

1. Tuning process.

① α - learning rate
 β momentum size.

$$\beta_1 = 0.9 \\ \beta_2 = 0.999 \quad \epsilon \rightarrow 10^{-8}.$$

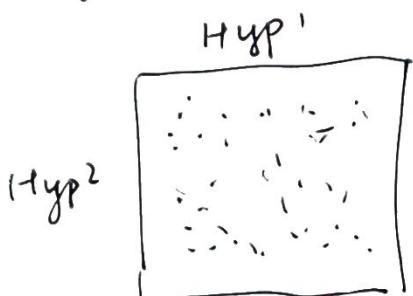
layers # hidden units:

③ learning-rate decay

[mini-batch size]

② second important

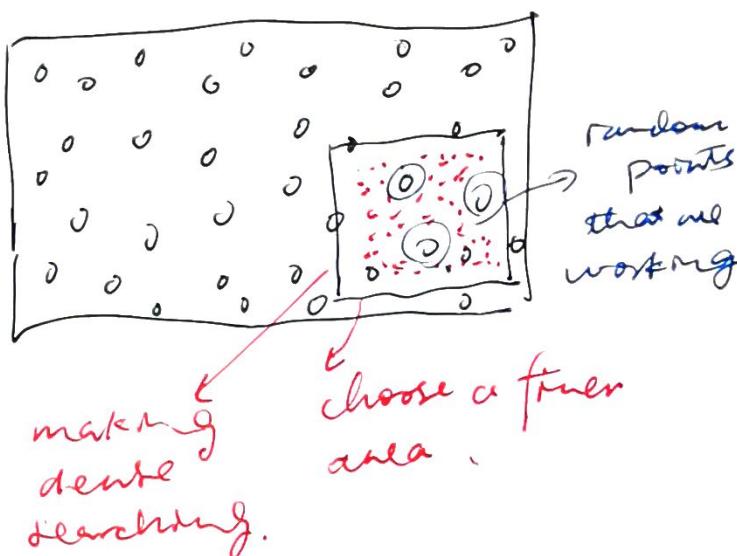
Try random values.



choose randomly ✓
say 25 points
try out 25 outcomes
of a single possible
hyperparameter reason

If use grid :: only
5 points of a grid the
hyp are tested out.

Coarse to fine.

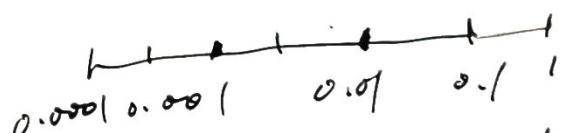


Using an appropriate scale to pick hyperparameters

$n^{[1]} = 50, \dots, 100$ $\frac{x}{50} \xrightarrow{100}$
layer $2 \sim 4$
uniformly at random over range of contemplating.

$$\alpha = 0.0001 \dots 1$$

logarithmic distribution



$\delta = -4 * np.random.rand()$
 $\alpha = 10^\delta \quad 10^{-4} \rightarrow 10^0$
 $\log_{10} 0.001 = -4$

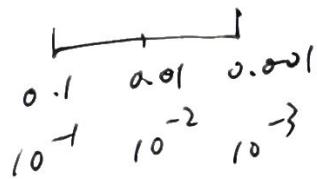
$$c = 0.$$

2.

momentum variable

$$\beta = 0.9 \sim 0.999$$

$$1 - \beta = 0.1 \sim 0.001$$



$$\gamma \in [-1, -3].$$

$$1 - \beta = 10^\gamma$$

$$\beta = 1 - 10^\gamma$$

Closer β is getting to 1.

It is more sensitive.

$$\beta : 0.9 \rightarrow 0.999 \sim 10 \text{ examples}$$

no big deal

$$\beta : 0.999 \rightarrow 0.9995$$

$\sim 1000 \text{ examples} \quad \sim 2000 \text{ examples}$
examples (averaging) over

Re-test hyperparameters occasionally.

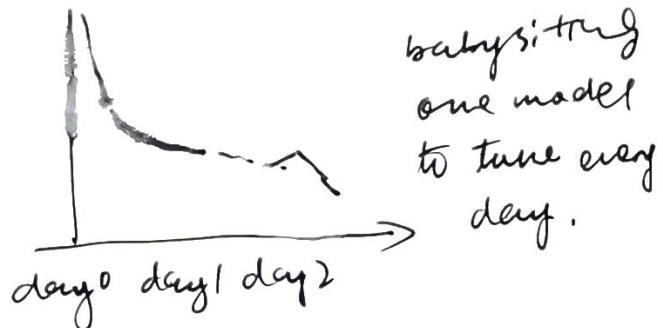
- NLP. vision.

Intuitively do get stable.

Re-evaluate occasionally.

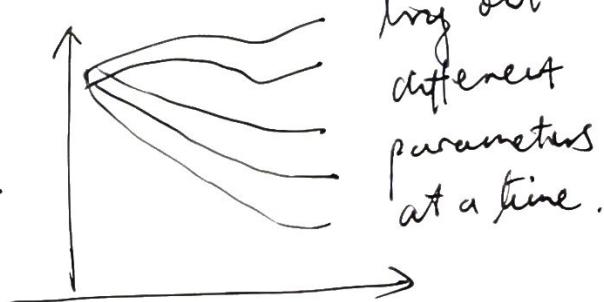
↑ ↓
new data. retest at least
new prob.... once several
 weeks.

Babysitting model



Huge data but no computation resources.

Train models in parallel



{ panda model
caviler model

? knowledge to
→ choosing computational resources.

Typically GPU? What capacity?

To handle what model?
How does it scale?
Does my laptop work?

3. Batch Normalisation

Sergey Tsofte & Christian Szegedy

Robustness.

Instead of only do norm
on inputs X . \rightarrow Try to
normalised middle layer
units $a \& (\bar{z})$.

Given some intermediate values
in NN $z^{(1)} z^{(2)} \dots z^{(m)}$ in L layers.

$$\mu = \frac{1}{m} \sum_i z^{(i)} \downarrow$$

$$\bar{J}^2 = \frac{1}{m} \sum_i (J_i - \mu)^2$$

$$\hat{\theta}_{\text{norm}}^{(i)} = \frac{\hat{\theta}^{(i)} - \mu}{\sqrt{\hat{\theta}^2 + \sigma^2}}$$

$$\sum_{m=1}^{(i)} \sim \sim (0.1)$$

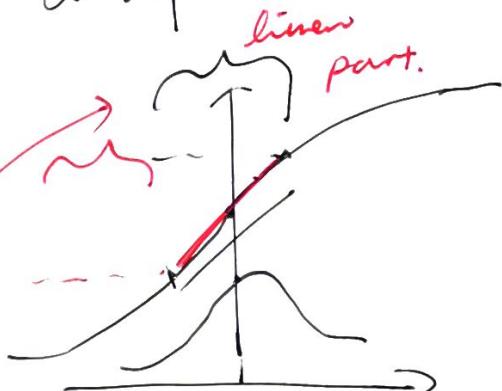
But sometimes do not want
to norm to always be ↗

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

Get updated
parameters
from w.b.

If $y = \sqrt{\alpha^2 + \beta}$ then $\frac{dy}{dx} = \frac{\alpha \cdot 2x}{\sqrt{\alpha^2 + \beta}}$

Use $\tilde{z}^{(i)}$ for later
computation of
consequent larger.
linear



might want functions
has a larger variance
& different means to
maybe use more
non-linear feature of
logistic function.

can be able to
control some means
& verbanues

Adding Batch Norms
to a network.

$$\begin{array}{c} x_1 \rightarrow \emptyset \rightarrow \emptyset \rightarrow \emptyset \rightarrow \hat{x} \\ x_2 \\ x_3 \rightarrow \emptyset \rightarrow \emptyset \rightarrow \emptyset \\ \hat{x} \rightarrow \tilde{x} \end{array}$$

$$4. \quad X \xrightarrow{w^{[l]}, b^{[l]}} z \xrightarrow{\text{BN}} \tilde{z} \xrightarrow{\beta^{[l]}, \gamma^{[l]}} \tilde{z}^{[l]}$$

Batch
Norm
(BN)

$$\rightarrow a = g(\tilde{z}^{[l]})$$

happens between computing \tilde{z} & a .

parameters $\left\{ w^{[l]}, b^{[l]}, \gamma^{[l]}, \beta^{[l]}, \mu^{[l]}, \sigma^{[l]} \right\}$

$$\frac{d\beta^{[l]}}{d\beta^{[l]}} \rightarrow \beta := \beta - \alpha d\beta$$

Tensorflow:

`tf.nn.batch_normalization`

Often applied with mini-batch

$$X \xrightarrow{\{w^{[l]}\}} z \xrightarrow{\{\beta^{[l]}, \gamma^{[l]}\}} \tilde{z} \rightarrow \dots$$

$$X \xrightarrow{\{b^{[l]}\}} z \xrightarrow{\{\mu^{[l]}, \sigma^{[l]}\}} \tilde{z}^{[l]}$$

$\tilde{z}^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$

\downarrow normalized.

(i) $\tilde{z}^{[l]}$ rescale by $\tilde{z}^{[l]}$

$\tilde{z}_{\text{norm}} \rightarrow \tilde{z}^{[l]}$

$w^{[l]}, b^{[l]}, \gamma^{[l]}, \beta^{[l]}$

cancel out the $b^{[l]}$. adding any constant do not change anything.

$$\tilde{z}^{[l]} = w^{[l]} a^{[l-1]}$$

$\tilde{z}^{[l]}$ *divide variance* \rightarrow *divide the mean*

$$\tilde{z}^{[l]} = \mu^{[l]} \tilde{z}_{\text{norm}} + \beta^{[l]}$$

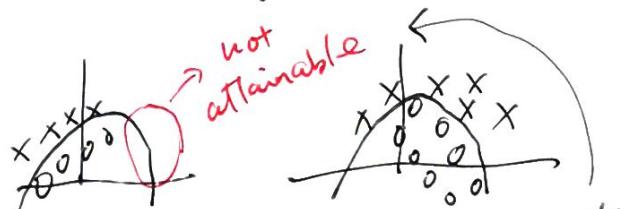
$$(\mu^{[l]}, \sigma^{[l]}) \quad (\mu^{[l]}, \sigma^{[l]})$$

just as $b^{[l]}$

Reason for BN to work.

- ① *using norm at hidden layers, instead only X.*
- ② makes weights in deep layer more robust to *st changing weight in earlier layers.*

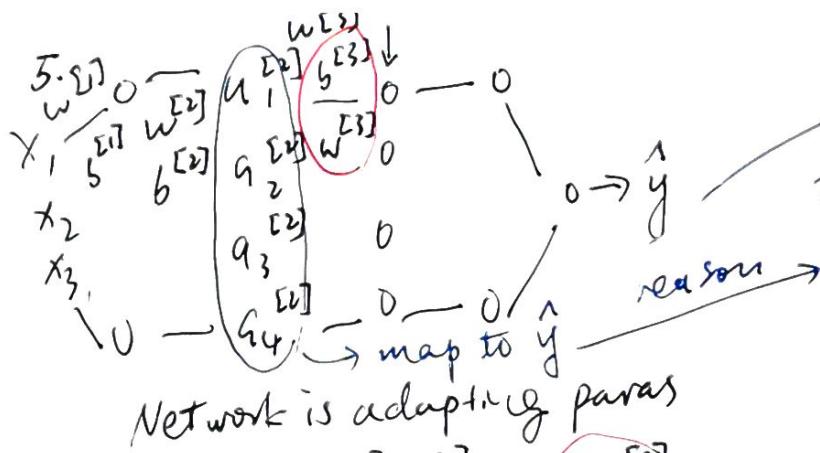
e.g.: train black cats pic & apply to colored cats.



train → could not get right part.

covariate shift TRUE func DID NOT CHANGE.

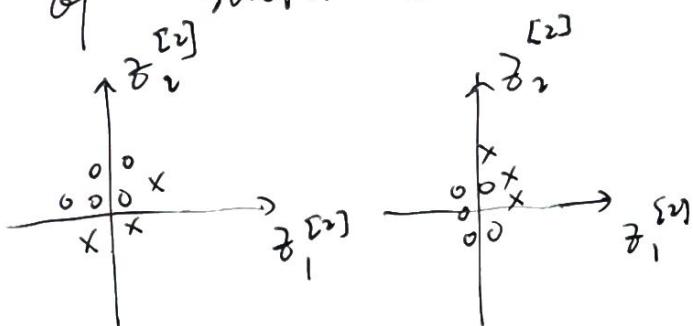
$X \rightarrow Y$
mapping trans changes
when X distri change



From perspective of third hidden layer:
changing all the time.
(like the change of input)
suffer from covariate shift.

$w^{[2]}, b^{[2]}, w^{[3]}, b^{[3]}$ → $\alpha^{[2]}$
change.

BN → Reduces the amount
of shifts around.



kind of control the
mean & variance of the
changing paras.

$$\beta^{[2]}, \gamma^{[2]} \\ \text{governed by}$$

makes the current layers
more stable. The
consequent layers
could stand more reliably

Each mini-batch
scaled by mean/var
computed on just the
same mini-batch

→ adds noise → rep.

similar to dropout
adds noise to each
hidden layer. small

(operating on a small
scale of data)

→ Regularization
effect
(side effect).

larger mini-batch size
→ reduce noise
reduce Regularization effect

At test time, μ, σ^2 as some sort of mean may
not be calculated on with many examples
calculate example once a time
→ EWA.

6. estimates - μ , σ^2 using weighted average (across mini-batches.)

$$x \cdot \underbrace{x}_{\mu_{\text{est}}}, \underbrace{x}_{\mu_{\text{est}}}, \dots \rightarrow \mu_{\text{est}}$$

Test time, process sample example at a time.

In practice, use EWA to keep track of μ & σ^2 . (running average)

$$\rightarrow \tilde{\mu}_t = \beta \tilde{\mu}_{t-1} + (1-\beta) \mu_t.$$

For multi-class classification

Softmax Regression
c = # classes = 4.

$$x \rightarrow \boxed{\quad} \rightarrow \boxed{\begin{matrix} 0 \\ 0 \\ 0 \\ 0 \end{matrix}} \rightarrow \hat{y} \rightarrow p(\text{other}|x)$$

on layer L:

$$z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]}$$

(4.1).

Activation Func. softmax

$$t = e^{(z^{[L]})} \quad (4.1)$$

$$a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{j=1}^4 t_j} \quad (4.1)$$

normalised.

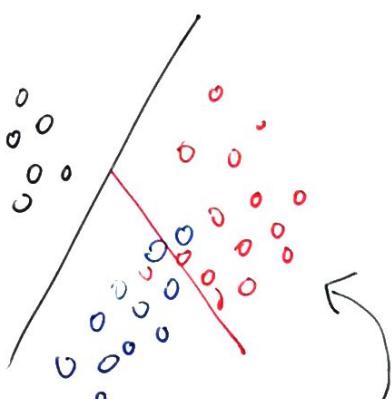
$$a_i = t_i / \sum_{i=1}^4 t_i$$

$$g: z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}, t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = [1]$$

$$\sum_{j=1}^4 t_j = 176.3, \quad a^{[L]} = \frac{t}{176.3}$$

$$\sum_{j=1}^4 \left\{ \begin{array}{l} 0 \\ 0 \\ 0 \\ 0 \end{array} \right\} \rightarrow a = \frac{e^5}{176.3} = 0.842$$

$$a = \frac{e^2}{176.3} = 0.042 \dots$$



softmax single layer result (linear boundary)
would achieve non-linearity through deep training

Hardmax

$$a^{[l]} = \begin{bmatrix} 0.842 \\ 0.092 \\ 0.002 \\ 0.114 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$z^{[l]} \rightarrow a^{[l]} = \hat{y} \rightarrow \delta(\hat{y}, y).$$

back prop:

$$\frac{\partial \delta^{[l]}}{\partial z^{[l]}} = \hat{y} - y.$$

(4.1) start at layer L to do backprop

Tensorflow helps with backprop. only do forward right is enough.

Loss function (entropy loss)

$$\hat{L}(y, \hat{y}) = -\sum_{j=1}^4 y_j \log \hat{y}_j$$

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \quad \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$L(\hat{y}, y) = -\log \hat{y}_2$$

→ make L ↓. $\hat{y}_2 \uparrow$

corresponding prob ↑
as high as possible

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

$$Y = [y^{[1]}, y^{[2]}, \dots, y^{[m]}]$$

$$= \begin{bmatrix} 0 & 0 & \dots \\ 0 & 1 & \dots \\ 0 & 0 & \dots \end{bmatrix} \rightarrow (4, m)$$

$$\hat{Y} = \begin{pmatrix} 0.2 & 0.5 & \dots \\ 0.3 & \vdots & \vdots \\ 0.6 & \vdots & \vdots \\ 0.4 & \vdots & \vdots \end{pmatrix} \rightarrow (4, m)$$

?

Carry out error analysis.

90% accurate

Sometimes look like
a cat 
but it is a dog.

Start cat classifier so
better or deg? randomly

Error analysis procedure:

1. Get ~100 mislabelled dev set examples.

2. count up how many are dogs.

$$5\% \rightarrow \frac{5}{100} \rightarrow 90\% \rightarrow 90.5\%$$

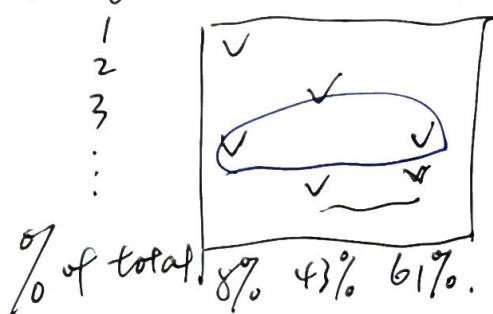
{ ceiling of improvement of performance.

$$50\% \rightarrow 90\% \rightarrow 95\%.$$

work several hypotheses in parallel.

{ dogs mislabelled by cat misrecg. blurry images.

Image: Dog DC Blurry cominit



Cleaning up incorrected 1 labelled data.
(mislabelled by human etc.)

DL alg are quite robust to random errors in training set.

But systematic errors ↗

ADD extra column for in dev/test set.

startling.

10%.

{ Errors due to mislabelled errors → 0.6%.

other causes → 9.4%.

6% significant.

2% { 1.4%.

↓

Instagram

a sense of overall improvement ceiling for each feature.

1.

2 A & B.

2.1% 1.9%. hard to judge the (not trust) then right one.

→ mislabeled data

→ fix mislabeled.

Correcting incorrect.

dev/test set examples.

① apply same process.

② to ensure same distribution

③ consider examples

trigly classified → not actually right.

→ otherwise → bias.

④ Train & dev/test data slightly may vary differently in term of distribution

Build new system suggestion:

build first system quickly

& start iterating

set up dev/test set

& metric

→ build first system quickly. & dirty

use variance/bias analysis & error analysis?

↓
figure out what direction to go is most worthwhile.

toss of directions

Train & dev/test on different distributions

Data from webpage	Data from mobile app.
clear. high resol.	blurry. low resol.
$\approx 200,000$	$n \approx 10,000$

and about performance of pred from mobile app.

OP1

↑ not good for this
combine them randomly shuffle.

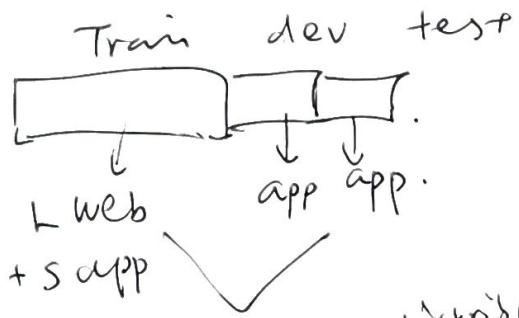
OP2

train	dev	test
200,000	5,000	app app
web.	mobile	

better for long term.

2.

Data mismatch problem.



not same distribution.

train error 1% → no longer
dev error 10% → concluded
+ e... as variance.
directly.

Since, training data
could be easy to train
& does not apply well on
another distribution.

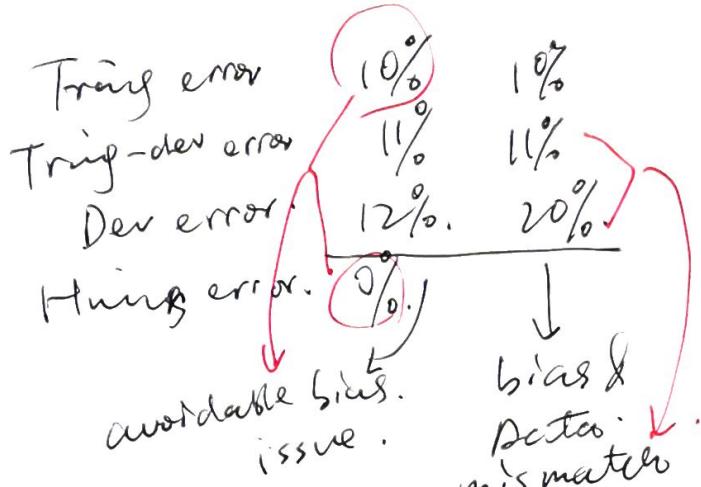
↓ solutions.

+ train/dev set

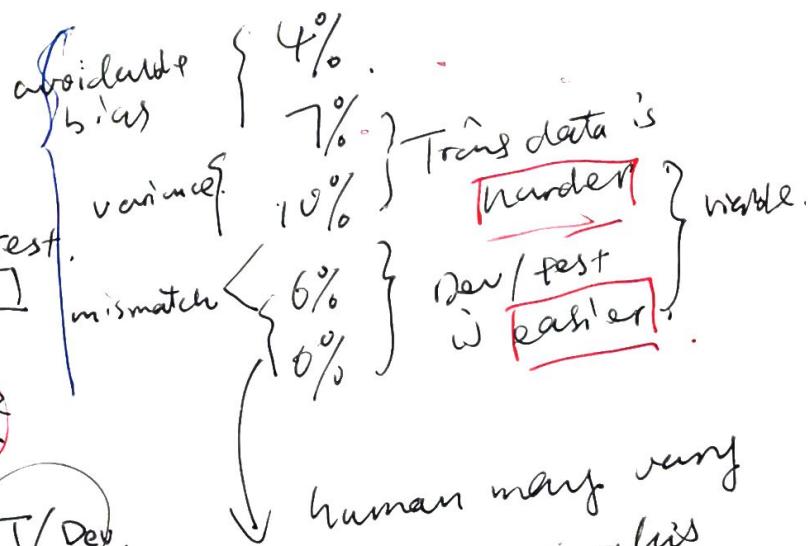


shuffle
Train &
Set data Same
and leave one block for T/Dev.
Same distri

Train error 1% 1% not training.
Train/Dev error 9% 1.5% } data
Dev error 10% 10% } mismatch
problem.
Variance prob



Human level error } avoidable
Train set error } bias
T-dev set error } variance.
Dev error } data
overfitting Test error mismatch



e.g.: 4% → { 6%, 2%

3.

Addressing Data mismatch problem.

- ④ Carry out manual error analysis to understand difference between training & dev/test sets.

Eg: noisy \rightarrow car noise.
for dev set.

several numbers.

insights:

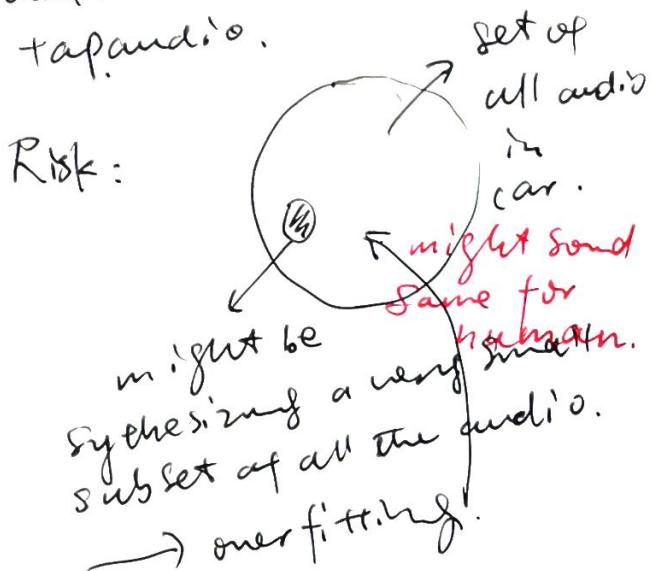
- ② Making training data more similar to Dev/test. or collecting more data. to dev/test sets.

Artificial Data Synthesis.

$$\text{eg: } \text{The quick brown fox jumps over the lazy dog} + \text{car noise} = \text{Synthesized}$$

"The quick brown fox jumps over the lazy dog" \rightarrow easy sentence with all alphabet.
car noise \rightarrow synthesized in-car audio.

replicates small subset of car audio 1000 times & add it with the normal car audio.

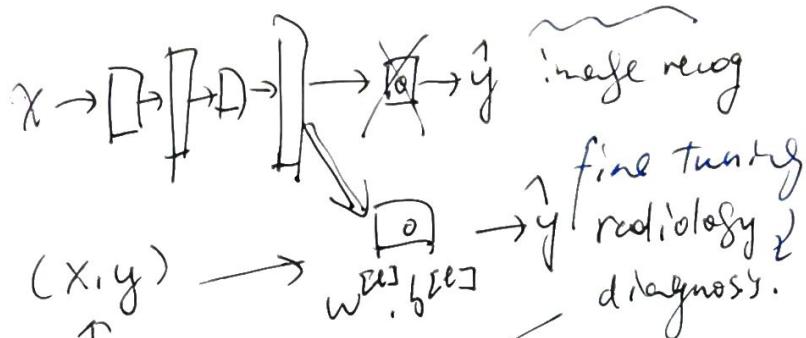


might overfit to the small amount of Train sub set. while dev/test \rightarrow very general cases.

Using knowledge trained by a task to apply to a second task — Secure transfer learning.

Transfer learning.

pre-training



new inputs. last few layers.
small set of data viable
to only train last layer.

large set of data. train all the layers of data.

or the layers targeted
to do the recog.

points \rightarrow edge \rightarrow shape \rightarrow facial
(curve)

may. \downarrow diagnoses

create several new

layers to tackle your
new problem.

should be
inter-connected
somehow.

Makes sense when:

a lot of data transferred
from.

limited data on
Targeted (transfer to)

not viable vice-versa. not needed but

less valuable.
& less helpful.

One image can have multiple
labels.
a car can have 4 feats.
a thing only be a car or dog
or cat..

5.

Multi-task learning.

self-driving cars.
detecting pedestrians.
other cars, stop-signs.
traffic lights ...

$$X^{(i)} \rightarrow \begin{pmatrix} y_1^{(i)} \\ y_2^{(i)} \\ \vdots \\ y_m^{(i)} \end{pmatrix} \in \mathbb{R}^m$$

$$Y = \begin{bmatrix} y_1^{(1)} & y_2^{(2)} & \dots & y_m^{(m)} \end{bmatrix} \quad 4 \times m.$$

$$\square \rightarrow \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \rightarrow \hat{y} \in \mathbb{R}^4$$

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

\rightarrow usual logistic loss.

multiple softmax
regression.

One image can have multiple
labels.

a car can have 4 feats.
a thing only be a car or dog
or cat..

Training one neural network to do four things might be better than training 4 network to do 4 things separately.

→ earlier feature in an can be shared.
→ multi-task learning advantage.

$$Y = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ ? & ? \end{bmatrix}$$

vable for some labels
missing for certain examples.

Sum over 1-4
vs sum over 1-2
over where labels present

End to end deep learning.

Eg: speech recog example.

X audio $\xrightarrow{\text{MFCC}}$ features $\xrightarrow{\text{ML}}$ phonemes "cat"

\rightarrow words \rightarrow transcript.

audio $\xrightarrow{\text{whole nn.}}$ y. \rightarrow transcript.

makes sense:

(1) Amount of data for each task is quite similar

A₁ 1000 } data.
A₂ 1000 } amount.
⋮
A 100-1000.

(2) training on a set of tasks that could benefit from having shared lower-level features.

(3) Training a big enough nn to do well on all tasks

\downarrow
hurt performance
only when nn. not big enough,

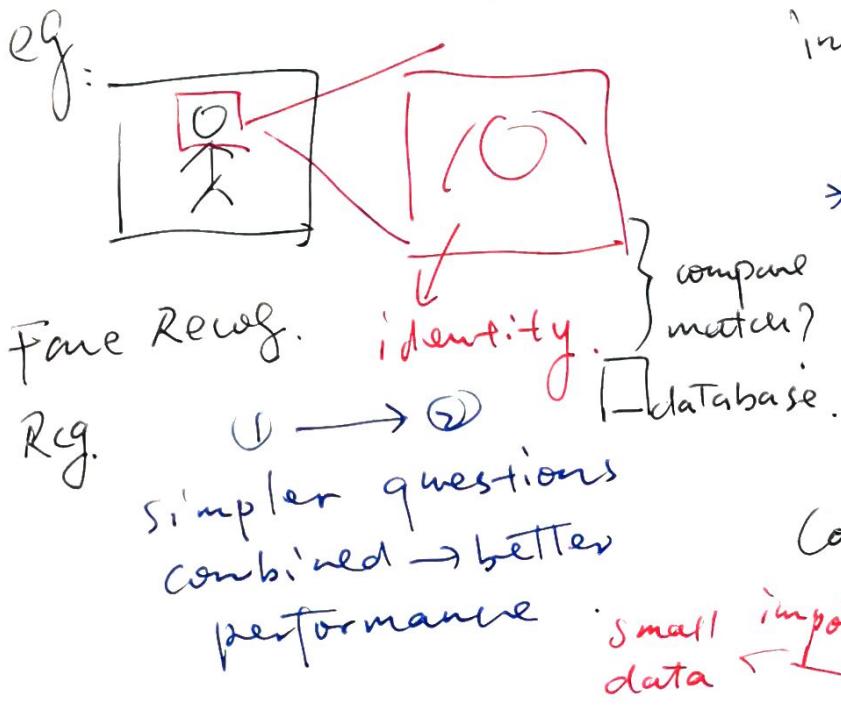
bypassing all intermediate steps. with one neural n.

traditional pipeline

End to End pipeline.

Large amount of data
→ End to End approach
better.

Small, medium size
→ traditional or
Small (partially)
end to end approach



A lot of data for
doing both ① & ②.
Not enough for directly
end to end.

But for Machine Translation
Eng → text... → ... → French.
Eng → → French.

End to End does
better in
machine trans.

Big data on (X, y)
Eng → French.

Data feed in the most
important, productive
element.

* Pros:
① Let the data speak, not enforce human intuition
② less hand-designing of components needed.

Cons:
① Large data needed.
② Excludes potentially useful hand-designed components.

Combined with other
subject - motion planning
control more promising

Core question: how to evaluate
amount of data needed
to map from $X \rightarrow Y$.
unsolved. 7.