# Deep Learning Specialization: Strategy

Why ML strategy?

→ If results are not satisfying, there are __LOTS__ of things you could try.

→ Many things will take time / be expensive

ML Strategy: WHAT should you try to improve?
Effective ML engineers are those who know this.
     ↳ much like in software!

## ORTHOGONALISATION

- Decompose potential actions into orthogonal directions
  ↳ Know what actions to take based on the problem you have.

4 Core assumptions in ML:

① Fit training set on cost function well
   · More powerful network (layers, units, BN, activations, ...)
   · Better optimisation (Algo, LR, longer runs, more compute

② Fit dev set on cost function well
   · Regularization
   · Bigger training set

③ Fit test set on cost function well
   · Bigger dev set?

④ Perform well in real world.
   · Different cost function?
   · Bigger dev set

# Single number evaluation metric

- Use a single number to evaluate & compare models.
  - ⇒ Makes selecting among multiple candidates much easier
  - ⇒ Faster to iterate and shoot for a target.
- Ex: Precision + Recall vs F1 score
  - ↳ Harmonic mean of P & R.

# Optimizing vs Satisficing metrics

Optimizing: Maximize/minimize this number
Satisficing: Guard rails. Meet a certain threshold.
Ex: Accuracy vs Running time

* When it is hard to hit all metrics you care about into a single metric, make one optimizing and the others satisficing.

## Train / dev / test set distributions

Make sure <u>dev</u>, <u>test</u> (and prod) come
from same distribution!

-> otherwise moving the target!

Examples:
- Dev set  US, Uk, Europe,  Test: India, China, Asia
- Dev: Middle income zip codes  Test: Low income zip
- Dev: Cat pics from internet  Test: Cat pics from users

Guideline: Choose dev <u>and</u> test set to reflect
data you expect and is important to do
well on.


## Dev & Test set sizes

- Dev set big enough to detect differences in
algos you train.

- Test set big enough to give high confidence
in perf of system.

## When to change your cost function
- If it does not capture well what you really want to achieve
- If doing well on metric in dev and test does not mean doing well in your application
  - → change it!

## Comparing to human level performance.

- Bayes error: Error on best possible performance
  → likely > 0%
- Human level performance:
  → Our best proxy on a lot of tasks
  → Humans are quite good at a lot of tasks, esp. natural perception.
- While worse than humans, you can
  → Get labled data from humans
  → Gain insight from manual error analysis
    → Why did a person get this right?
  → Better analyse bias/variance

## Avoidable Bias

Cat classification         Error

Humans (proxy for Bayes)       0.5 %

                                 $\updownarrow$ 7.5% avoidable bias

Training   error               8 %

                                 $\updownarrow$ 2% variance

Dev     error               10%


## Human Level Error

Ex:
- Typical human:     3%
- Typical doctor:     1%
- Expert doctor     0.7%
- Team of experts    0.5%

       $\hookleftarrow$ Proxy for Bayes error

# Surpassing human level performance

→ Not impossible
→ Can be hard, depending on problem

Exmples
- Online advertising
- Product recommendations
- Logistics ( prod. transit time)
- Loan approvals

→ Lots of structured and quantidative data

# Improving your model performance
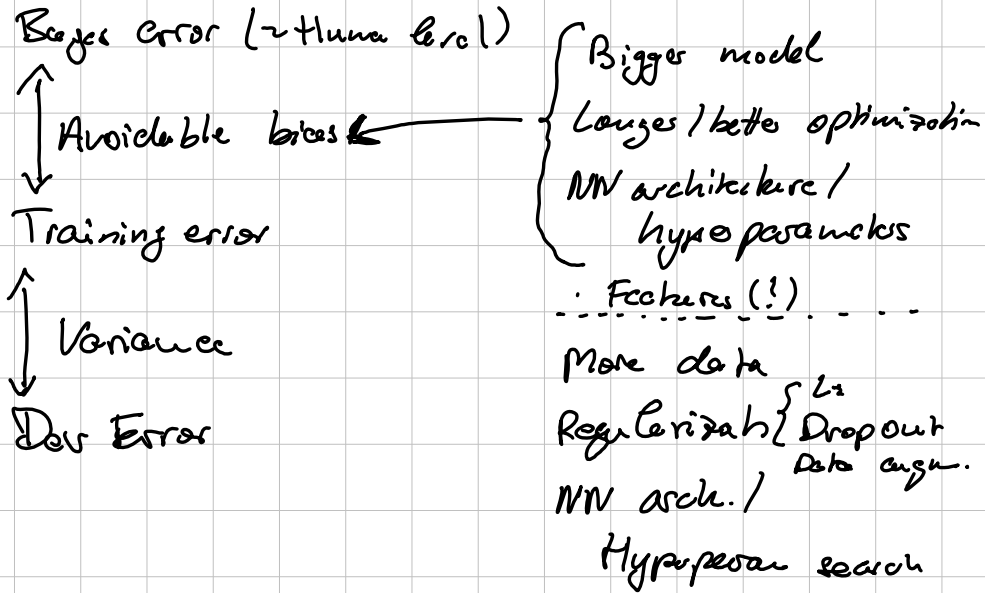
Two fundamental assumptions of supervised
learning:
1) You can fit training set well
   => low avoidable bias
2) Training set performance generalizes
   well to the dev & test set
   => low variance

# Reducing (avoidable) Bias and Variance

Bayes error (~ Human level)

↑↓ Avoidable bias ←————————— {
Bigger model
Longer / better optimization
NN architecture /
    hyperparameters

Training error

· Features (!) · · · · · · · ·

↑↓ Variance

More data
Regularization { L₂
         Dropout
         Data augm.

Dev Error

NN arch. /
   Hyperparam search

# Week 2: Error Analysis

Q: Should you work on making your cat
   classifier better on dogs?
   → Evaluate ROI!
      Look at 100 misلabeled examples →
      how many are dogs. Is it worth it?

· Evaluate multiple ideas in parrallel
   (e.g. dogs, big cats, blurry pictures, ...)
· Classification spread sheet.

## Cleaning up wrong lables

Q: Should clean up wrong lables in
   training set?

A: DL algorithms robust to _random_ errors in
   training set. (Systematic errors are a problem)

On dev set: Impact analysis: Will fixing this help
   selecting between two different algos?

- Apply same method of correcting to dev and test set. These need to come from the same distribution!

---

Strategy: Build first system quickly, then iterate!

---

## Training & testing on different distributions

- Can be ok. But algo might struggle on data it has not seen at all
- Can augment the train dataset?
  - be careful when synthesizing. Can work but ensure it's not too monotonic.

How to ass.s bias / variance when training &
testing on different sets?

→ Train - dev set: ¬ same distribution as
  train set, but used for testing
  → Shows variance of model on distribution
    it was trained on
  → Ass.es: overfitting or learning doesn't generalize
    ?

Bayes error
            } avoidable bias
Training error
            } variance
Training - dev error
            } data mismatch
Dev error
            } overfitting the dev set
Test error

|  | Train data | Real data |
| --- | --- | --- |
| Human level | Human level 4% | 6% |

Human level       Human level 4%     6%

↑3%           } avoidable bias

Error on data      Train error 7%     6%
trained on

↑3%           } variance

Error on data       Train-dev 10%    Dev/Test 15%
<u>not</u> trained on       error          error

←—————→
data
mismatch
5%

## <u>Addressing data mismatch</u>

- Carry out error analysis to understand difference between training & dev/test sets
- Make training data more similar

# Transfer Learning

Swap out last layer *or maybe add a few layers* of an already trained neural net and train (only last, or all) on new task.

Depending on how much new data you have (for target task) train only last or all layers.

little data — lots of data

Makes sense if:

- Input data to the tasks is the same
- You have (a lot) more samples for the task you're transferring from.
- Tasks are similar / some reason to suspect things will transfer. low level features from Task A could be helpful for Task B.

# Multi task Learning

- Learn multiple tasks at the same time,
  e.g. multi-object classification.

- Multi-output layer.

Loss function: sum over individual losses
$$\frac{1}{m} \sum_{i=1}^{m} \sum_{j=1}^{4} \mathcal{L}(y_j^{(i)}, \hat{y}_j^{(i)})$$

- If you have some examples where you only
  have one / a few things labeled, then
  don't add the loss for those examples.

## When it makes sense

- The different tasks could benefit from
  similar low-level features
- Amount of data for each task is
  ~ similar
- Can train NN big enough to handle all
  tasks.