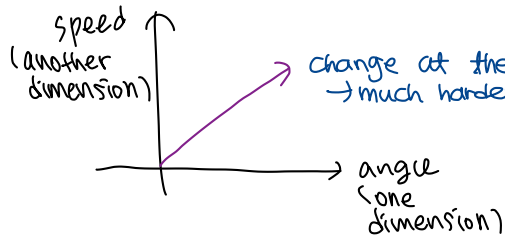


Orthogonalization

<TV / Car Example>

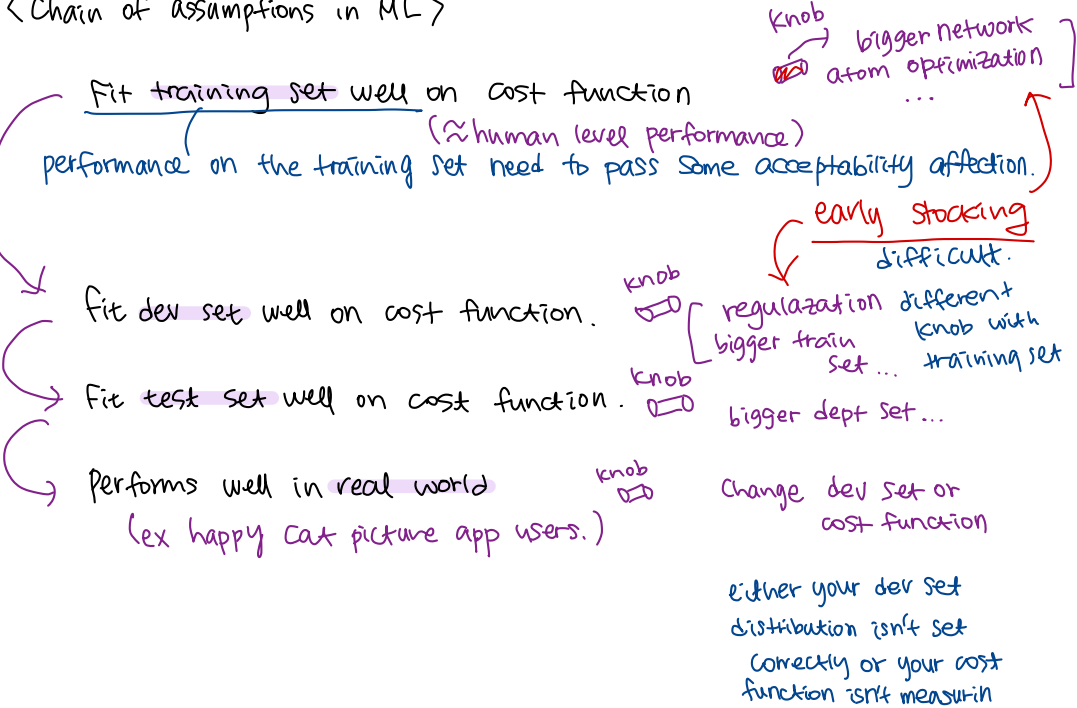
Q2

knobs → change



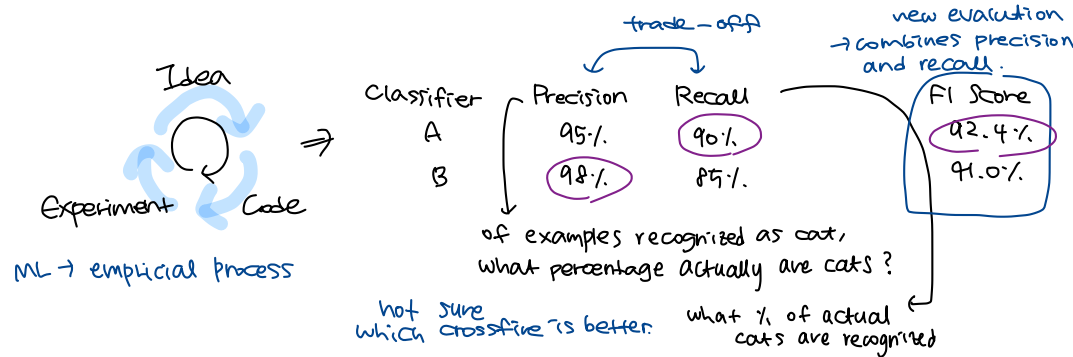
in theory, by tuning these two knobs
→ get your car to steer the angle and
the speed you want much harder
if you just have one
single control for
controlling a steering angle

<Chain of assumptions in ML>



Single Number Evaluation Metric

<Using a single number evaluation metric>



F1 Score = "Average" of precision P and recall R

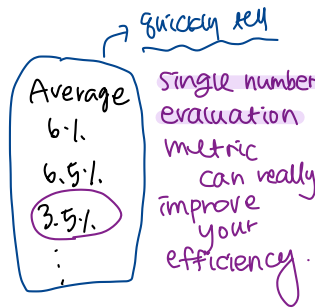
$$\left(\frac{2}{\frac{1}{P} + \frac{1}{R}} \right), \text{ "Harmonic mean" }$$

Dev Set + Single number evaluation metric
real ⇒ Speed up iterating.

<Another Example>

Algorithm	US	China	India	Other
A	3%	7%	5%	9%
B	5%	6%	5%	10%

testing a lot of different compilers
→ difficult to look at all these numbers
and quickly pick one.



Satisficing and Optimizing metrics.

more important: optimizing
less : satisficing

< Another Cat Classification Example >

optimizing metric

you also care about it..

Classifier	Accuracy	Running Time	Satisficing metric
A	90%	60 ms	
B	92%	95 ms	
C	95%	1500 ms	

Cost = accuracy - 0.5 x running Time

maximize accuracy
subject to running time \leq 100 ms.

be less than 100 ms and beyond that you don't really care / don't care that much.

N metrics: 1 optimizing
N-1 Satisficing

ex) Wakewords / Trigger words - Alexa
OK Google
accuracy
false positive
Hey Siri...

maximize accuracy - optimizing metric
s.t. \leq 1 False positive - satisficing metric
every 24 hours.

Train / dev / Test Distribution

how you set up your dev and test set.

< Cat classification dev / test sets >

Regions:

- US
- UK
- Other Europe
- South America
- India
- China
- Other Asia
- Australia

development set, hold out cross validation

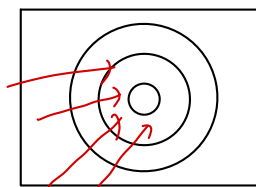
very bad idea \rightarrow dev & test come from very different distributions.

Dev

very different

Test

not giving you a good performance on the test set (different location)



dev set + evaluation metric : tell you a goal

take all these data preventing the shuffle data into the dev / test set
Randomly shuffle into dev / test

< True Story (details changed) >

[Optimizing on dev Set on loan approvals for medium income zip codes

X (about loan application)

↓ predict

Y (whether or not they'll be paying alone)
(repay loan?)



Tested on low income zip codes didn't work at all.

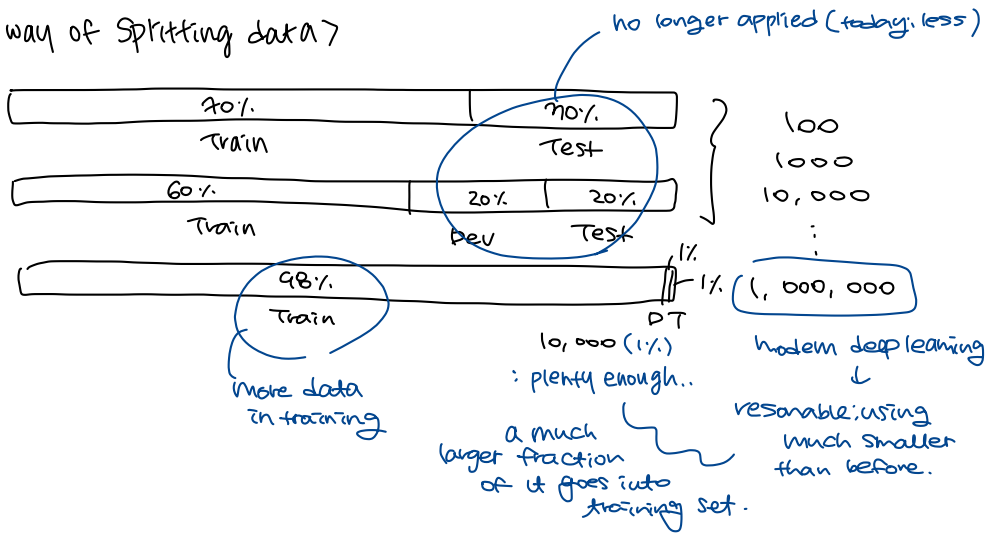
< Guideline 7

Same distribution.

Choose a dev set and test set to reflect data you expect to get in the future and consider important to do well on.

Size of Dev and Test Sets

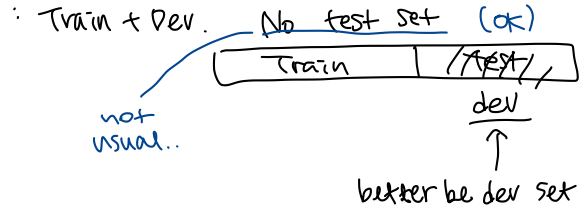
<Old way of Splitting data>



<Size of Test Set>

→ Set your test set to be big enough to give high Confidence in the overall performance of your system.
 10,000
 100,000 ...
 to your purpose..

Some applications: don't need high confidence



When to Change dev/test sets and metrics

<Cat dataset Examples>

→ Metric: classification error
 Algorithm A: 3% error
 ✓ Algorithm B: 5% error
 B is better: not learning though any pornographic images.
 (even A has better evaluation metric)
 time to think about defining new evaluation metric
 → pornographic images (totally unacceptable)
 no longer directly rank ordering performance. miss predicting

Metric + dev: prefer A (lower error)
 You/Users: prefer B

(just one possible way)

Error:
$$\frac{1}{N} \sum_{i=1}^N w^{(i)} \mathbb{I}(\hat{y}^{(i)} \neq y^{(i)})$$

 This formula just count up the number of misclassifying examples.
 → $w^{(i)} = \begin{cases} 1 & \text{if } x^{(i)} \text{ is non-porn} \\ 10 & \text{if } x^{(i)} \text{ is porn} \end{cases}$
 give a much large weight to examples that are porn.

<Orthogonalization for cat pictures: anti-porn>

1. So far we've only discussed how to define a metric to evaluate classifiers. place the target. ↻ distinct step
2. Worry Separately about how to do well on this metric. ↻ aim/shoot at the target.

→
$$J = \frac{1}{N} \sum_{i=1}^N w^{(i)} \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

 changing cost function J that your network is optimizing.

✓ If doing well on your metric + dev/test set does not correspond to doing well on your application, change your metric and/or dev/test set.