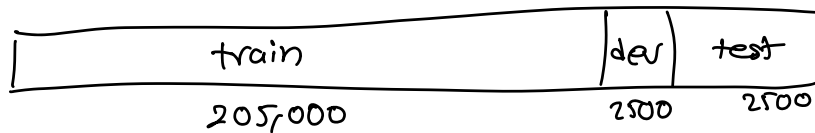## Cat App example:-

training data
- data from mobile app (blurry, low res) $\approx$ 10,000
- data from web (high quality, high res) $\approx$ 200,000

test data: data from mobile app.

## Will data from web help: ?

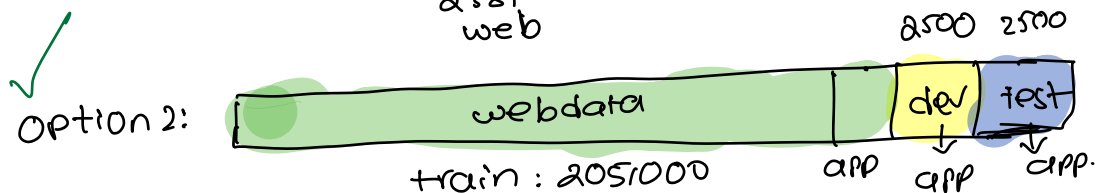✗ Option 1: combine $m = 210,000$, random shuffle (bad dev set)

| train | | dev | test |
|---|---|---|---|
| 205,000 | | 2500 | 2500 |

adv: train, dev, test have same distribution

disadv: dev set major data will come from web page data.

dev $m = 2500$

2381 web    119 mobile

✓ Option 2:

| webdata | | dev | test |
|---|---|---|---|
| train: 205,000 | | 2500 app | 2500 app. |
| | | app | |

train set: $200,000 + 5000$
            web      app

dev: 2500 app

test: 2500 app.

adv: dev set is accurate to test set

dis: training dist is different from dev, test set

## speech recognition example!

— speech activated rear view mirror in car.

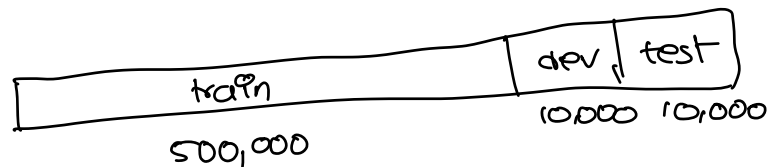**Training:**

purchased data $(X, y)$
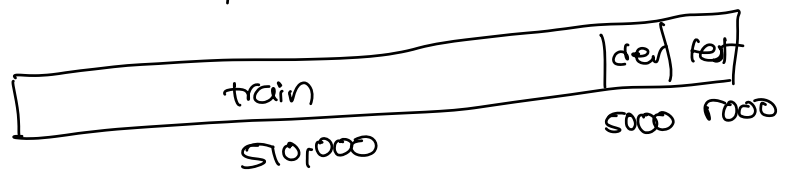
smart speaker control

voice keyboard

500,000

**Dev|test**

speech activated
rearview mirror

20,000



train 500,000

dev | test
10,000  10,000



train 510,000

dev | test
5000  5000

train: 500,000 from others
10,000 from rearview
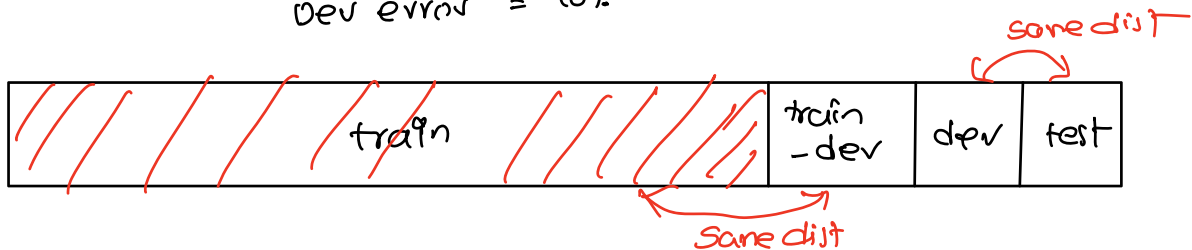
dev:  5000 from rearview
test:  5000 from rear view.

---

## rec2: Bais & variance with Mismatched Data distribution:-

Assume human error on cat classification = 0%

Training error = 1%
Dev error = 10%



train | train -dev | dev | test

same dist (between dev and test)

same dist (between train and train-dev)

if Train dist = Dev dist
   problem is variance, more data, $\lambda 2$ reg, dropout.

if Train dist ≠ Dev dist

→ maybe it is doing fine on dev set.
→ training set may be easy (high res, high quality)
→ dev set may be hard (low res, low quality)

| | | | | |
|---|---|---|---|---|
| Human error | 0% | 0% | 0% | 0% ⎤ avoidable bias |
| Training error = | 1% ⎤ | 1% | 10% ⎤ | 10% ⎦ variance |
| Training dev set = | 9% ⎦ | 1.5% ⎤ | 11% | 11% ⎤ data mismatch |
| Dev set = | 10% | 10% ⎦ | 12% | 20% ⎦ degree of overfit to dev set. (find bigger dev set) |
| Test set = | | | | x% ⎦ |

This is variance problem
→ not generalising well on training-dev set which comes from same train dist

low variance
→ data mis match problem.
→ get more train data dist to dev,test

→high bias

→high bias
→ data mismatch problem.
→ low variance

| | | |
|---|---|---|
| Human level error : | 4% ⎤ avoidable bias | 4% |
| Training set error : | 7% ⎦ variance | 7% |
| Training dev set error : | 10% ⎤ data mismatch | 10% |
| Dev set error : | 12% ⎦ degree to overfit dev set | 6% |
| Test set error : | 12% ⎦ | 6% |

training data is harder than dev, test set.

## More General Formulation:-

|  | General speech recognition | rear view mirror speech data |
|---|---|---|
| Human level | human level: 4% | human level: 6% |
| Error on examples NN trained on | Training error: 7% | Training error: 6% |
| Error on examples NN has not trained | Training-dev : 10% error | Dev/test : 6% error |

avoidable bias

variance

data mismatch

## Lec3: Data Mismatch Problem:-

· carry out manual error analysis to try understand difference between training and dev/test sets.

Eg:  1. rear view number - car noise in dev set

2. rear view number - misrecoginizing street numbers in dev set.

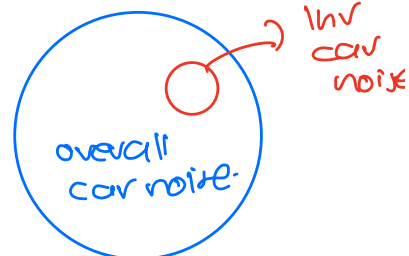Sol: collect more data similar to dev/test set.

→ add car noise to original avdio

→ add more numbers avdio in train set.

# Artifical data synthesis:-

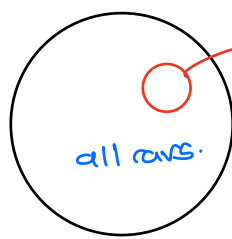normal avido + "car noire" = synthesized in car audio
10,000hr            1hr

   disadu : overfit data 1hr car noire

   10,000 hr unique car noire can
       be even better performance.

overall
car noire.

1hr
car
noire

## Car recognition:-

computer graphic generation of car image.

all cars.

synthesibed ≈ 20

model might overfit training
synthesized data.

Trained NN : $(x, y)$     $x \to$ image
$y \to$ cat, dog, bird

Transfer learning   $x \to$ image ( radiology diagnosis)

→ remove last layer weights & ouput radiology output.

→ retrain the last layer.

→ fine tuning.

```
┌──────────┐  ┌─────────────────┐
│          │  │  ┌───┐          │
│    NN    │ ⇒│  │ O │ → y      │
│          │  │  └───┘          │
└──────────┘  └─────────────────┘
  Pretrained.    new weights
                   $w^{[l]}, b^{[l]}$
```

① if you have small data, then retrain last layer weights.

② if you have more data, then retrain complete NN.

**Use case:-**

→ If you have alot of data, for the problem you are transfering
from $(m = 1,000,000)$

→ and less data, for the problem, you are transfering to.
        $(m \ni 100 \leq 1000)$

**Transfer from $A \to B$**

① Task A,B have same input $X$

② You have lot more data Task A than Task B.

③ low level features from A could be helpful for learning B.
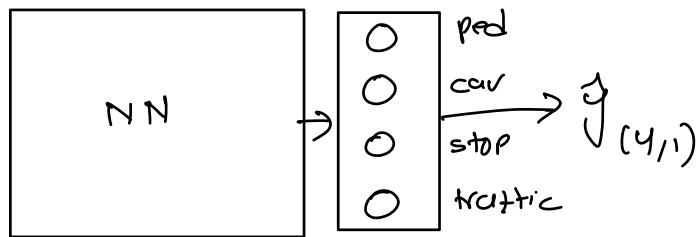
# Lec 5: Multi Task Learning:

Autonomous vehicle:-

detect:
- pedestrians
- other cars
- stop signs
- traffic lights

input   $x^{(i)}$ (image) $\longrightarrow$ $y^{(i)}$ [ped, car, stop, traffic lights] $(4,1)$

$$Y = \left[ y^{(1)} \; y^{(2)} \ldots y^{(m)} \right] \; (4, m)$$



| NN | |
|---|---|

ped
car
stop
traffic

$\hat{y} \; (4,1)$

Loss: given $\hat{y}^{(i)}_{(4,1)}$

$$Loss = \frac{1}{m} \sum_{i=1}^{m} \sum_{j=1}^{m} \mathcal{L} \left( \hat{y}^{(i)}_j, \hat{y}^{(i)}_j \right) \quad \text{(usual logistic loss)}$$

$$\mathcal{L} = - y^{(i)}_j \log \hat{y}^{(i)}_j - (1 - y^{(i)}_j) \log (1 - \hat{y}^{(i)}_j)$$

$\Rightarrow$ This is multi task learning.

unlike softmax regression:
one image has multiple labels.

$$Y = \begin{bmatrix} 1 & 1 & 0 & ? \\ 0 & 1 & ? & 1 \\ ? & ? & 1 & 2 \\ 2 & ? & 1 & ? \end{bmatrix} \rightarrow \text{Still works for multi task learning.}$$

(sum only over labelled class).

When Multi task learning Works:-

① training on set of tasks that could benefit from learning shared lower level features.

② Usually: amount data you have for each task is quite similar.

Transfer learning

A (1000,000)

↓

B (1,000)

Multi task learning

A1  1000
A2  1000
:
A100  1000
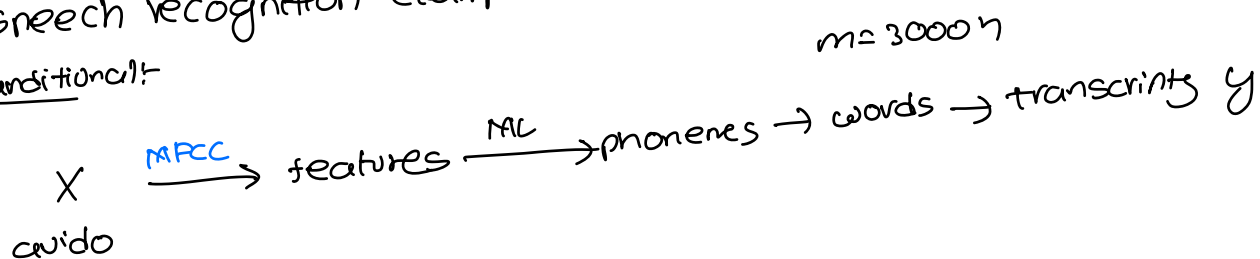
→ A100 now has 99,000 more data point which can be used for low level feature extraction.
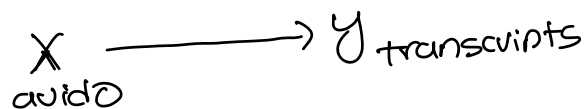
③ can train a big enough NN to do well on all tasks.

## Lec6: End to End Deep Learning:-

Speech recognition example:

traditional:-

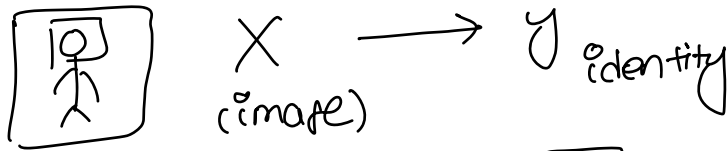$X$ $\xrightarrow{\text{MFCC}}$ features $\xrightarrow{\text{ML}}$ phonemes → words → transcripts $y$

audio                                                    m= 3000 h

End to End:-

$X$ ——————→ $y$ transcripts

audio

→ need a lot of data for end to end approach.

m = 10,000 h – 100,000 h

# Face recognition :-

 X $\longrightarrow$ Y identity
(image)

① detect face and crop it 🙂 ⎤
                              ⎬ works much better.
② 🙂 → identity.      ⎦
       ↘ 🙂 a very data base.

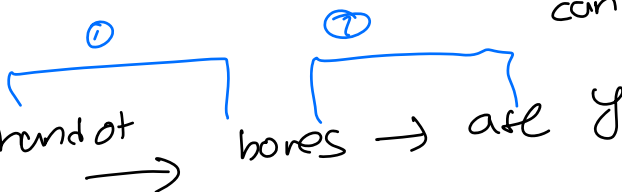→ 2 seperate easy problems.

→ Task A, Task B have individual more data.

## Machine translation :-

before: English → text analysis ----→ french

today: English (X) $\longrightarrow$ French (y) work very well.

## Estimate child's age :-

before                    ①        ②        can work better.
X ( xray image of hand of ⌐‾‾‾⌐  ⌐‾‾‾‾⌐
   child)         $\longrightarrow$ bores → age Y

today        X $\longrightarrow$ age (y) (doesn't work very well
                                          because of less data).

## lec7: Whether to use End to End Deep learning :-

**Pros:**
- lets the data speak (hopefully NN learns best $f(x) = y$)
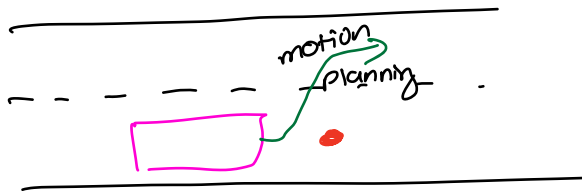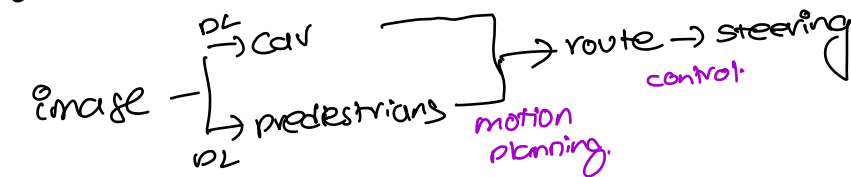- less hand-designing of components needed.

**cons:**
- May need large amount data which is hard to collect.
- Excludes pontentially usefully hand designed components. which can be very helpful and can be cheap.

**key question!**

Do you have sufficient data to learn a function of the complexity needed to map $x \to y$?

Autonomous car :-

image — [ DL → car,  DL → predestrians ] → route → steering control. motion planning.



- use DL to learn individual components
- carefully choose $x \to y$ depending on what you can data for.

image ⟶ steering ( ? ) (work in progress).