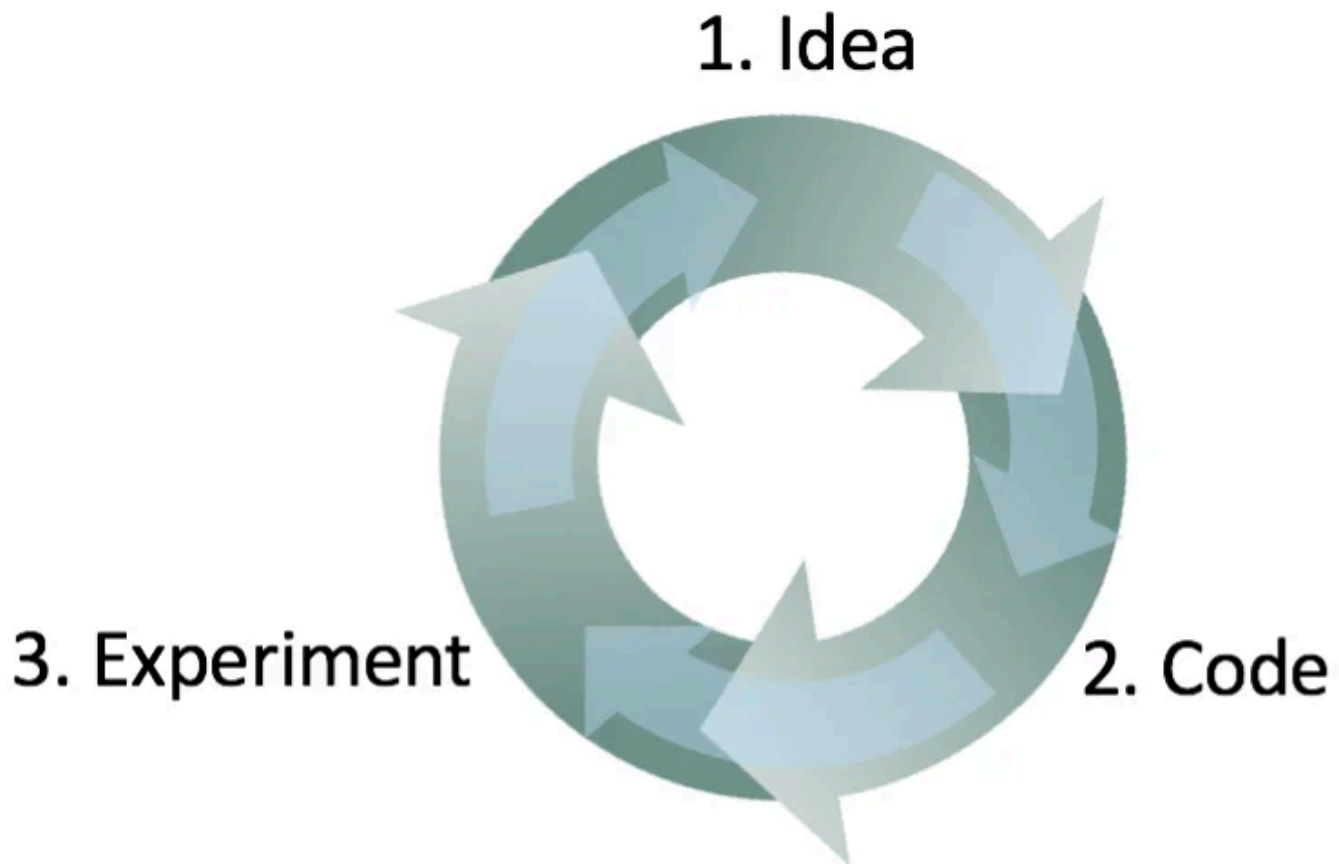


Why ML Strategy

- Ideas:
 - Collect more data
 - Collect more diverse training set
 - Train algorithm longer with gradient descent
 - Try Adam instead
 - Try bigger network
 - Try smaller network
 - Try dropout
 - Add L2
 - Change network architecture:
 - Activation Functions
 - # hidden units



Machine Learning Project Iteration Process

Orthogonalization

Orthogonalization is the process to adjust a given set of parameters to achieve a desired objective without changing the rest of the parameters. Orthogonalization can also be linked to interpretation of the model since it's pretty related to how much the model reacts to the inputs to model the output separately.

The level performance (at least on the training phase) should be at least as good as the human performance or close to it. Then the performance should fit well on the dev set to ensure that the training was correctly done and generalized well. Then the test set must also be tested to ensure that the performance on the dev set wasn't simply overfitted. Finally the performance in the real world should be assessed to ensure that the model really works.

To knob-adjust each step we can:

- Fit training set:
 - Train bigger network
 - Change the optimizer
- Fit the dev set:
 - Add regularization
 - Increase the training set size
- Fit the test set:
 - Increase the dev set so it really replies the distribution of the test set
- Perform well in the real world:
 - Change dev set
 - Change cost function

Single number evaluation metric

It's much better to have a single number to evaluate the performance of a given model. This number helps comparing across the models quickly. If you have to compare across multiple combinations or sources it might be useful to observe the model mean KPI, for example.

Satisficing and Optimizing metric

A Satisficing metric is a metric that cannot be overcome (as running time of the model inference) after the threshold you don't really care about the improving performance. Optimizing metric is the metric you want to maximize (or minimize if it's an error, for ex)

Train/dev/test distributions

The sets MUST have similar distributions. It's critical to have the sets distribution verified since the beginning and rechecked if the data changes over time to ensure that the distributions are not far in each set.

Size of the dev and test sets

Classic ML:

- 70/30 or 60/20/20

Deep Learning:

- 98/1/1

When to change dev/test sets and metrics

When the evaluation metric is no longer compatible with the chosen algorithm it's time to change the metric.

- $\text{Error} = \frac{1}{n_{\text{dev}}} \sum_{i=1}^{n_{\text{dev}}} L\{y_{\text{pred}}^{(i)} \neq y^{(i)}\}$
- Error weighted on pornographic images = $\frac{1}{w^{(i)}} \sum_{i=1}^{n_{\text{dev}}} w^{(i)} L\{y_{\text{pred}}^{(i)} \neq y^{(i)}\}$ where $w^{(i)}$ is defined as 1 if $x^{(i)}$ is non-porn and 10 if $x^{(i)}$ is porn

If doing well on your metric and your current dev sets or dev and test sets' distribution, if that does not correspond to doing well on the application you actually care about, then change your metric and your dev test set.

Why human-level performance?

Bayes optimal error is the very best theoretical error using a mapping function mapping X to Y.

Why compare to human-level performance?

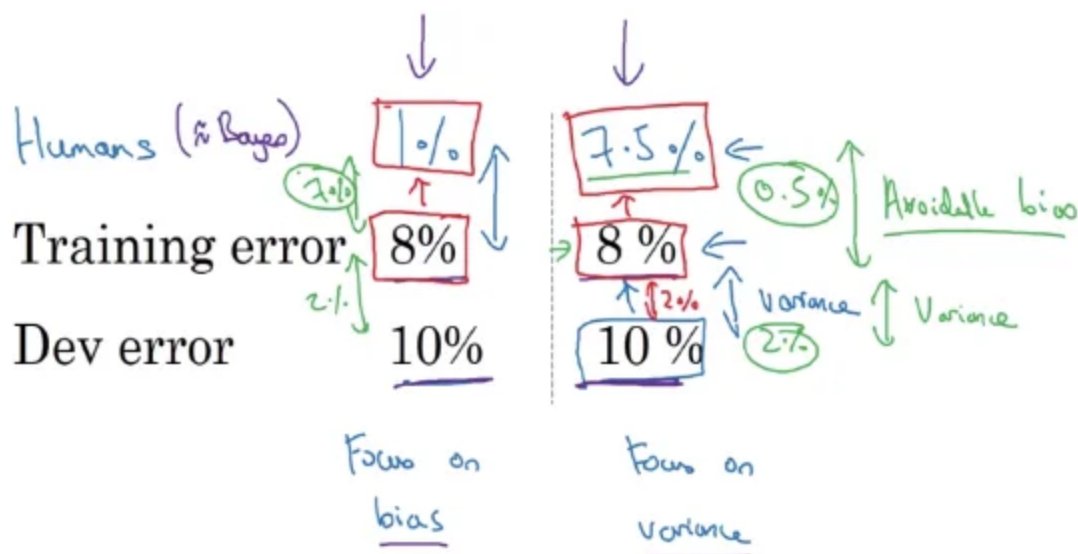
- Get labeled data from humans
- Gain insights from manual error analysis: Why did a person get this right?
- Better analysis of bias/variance

Avoidable bias

Human level error can be used as a proxy for **Bayes error**. If the human error is smaller than the current training error it might be worth focus in bias reduction. If the human level error is close to the training but far from the dev error there's much more room for improvement in the variance.

Avoidable Bias: The difference between the training error and the human-level performance. You don't really estimate to do much better than the bayes error without overfitting your data

Variance: Difference of error between the training error and the dev error



Avoidable bias and room for improvement in different scenarios

Understanding human-level performance

Humans tend to be very good at natural perception tasks. ML can surpasses human-level performance on structured data with huge amounts of data available

Improving your model performance

1. You can fit the training set pretty well, in other words, low avoidable bias
 - Train bigger model
 - Train longer
 - Better optimization algorithms
 - NN architecture search
 - Hyperparameter search
2. The training set performance generalizes pretty well to the dev/test set, in other words, low variance
 - More data
 - Regularization (L2, Dropout, Data Augmentation)

Quiz

1. Problem Statement This example is adapted from a real production application, but with details disguised to protect confidentiality. You are a famous researcher in the City of Peacetopia. The people of Peacetopia have a common characteristic: they are afraid of birds. To save them, you have to build an algorithm that will detect any bird flying over Peacetopia and alert the population.

The City Council gives you a dataset of 10,000,000 images of the sky above Peacetopia, taken from the city's security cameras. They are labelled:

$y = 0$: There is no bird on the image $y = 1$: There is a bird on the image

Your goal is to build an algorithm able to classify new images taken by security cameras from Peacetopia.

There are a lot of decisions to make:

What is the evaluation metric? How do you structure your data into train/dev/test sets? Metric of success The City Council tells you that they want an algorithm that

Has high accuracy Runs quickly and takes only a short time to classify a new image. Can fit in a small amount of memory, so that it can run in a small processor that the city will attach to many different security cameras. Note: Having three evaluation metrics makes it harder for you to quickly choose between two different algorithms, and will slow down the speed with which your team can iterate. True/False?

- ☒ True
- ☐ False

2. After further discussions, the city narrows down its criteria to:

“We need an algorithm that can let us know a bird is flying over Peacetopia as accurately as possible.” “We want the trained model to take no more than 10sec to classify a new image.” “We

want the model to fit in 10MB of memory.” If you had the three following models, which one would you choose?

- ☐ Test Accuracy Runtime Memory size 97% 1 sec 3MB
- ☐ Test Accuracy Runtime Memory size 99% 13 sec 9MB
- ☐ Test Accuracy Runtime Memory size 97% 3 sec 2MB
- ☒ Test Accuracy Runtime Memory size 98% 9 sec 9MB

3. Based on the city’s requests, which of the following would you say is true?

- ☒ Accuracy is an optimizing metric; running time and memory size are a satisficing metrics.
- ☐ Accuracy is a satisficing metric; running time and memory size are an optimizing metric.
- ☐ Accuracy, running time and memory size are all optimizing metrics because you want to do well on all three.
- ☐ Accuracy, running time and memory size are all satisficing metrics because you have to do sufficiently well on all three for your system to be acceptable.

4. Structuring your data Before implementing your algorithm, you need to split your data into train/dev/test sets. Which of these do you think is the best choice?

- ☐ Train Dev Test 6,000,000 1,000,000 3,000,000
- ☒ Train Dev Test 9,500,000 250,000 250,000
- ☐ Train Dev Test 3,333,334 3,333,333 3,333,333
- ☐ Train Dev Test 6,000,000 3,000,000 1,000,000

5. After setting up your train/dev/test sets, the City Council comes across another 1,000,000 images, called the “citizens’ data”. Apparently the citizens of Peacetopia are so scared of birds that they volunteered to take pictures of the sky and label them, thus contributing these additional 1,000,000 images. These images are different from the distribution of images the City Council had originally given you, but you think it could help your algorithm. You should not add the citizens’ data to the training set, because this will cause the training and dev/test set distributions to become different, thus hurting dev and test set performance. True/False?

- ☐ True
- ☒ False

6. One member of the City Council knows a little about machine learning, and thinks you should add the 1,000,000 citizens’ data images to the test set. You object because:

- ☐ A bigger test set will slow down the speed of iterating because of the computational expense of evaluating models on the test set.
- ☒ This would cause the dev and test set distributions to become different. This is a bad idea because you’re not aiming where you want to hit.
- ☒ The test set no longer reflects the distribution of data (security cameras) you most care about.
- ☐ The 1,000,000 citizens’ data images do not have a consistent $x \rightarrow y$ mapping as the rest of the data (similar to the New York City/Detroit housing prices example from lecture).

7. You train a system, and its errors are as follows (error = 100%-Accuracy):

Training set error 4.0% Dev set error 4.5% This suggests that one good avenue for improving performance is to train a bigger network so as to drive down the 4.0% training error. Do you agree?

- ☐ Yes, because having 4.0% training error shows you have high bias.
- ☐ Yes, because this shows your bias is higher than your variance.
- ☐ No, because this shows your variance is higher than your bias.

- ☒ No, because there is insufficient information to tell.

8. You ask a few people to label the dataset so as to find out what is human-level performance. You find the following levels of accuracy:

Bird watching expert #1 0.3% error Bird watching expert #2 0.5% error Normal person #1 (not a bird watching expert) 1.0% error Normal person #2 (not a bird watching expert) 1.2% error If your goal is to have “human-level performance” be a proxy (or estimate) for Bayes error, how would you define “human-level performance”?

- ☐ 0.0% (because it is impossible to do better than this)
- ☒ 0.3% (accuracy of expert #1)
- ☐ 0.4% (average of 0.3 and 0.5)
- ☐ 0.75% (average of all four numbers above)

9. Which of the following statements do you agree with?

- ☒ A learning algorithm’s performance can be better than human-level performance but it can never be better than Bayes error.
- ☐ A learning algorithm’s performance can never be better than human-level performance but it can be better than Bayes error.
- ☐ A learning algorithm’s performance can never be better than human-level performance nor better than Bayes error.
- ☐ A learning algorithm’s performance can be better than human-level performance and better than Bayes error.

10. You find that a team of ornithologists debating and discussing an image gets an even better 0.1% performance, so you define that as “human-level performance.” After working further on your algorithm, you end up with the following:

Human-level performance 0.1% Training set error 2.0% Dev set error 2.1% Based on the evidence you have, which two of the following four options seem the most promising to try? (Check two options.)

- ☒ Train a bigger model to try to do better on the training set.
- ☒ Try decreasing regularization.
- ☐ Get a bigger training set to reduce variance.
- ☐ Try increasing regularization.

11. You also evaluate your model on the test set, and find the following:

Human-level performance 0.1% Training set error 2.0% Dev set error 2.1% Test set error 7.0% What does this mean? (Check the two best options.)

- ☒ You have overfit to the dev set.
- ☐ You have underfit to the dev set.
- ☐ You should get a bigger test set.
- ☒ You should try to get a bigger dev set

12. After working on this project for a year, you finally achieve:

Human-level performance 0.10% Training set error 0.05% Dev set error 0.05% What can you conclude? (Check all that apply.)

- ☒ It is now harder to measure avoidable bias, thus progress will be slower going forward.

- ☐ This is a statistical anomaly (or must be the result of statistical noise) since it should not be possible to surpass human-level performance.
- ☒ If the test set is big enough for the 0.05% error estimate to be accurate, this implies Bayes error is ≤ 0.05
- ☐ With only 0.09% further progress to make, you should quickly be able to close the remaining gap to 0%

13. It turns out Peacetopia has hired one of your competitors to build a system as well. Your system and your competitor both deliver systems with about the same running time and memory size. However, your system has higher accuracy! However, when Peacetopia tries out your and your competitor's systems, they conclude they actually like your competitor's system better, because even though you have higher overall accuracy, you have more false negatives (failing to raise an alarm when a bird is in the air). What should you do?

- ☐ Look at all the models you've developed during the development process and find the one with the lowest false negative error rate.
- ☐ Ask your team to take into account both accuracy and false negative rate during development.
- ☒ Rethink the appropriate metric for this task, and ask your team to tune to the new metric.
- ☐ Pick false negative rate as the new metric, and use this new metric to drive all further development.

14. You've handily beaten your competitor, and your system is now deployed in Peacetopia and is protecting the citizens from birds! But over the last few months, a new species of bird has been slowly migrating into the area, so the performance of your system slowly degrades because your data is being tested on a new type of data. You have only 1,000 images of the new species of bird. The city expects a better system from you within the next 3 months. Which of these should you do first?

- ☒ Use the data you have to define a new evaluation metric (using a new dev/test set) taking into account the new species, and use that to drive further progress for your team.
- ☐ Put the 1,000 images into the training set so as to try to do better on these birds.
- ☐ Try data augmentation/data synthesis to get more images of the new type of bird.
- ☐ Add the 1,000 images into your dataset and reshuffle into a new train/dev/test split.

15. The City Council thinks that having more Cats in the city would help scare off birds. They are so happy with your work on the Bird detector that they also hire you to build a Cat detector. (Wow Cat detectors are just incredibly useful aren't they.) Because of years of working on Cat detectors, you have such a huge dataset of 100,000,000 cat images that training on this data takes about two weeks. Which of the statements do you agree with? (Check all that agree.)

- ☐ Having built a good Bird detector, you should be able to take the same model and hyperparameters and just apply it to the Cat dataset, so there is no need to iterate.
- ☒ Needing two weeks to train will limit the speed at which you can iterate.
- ☒ If 100,000,000 examples is enough to build a good enough Cat detector, you might be better off training with just 10,000,000 examples to gain a $\approx 10\times$ improvement in how quickly you can run experiments, even if each model performs a bit worse because it's trained on less data.
- ☒ Buying faster computers could speed up your teams' iteration speed and thus your team's productivity.