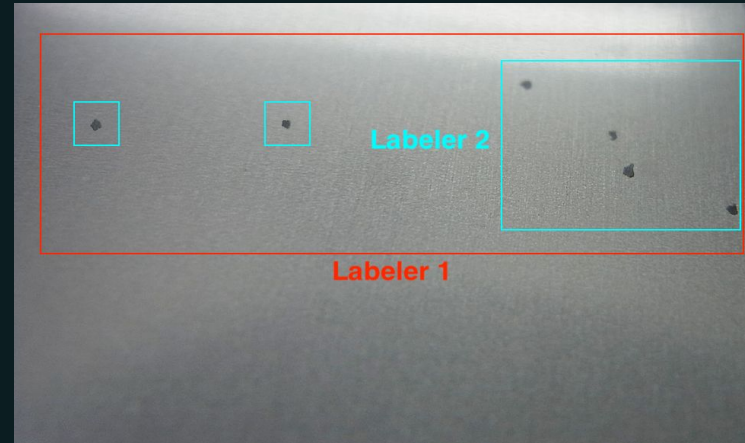
# "Old School" ML Vs Deep Learning

|  | Old School ML | Deep Learning |
|---|---|---|
| Example Tasks | Forecasting, recommendations | Object detection, seq2seq |
| Data types | Structured / tabular data | Unstructured data (imagery, audio, etc.) |
| Labeling | Labels come "for free" | Pay people to label data |
| Algorithms | Logistic regression, SVMs, random forests | Neural networks |
| Development Emphasis | - Data pipelines + infrastructure<br>- Feature engineering<br>- Model experimentation (ex: sparsity) | - Data pipelines + infrastructure<br>- Fine tuning pretrained models<br>- Improving quality + variety of datasets |

# How Do You Improve Your Data?

- Find problems in the data / model performance

- Figure out why the problems are happening

- Modify your dataset to fix the problems

- Make sure the problems are fixed as you retrain your model on the new dataset

- Deploy new model + repeat

# Types Of Data Problems

- Invalid data

- Labeling errors, ambiguities

- Difficult edge cases

- Out of sample data

# How Do You Improve Datasets Efficiently?

- Problem: Identifying failure cases in the data / model performance

  - Lots of labeled data, only a few examples are problematic. Labor intensive to dig through haystack looking for the needle

  - Example: Triple QA finds many issues but can 3x your labeling cost

  - Example: Hard to understand model failure modes without metadata to slice on

# How Do You Improve Datasets Efficiently?

- Solution: Get feedback signal that tells you where to look

  - Feedback from double-checking

    - Human-check prod model outputs (example: customer feedback)

    - Check disagreement between automated systems

  - Feedback from model

    - High loss disagreements with labels (tend to be labeling errors)

    - Error patterns vs labels (in metadata + raw data)

    - Distributional shifts between training + prod environments

# Example: KITTI

# Example: KITTI

# Example: KITTI

# Long Tail Is Long



Faces [Zhang et al. 2017]

Places [Wang et al. 2017]

Species [Van Horn et al. 2019]

Actions [Zhang et al. 2019]

*Source: Z. Liu, Z. Miao, et al*

# Example: Oxford IIT Pets

# Example: Oxford IIT Pets

# Example: Oxford IIT Pets

# Example: Oxford IIT Pets

# Aquarium

Make it easier to build and improve production ML systems!

Q&A