



deeplearning.ai

Mismatched training and dev/test data

Addressing data mismatch

*You've done your "Error Analysis" & see
most of it is because of data mismatch
ie, $\text{dist}(\text{Train}) \neq \text{dist}(\text{dev/Test})$
So what can we do?*

Addressing data mismatch

- Carry out manual error analysis to try to understand difference between training and dev/test sets

eg car Rearview mirror Assistant
= one way data differs → car samples are noisy [car noises in speech samples]
compared to training data
↳ traffic
↳ inside/Engine noise
↳ Honk

- Make training data more similar; or collect more data similar to dev/test sets & put this in training set

Artificial data synthesis



+



=



“The quick brown
fox jumps
over the lazy dog.”

Clean Audio

Car noise

Synthesized
in-car audio

Caveat

↳ You have 10 K hours of speech (clean, noiseless speech)
↳ You have 1 hour of car noise
⇒ You generate all synthesised car Audio Speech using
this 1 hour car noise, Repeated over 10 K hours
⇒ overfitting on that 1 hour of car noise

Artificial data synthesis

Car recognition: *For Self Driving car, You need to Recognize nearby cars*



- Synthesize Car Images instead of taking 1M pictures of cars on the street, why not?

