



deeplearning.ai

Mismatched training
and dev/test data

Training and testing
on different
distributions

Cat app example

Data from webpages



≈ 200,000 Images

Dist of training is diff than dist of test

① option 1: combine all Images (210,000)

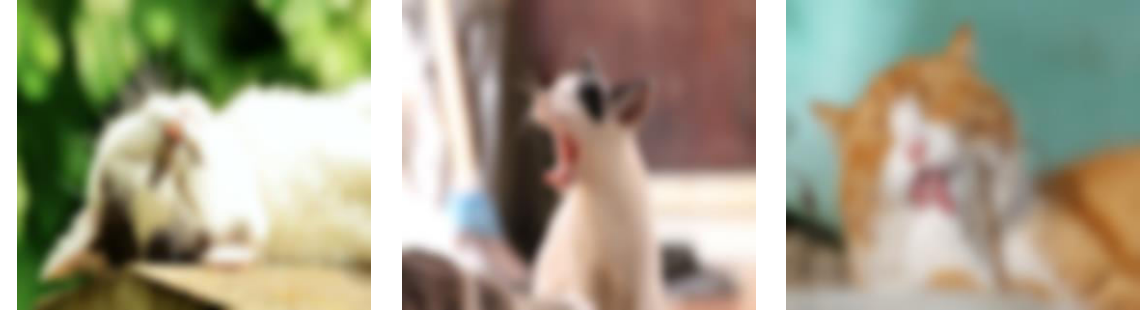


Disadv: dev/Test set has too many Images that are not of the "true" test set dist, ie, mobile uploads

out of 2.5K, $\frac{200K}{210K} \times 2.5K = 2381$ will be high quality non real test Images

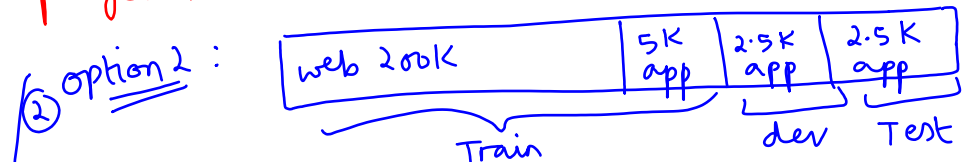
General Approach, just get data from Anywhere & shove it into training

Data from mobile app



≈ 10,000 Images

(You care about this, doesn't matter if you've trained on high quality Images, if you can't perform in real world, your App fails)



- Train = 200K web + 5K app = 205K
 - At least this way you're optimizing for the correct dev set Images
 - Disadv: Train dist is diff from dev/test
- Go w/ this split

Speech recognition example

Alexa located
in the rear
view mirror
of your car



Training

Purchased data

Smart speaker control
(Alexa)

Voice keyboard

- Say you collect 500K samples from
all training example sources

...

The problem is that none of the trained examples
will be specific to the context of the car →

Alexa, drive me to the nearest gas station
Alexa, open the trunk of the car
⇒ train is diff from dev/test

Dev/test

Speech activated
rearview mirror

20K Samples

∴ the split could
look like

