# Structuring ML Projects :

⊕ Introduction to ML stratergy.

*orthogonalisation

— what to tune in order to acheive one effect, people are clear about t
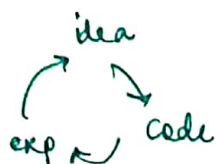
tune ~~measure~~ different parameters differently

chain of assumptions in ML :

{ different
tuners
for
different
assumptions
independent
of each other }

- Fit training set well on cost funcⁿ. ⇄ → bigger N/w
  → Better opt. algo
  ↓
  Fit dev set well on cost func ⇄ → Regularization
  → Bigger training set
  ↓
  Fit test set well on cost funcⁿ. ⇄ → bigger dev set
  ↓
  Performs well in real word, ⇕→ change dev set or cost fun

↓↓
orthogonalisation.

* Single number evaluation metric.

idea



exp ← code

→ ? examples as cats, what % are cats
→ what % of actual cats recognise

precision  recall

| | precision | recall |
|---|---|---|
| A | 95% | 90% |
| B | 98% | 85% |

f Score = "Avg" of precision P &
recall R

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

\* Satisficing & optimizing metric :

Accuracy → satisficing

Run time → optimizing

| | Accuracy | Run time |
|---|---|---|
| A | 90% | 80 ms |
| B | 92% | 95 ms |
| C | 95% | 1500 ms |

let cost = accuracy - 0.5 × running time

maximise accuracy

t ≤ 100 ms
⇕
has to be satisfied.

N metrics : 1 optimizing
N-1 satisficing } best.

wakewords (trigger words

Alexa, ok google, Hey sri , ----

accuracy ⟶ maximize
\# false positive ⟶ atmost ≤ 1 per 24 hrs. (satisficing)

\* <u>Train, dev, test set distributions</u> :

dev / test
⇕
develop set
holdout cross valid$^n$ set

Regions :

{ US
UK
Other Europe
South America } DEV

{ India
China
Other Asia
Australia. } Test.

dev set + metric
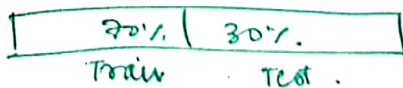
due to diff in data, optimisation for 1 can degrade other.

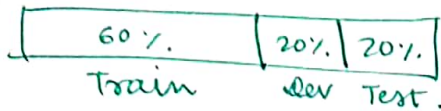⟶ Randomly shuffle data, so that better distribution exists.

- choose dev set and test set to reflect data you expect to get in future &
consider important to do well on
dev & test st shud be from same distribution.

* size of dev & test set :
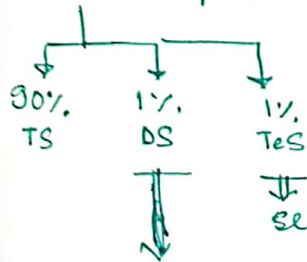
How long should they be.

- old way of splitting data.

| 70% | 30% |
|-----|-----|
| Train | Test |

reasonable if ~ 1000 examples.

| 60% | 20% | 20% |
|-----|-----|-----|
| Train | dev | Test |

- million examples

90% TS    1% DS    1% TeS

set big enough to give high confidence in overall performance

Some appln : Train + Dev, no Test set.

* when to change dev/test set metrics :

Metric : classif$^n$ error

A : 3% error → has porn (better on eval metric, but worse)

B : 5% error → doesn't have porn (poor, but good as no porn)

metric + eval → A
    user → B

∴ change evaluation metric.

$$Error = \frac{1}{m_{dev}} \sum_{i=1}^{m_{dev}} w^{(i)} \{ y_{pred}^{(i)} \neq y^{(i)} \}$$

(0/1)

$$w^{(i)} = \begin{cases} 1 & ; x^{(i)} \text{ not porn} \\ 100 & ; x^{(i)} \text{ porn}. \end{cases}$$
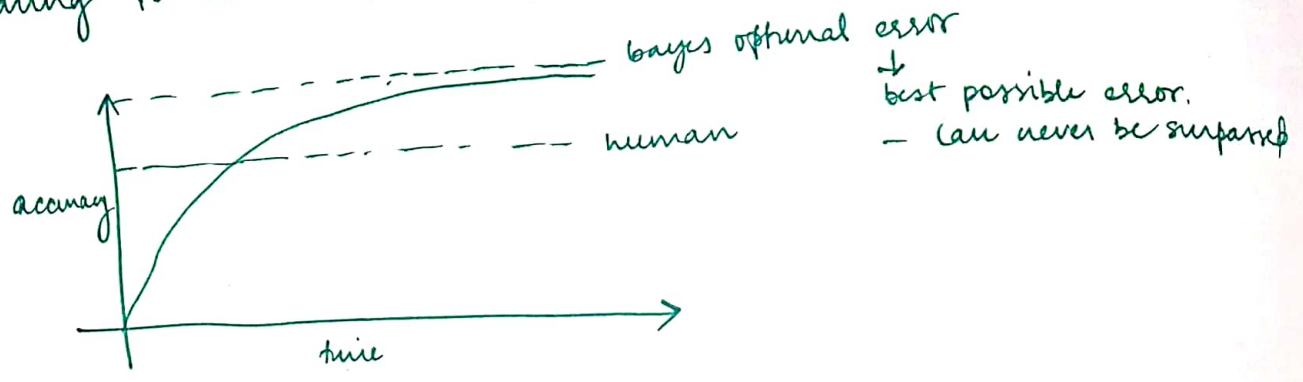
∴ when weights introduced.

new metric defined.

$$Error : \frac{1}{\sum w^{(i)}} \sum_{i=1}^{m_{dev}} w^{(i)} \{ y_{pred}^{(i)} \neq y^{(i)} \}$$

$$J = \frac{1}{\sum w^{(i)}} \sum_{i=1}^{m} w^{(i)} L(\hat{y}^{(i)}, y^{(i)})$$

- if doing well on metric + dev/test set doesn't correspond on doing well on your application, change your metric.

⊕ Comparing to human level performance :-



bayes optimal error
↓
best possible error.
- can never be surpassed

$X \longrightarrow y$
audio    transcript.

- progress slows down after human level performance.
- progress never exceeds bayes optimal error.

# Avoidable bias

| | | |
|---|---|---|
| Human | 1%. | ↑ to reduce |
| Training error | 8%. | |
| Dev error | 10%. | |

Focus on <u>reducing bias</u>.

| | | |
|---|---|---|
| 7.5%. | | ↑ Avoidable bias |
| 8%. | ↑ to reduce | ↕ |
| 10%. | | ↓ variance. |

Focus on variance +
by regularizⁿ,
more training data.

Human level error as a proxy for bayes error.

# Understanding human level performance :

Ex:   Typical human   — — —   3% error
    "   Doctor   — — —   1% error
    . experienced doctor   — — — 0.7% error
    Team of exp. doctors   — — — 0.5%.

~~what is~~ "human level" ⟵ Baye's error ≤ 0.5%.
     error : 0.5%

Human (proxy for bayes error) = 1%, 0.7%, 0.5%.

Avoidable ↑ bias | Training error = 5%          Focus on bias red^n :
variance ↓        | Dev       error = 6%.

human error = 0%.
      (bayes error)        1%, 0.7%, 0.5%          TE: 1%.          DS: 5%.
Avoidable bias                      ⟵————————⟶      ⟵————————
becomes bias.                                              ↓se the
                                                     variance (regulaiz^n,
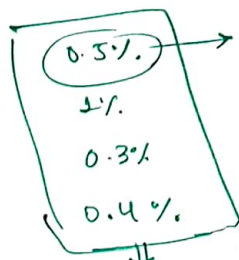                                                           more Train
                                                             data)

* Surpassing Human level performance :-

                                        ┌─────────┐ → progress options
Team of humans      0.5%                │  (0.5%) │    less clear.
one human  ↑A-B↓    1%                  │   1%    │
                    0.6%                │   0.3%  │
TE                  0.8%                │   0.4%  │
                                        └────↧────┘
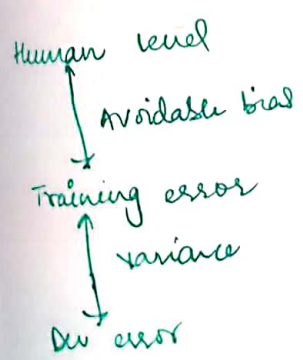DE                                      better performance than humans.
                                online advert, product recommend^n, logistics, loan
                                                                    approvals.

* Improving model performance :

1. Fit TS well ⟷ achieve low avoidable bias
2. TS performance generalises well to dev & test set ⟶ Variance

Human level                            ⎧ Train bigger model
    ↑                                  ⎪ Train longer / better optin's algo
    │ Avoidable bias      ⟶            ⎨ NN architecture, ---
    ↓                                  ⎩
Training error
    ↑
    │ variance            ⟶            ⎧ more data
    ↓                                  ⎨ Regularisation .
Dev error                              ⎩

⊕ Error analysis :

* Carrying out error analysis :

    Ex: 90% accuracy
        10% error .

    Error analysis :                                    " gives a ceiling as to
    manually · get ~100 mislabeled dev set examples        how best can
    takes 5-10 · count how many dogs.                      algo work "
                                out of 100
                                  50% → dog          10% → 5%

Scanned by CamScanner

Evaluate multiple ideas in 11d:

| Img | Dog | Cat | Great Cats | Blurry | Insta | Comments |
|-----|-----|-----|-----------|--------|-------|----------|
| 1 |  |  |  |  | ✓ | Pitbull. |
| 2 |  | ✓ |  | ✓ |  | Rain at zoo |
| 3 |  |  |  |  |  | : |
| : |  |  |  |  |  |  |
| % of total | 8% |  | 93% | ··· | 61% | 12% |

can improve performances through these.

- find a set of mislabelled examples in dev set & look for false +ue & false -ve.

* cleaning up mislabeled data :
- if errors are random, then its OK (% not too high)
- robust to random errors
- Systematic label errors are bad.

⇒ Add one "incorrect label examples" column.
    " 6%."

if makes significant difference, then fix the error, else don't bother.

   Ex : overall dev set error   10%.
             incorr labels = 6% of 10% = 0.6%.
          then : 9.4%.
            ∴ don't fix labels.

Correcting dev set examples (which R incorrect)
   • dev & test set should be processed similarly from same   distribution
                          (apply same process to both)
   • consider examples that algo got right, as well as wrong. (isn't always
                                                                    done)
   • Train & dev/test may come from different distribution.

* Build your first system quickly, then iterate :
   ⎰ • set up dev/test set & metric
   ⎱ • Build initial system quickly.
     • Use Bias/variance analysis & error analysis to prioritise next step.
   ↳ Build first system quickly, then iterate .

training & testing on different distributions:

ex: webpage : 200,000 (good qlty) ⎫ mix them
    phone app: 10,000 (bad qlty) ⎭ 210,000
                                        ↓
                                     shuffle,
                                  split them randomly
                                     in train, dev, test set.

ex: dev - 2500
    test - 2500

dont use this , as algo will optimize for web images most of the time.
Instead ,

    Training 200,000 → web + 5000 → app.
        dev/test: all mobile app.
        advantage: aiming the tgt.

            as we have to make mobile app
                        ∴ this is better.

• Speech recog examples.

Training                          Dev/Test
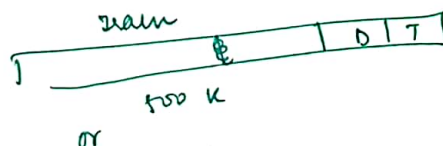
Purchased data                    speech activated ⎫
smart speaker control             Reartrived mirror ⎭
Voice keyboard.                   → 20000
. . .
    500,000 occurences.
                                  train
                                  ├──────────────┼───┤
                                  │              │ D │ T │
                                                500 K
                                  or

                                  ├──────────────────────┼────────┤
                                         500 K + 10K
                                              ↓         → also a
                                            from          good
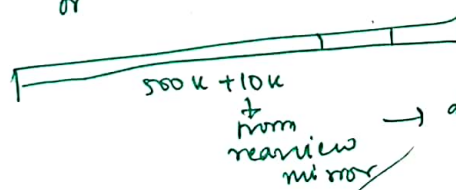                                          rearview        thing.
                                           mirror

* Bias and Variance with mismatched
  data distribution .

ex : Assume humans ≈ 0%.
     train err : 1% ↓ distribution is different ∴ difficult to analyze err.
     Dev err : 10%           ↓
                          define  Training-dev set : same dist as Tr.set
                                              but not used for
                                                   training.

                randomly shuffle Tr.S and
                           came out small
                            tr. dev set .

Eg: Training error = 1%.
Training -dev err = 9%.
Dev error = 10%.

↑ variance prob (as same distrib'n)

↓ 0% (Bayes)  ↑ Avoidable bias prob.
↓ 10%.
11%.
12%.

Tr    1%.    ↓ low variance prob
Tr-Dev 1.5%. ↑ data mismatch prob.
Dev   10%.   ↓

0%.  ↑ Avoidable bias
10%.     as well as
11%. ↓ large data mismatch.
20%.

Thus,

Human level error    4%   ↑ Avoidable bias
Tr set error         7%.  ↓ variance
Tr - Dev set error.  10%  ↓ mismatch
Dev error            12%  ↓ degree of overfitting
Test -error          12%. ↓

{ → same distribution.

ut.    Human    ---    4%.
       Tr set   ---    7%.
       Tr - Dev ---    10%.
       Dev      ---    6%.  } If dev/test set examples are
       Test     ---    6%.  }          much easier.

more general formulation: (Rear view mirror)

|                        | general speech recog | Rearview mirror speech data |
|------------------------|----------------------|------------------------------|
| Human level            | 4%. ↓ avoidable bias | 6%.                          |
| error on trained orea  | "Tr error" 7% ↓ variance | 6%.                      |
| err on not trained on ex. | "Tr dev set" 10%  data mismatch | "Dev/Test" error 6% |

Already doing quite well.

* Addressing data mismatch :
  • carry out manual error analysis to try to understand difference between training & dev/test set.
  • Make data more similar ; or collect more data similar to dev/test set.

# Artificial data synthesis:

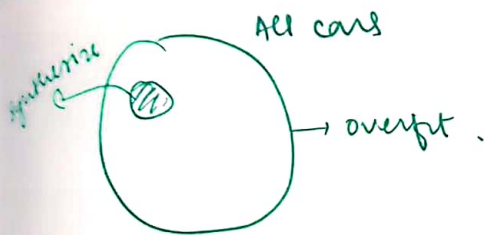"the quick brown fox jumps over the lazy dog" + "car noise" = synthesised in-car audio.

↑
10,000 hrs

↑
1 hour

sol^n: use 10,000 hrs car noise (possible).

opt 1: 1 hr, 10 000 times repeated. learning algo can overfit to 1 hr synthetic car noise.
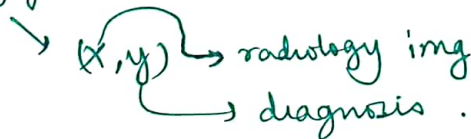
synthesise → all cars → overfit.

car recognition

synthesise → All cars → overfit.

- if data mismatch, do error analysis, get more Tr. data, etc...

## Learning from multiple tasks:
### Transfer learning:-

ex: Trained on img recog.

↓ change last layer radiology.



→ (x,y) → radiology img
→ diagnosis.

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

retrain the NN.

retrain $w^{[L]}$ a $b^{[L]}$ or if enough data, retrain whole model

→ have more radiology imgs

→ pretraining NN.

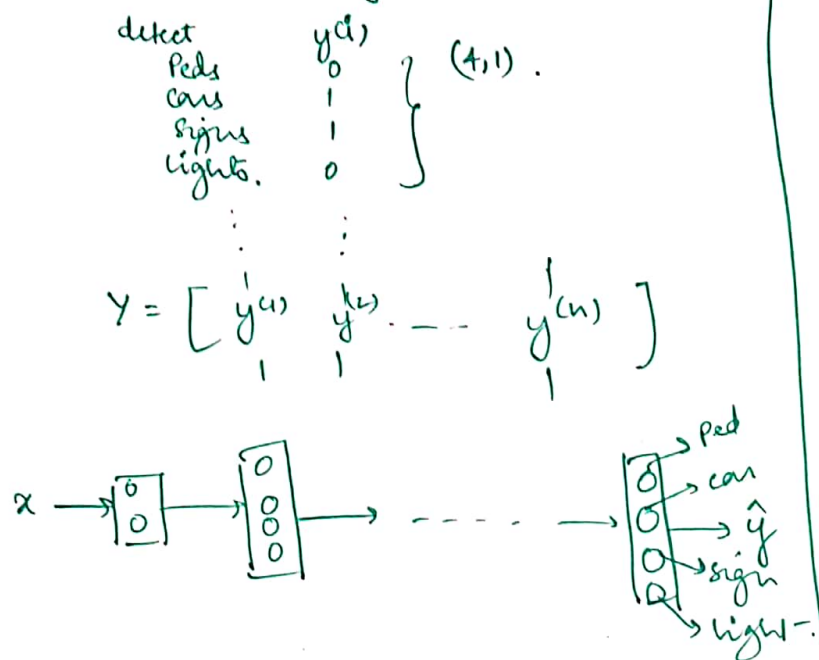can add several new layers instead of L-1 layer & retraining them.

### when Transfer learning makes sense?

Transfer A → B
- A & B have same i/p x.
- You have lot more data for A than B.
- low level features of B can be learned from A.

# * Multi-task learning:

ex: autonomous driving detector.

detect
Peds      $y^{(1)}$
Cars      0
Signs     1    } (4,1).
Lights    1
          0

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \cdots & y^{(n)} \\ 1 & 1 & & 1 \end{bmatrix}$$



→ Ped
→ car
→ $\hat{y}$
→ sign
→ light.

Loss  $\hat{y}^{(i)}_{(4,1)}$  $\frac{1}{m}\sum_{i=1}^{m}\sum_{j=1}^{y} L(\hat{y}^{(i)}, y^{(i)})$

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

Sum only when
val = 0/1

unlike softmax regr.
one img can have
multiple label.
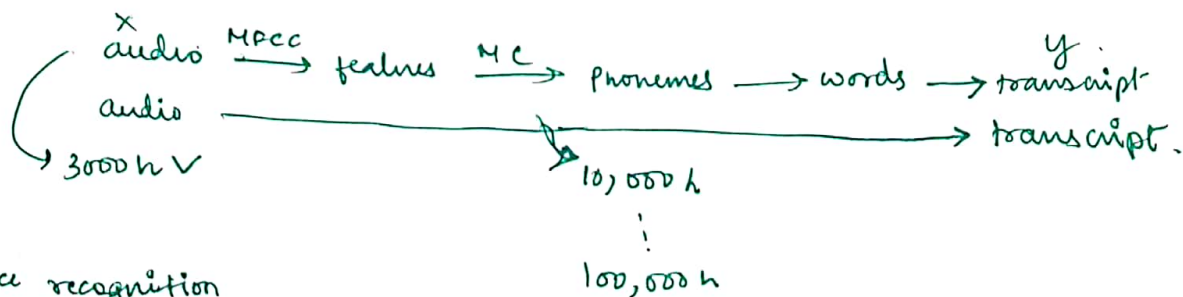
→ Multitask learning.
as 4 things are
done by 1 NN.

## When makes sense?

· Training for shared low level features.
· Usually: Amt of data for each task is quite similar.
° Can train a big enough NN to do well on all tasks.

→ if NN isn't big enough, then multitask learning can hurt model performance.
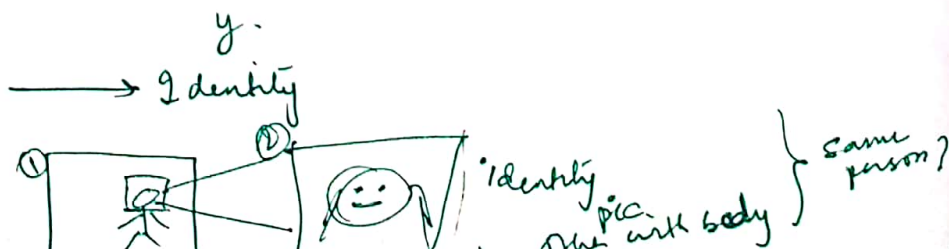
## ⊕ End to end deep learning:

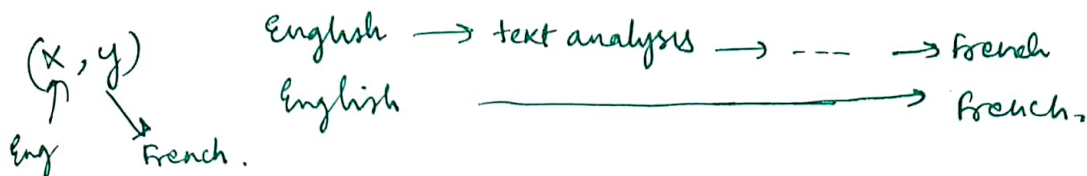Takes multiple stages and replaces with a NN.

· Speech recog example

$\overset{x}{\text{audio}} \xrightarrow{MFCC} \text{features} \xrightarrow{ML} \text{Phonemes} \longrightarrow \text{words} \longrightarrow \overset{y}{\text{transcript}}$

audio $\longrightarrow$ transcript.

→ 3000 h ✓

10,000 h
⋮
100,000 h

· Face recognition
Imag (x)



$\overset{y}{\longrightarrow} $ Identity

"Identity"
one pic with body  } same person?

∵ lot of data for 2 subtasks.

More examples.
- Machine translation

$(x, y)$

English → text analysis → --- → french
English ───────────────────→ french.

Eng     French.

Estimating child'sage : Image ①→ bones ②→ age.
                        Image ──────────→ age.

\* whether to use end to end deep learning?

● Pros & cons :

Pros :    $x \longrightarrow y$

• if big enough NN trained, it will figure it out.
• captures statistics in data.
Ex: phonemes   क, का , c  for cat. can be omitted

○ less hand-design of comp. needed


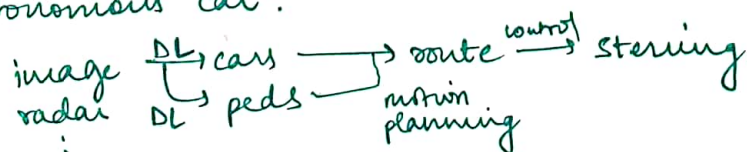Cons :
• need a lot of data for $x \rightarrow y$      $(x, y)$
                        i/pend o/pend.
• Excludes useful hand - designed components
                        By humans,   can be useful or harmful.

Ex: autonomous car :

image  DL→ cars ──→ route  control→ steering
radar  DL→ peds ──→ motion
                   planning
 ⋮
• use DL to learn individual comp.
• Carefully chose $x \rightarrow y$ depending what tasks you chose data for.