

Course 3

STRUCTURING ML Project

Orthogonalization

- many potential tools

Fit training set well on cost function

- bigger network

- Adam

- :

↗
X early stopping
↖ X.

Fit dev set well on cost function

- Regularization

- bigger training set

Fit test set well on cost function

- Bigger dev set

perform well in real-world

- change dev set or cost function

Single Number Value Evaluation

precision → recognized cat, % of them as cat

recall → % of actual cat got recognized

Competing... hard to evaluate

⇒ F1 score ~ average of P & R

$$\frac{2}{P+R} \text{ "Harmonic Mean"}$$

Satisficing & Optimizing Metric

e.g.

maximizing accuracy

subject to running time < 100ms

Optimizing

Satisficing

e.g. wake word / trigger word

'hi siri'

maximize accuracy

s.t. ≤ 1 false positive every 24 hrs.

Train / Dev / Test set

Dev / Test should have same distribution

choose Dev/Test that reflect your future target

Size of dev/test sets

Old

170%	30%
Train	Test
60%	20%
Train	D

100
S
10,000

1,000,000 data	98%	12%
Train	Dev	test

high

size of test set ~ big enough to give confidence on performance

- [train | dev] ← sometimes not recommended
- [train | test] X not good naming...

When to change Dev/Test and Metrics

real data
different from expectation

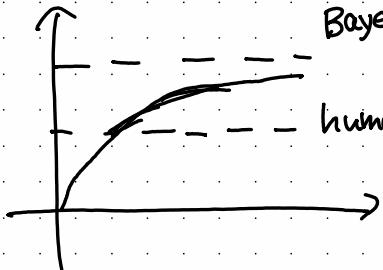
- [metric + dev prefer A → time to change
your user prefer B e.g. bigger penalty on
undesired image]

orthogonality
- place target (define metrics)

- shoot at it (change cost func... etc)

human-level performance

Bayes optimal error \leftarrow best possible



\rightarrow progress slows when
surpass human

if ML is worse

- get more labels
- gain insight from manual analysis
- better analysis of bias / variance

avoidable bias

\triangle assume human-level
error \approx bayes error
(esp. CV tasks)

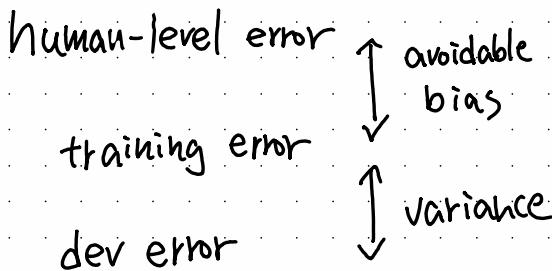
\times human	1%	7.5%
training error	8%	8%
dev error	10%	10%
	↑	↑
reduce bias		reduce variance

Avoidable bias

define human-level error

human-level error as proxy to bayes error.

e.g. human doctor
theoretical upper bound
group of doctors ← closer to bayes error



surpassing human-level error

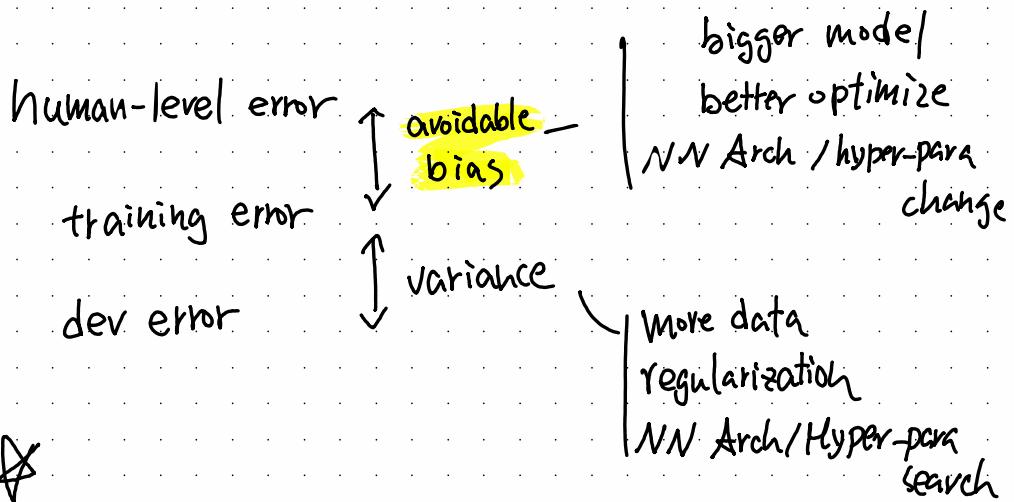
- online ad
- product recommend
- logistic

speech recog
medical data
:

learned from structure data
not involve perception
lots of data

fundamental assumption of supervised learning

- one can fit training set well
- training set generalizes to dev/test.



Error Analysis

cat vs dog. if find mislabelled ones are dog.

Analysis

1. get ~100 mislabeled dev set example
2. count how many are dogs
 - if low, probably not worth it.

▷ build table for incorrectly labelled dev set.

Image Dog big cat blurry .. ↪ ih

1	✓
:	
n	✓

%



based on percentage, decide what's next step.

Incorrect Labelled Data

- robust at "random" error in training data
- for dev/test, label accuracy might be less important.
(based on the err analysis table)

Correcting dev/test labels

- apply same process for both dev/test set
(ensure same distri)
- consider examples your algorithm got "right" as well as "wrong"
- note training and dev/test might come from different distri

Build system quickly then iterate

- set dev/test set early
- build
- bias/var/error analysis
 - esp for new applications

Training / Dev from different distri'

e.g.

Cat from webpage web from mobile (blur etc)

~200,000

~10,000

what we care

option:

1. $200k + 10k = 210k$ (X)

train dev/test

$$2,500 \rightarrow 2500 \times \frac{10}{210} = 119 \text{ mobile image.}$$

2.

train dev/test (v)

200k + 5k
web mobile

2500 2500
mobile photos

Bias & variance with mismatched data dist.

training

train dev/test

train on this

training-dev set

human err 0%

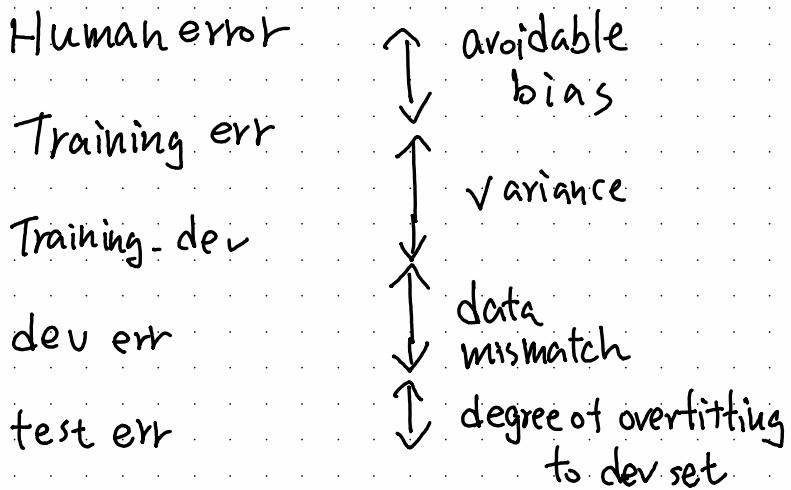
Train ↑ 1% 1% 10% 10%

same dist as training set

train-dev ↓ 9% 1.5% ↑ 11% 11%

dev 6% 10% ↓ 12% 20%

Variance data bias bias
problem mismatch problem + data mismatch



Sometimes dev performs better than training

- e.g. dev/test are of different distribution
- additional label on the dev/training might better understand human performance

Addressing Data Mismatching

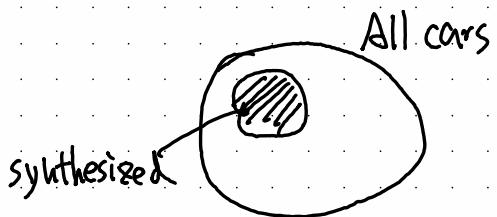
- manual error analysis to understand the diff between training, dev/test.
 - making training data more similar

Artificial Data Synthesis

e.g. speech + car noise = synthesised in-car speech

note ↑ ↑ ↑
10,000 hrs 1 hr might be overfitting
the 1 hr car noise

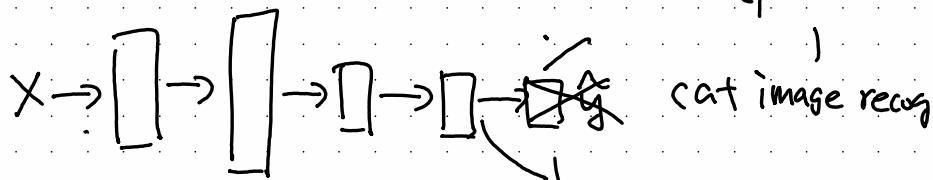
e.g. car photo



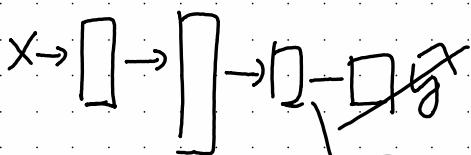
8 if using synth data, make sure it capture enough variety.

Transfer Learning

(pre-training)



or

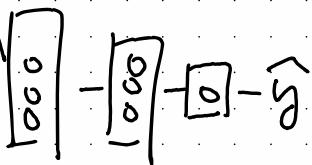


\hat{y} radiology

$w^{[L]}, b^{[L]}$

(fine-tuning)

Working bcz low-level
feature might be
similar.



Transfer from A \rightarrow B, useful when.

- task A & B have the same input
- much more data in A
- lower level feature of A is helpful for B

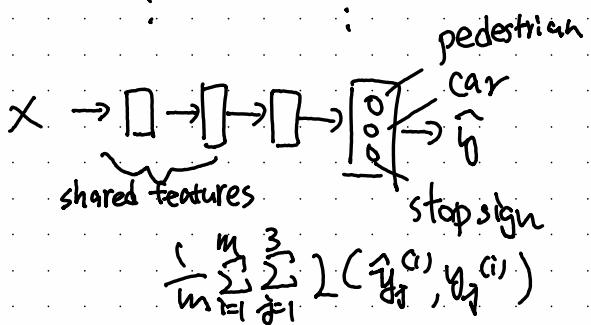
Multi-task Learning

(less often used)
except CV. compared to TL

e.g. autonomous driving

pedestrian	1
Cars	0
stop sign	0
:	:
:	:

▷ one image with
multiple labels



$$\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^3 L(\hat{y}_j^{(i)}, y_j^{(i)})$$

$$-y_j^{(i)} \log \hat{y}_j^{(i)} - (1-y_j^{(i)}) \log (1-\hat{y}_j^{(i)})$$

▷ might still work when some images only have
partial labels

MAKE SENSE WHEN

- tasks share similar low-level features
- amount of data for each feature are similar
- can train a big enough NN to do well on all tasks

End-to-End DL

X Audio $\xrightarrow{\text{MFCC}}$ features $\xrightarrow{\text{ML}}$ phonemes \rightarrow words \rightarrow transcript
audio \longrightarrow transcript

need a lot of data!

e.g.

Face recognition

image $\xrightarrow{\text{tracking}}$ face $\xrightarrow{\text{ML}}$ recognition
face detection

& have lot data for each of 2 sub-tasks.
even more data for end-to-end data.

Machine Translation

X Eng \rightarrow text analysis $\rightarrow \dots \rightarrow$ French

✓ Eng \longrightarrow French

Estimate child age

✓ image $\xrightarrow{\text{①}}$ bone $\xrightarrow{\text{②}}$ age

X image \longrightarrow age

wheather to use End-to-End DL

Pro :

- let the data speak
- less hand-design features

Con :

- data hunger
- exclude hand-designed features

↳ hand-design feature helps when training set is small

KEY Q:

sufficient data to learn a function maps X to Y ?