

Cat app example

Data from webpages



care about this

Data from mobile app

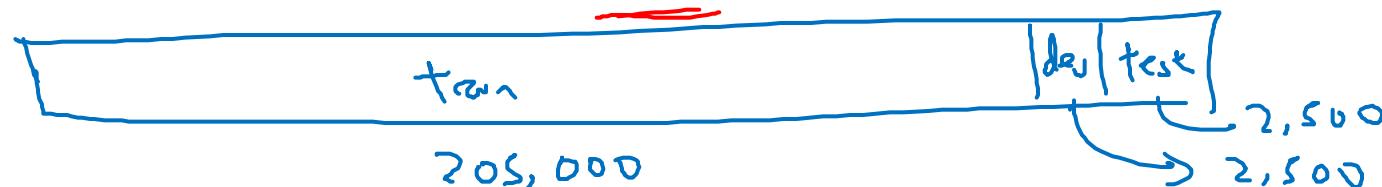


$\rightarrow \approx 200,000$

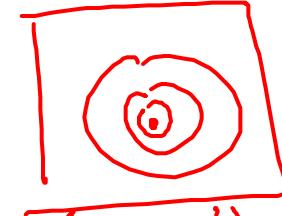
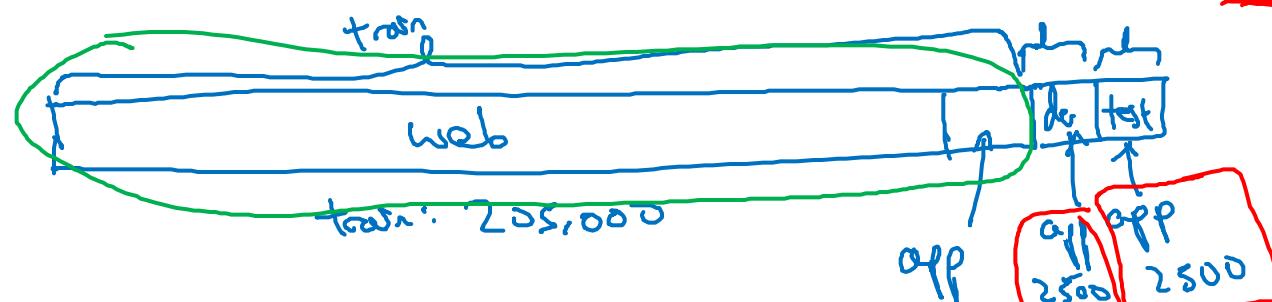
$210,000$
shuffle

$\rightarrow \approx 10,000$

X Option 1:

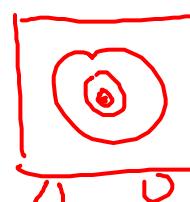


Option 2:



$\frac{200K}{210K}$

2381 - web
119 - mobile app



Speech recognition example

Speech activated rearview mirror



Training

Purchased data $\downarrow \downarrow$
 x, y

Smart speaker control

Voice keyboard

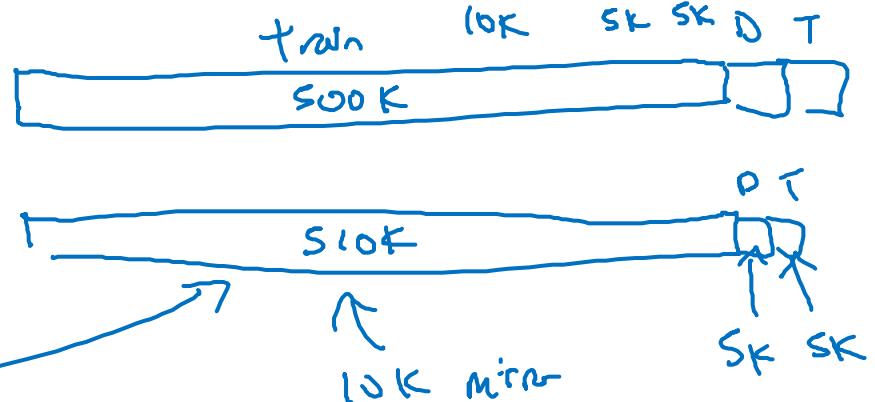
...
500,000

utterances

Dev/test

Speech activated
rearview mirror

$\Rightarrow 20,000$





deeplearning.ai

Mismatched training
and dev/test data

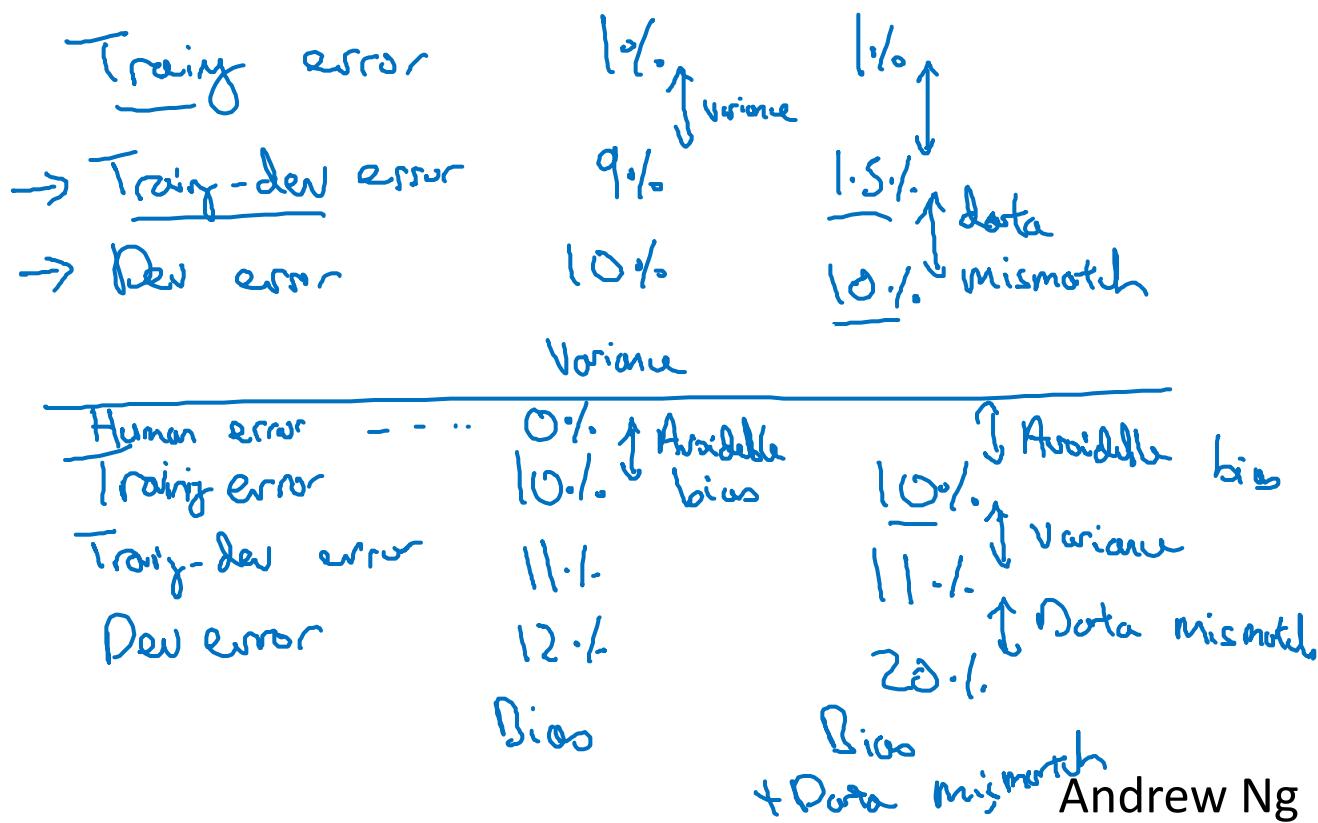
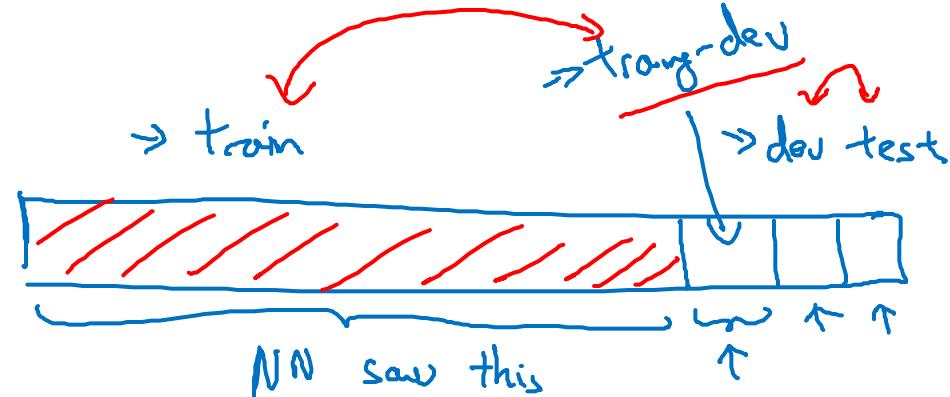
Bias and Variance with
mismatched data
distributions

Cat classifier example

Assume humans get $\approx 0\%$ error.

Training error 1% \downarrow 9%
Dev error 10% \uparrow

Training-dev set: Same distribution as training set, but not used for training



Bias/variance on mismatched training and dev/test sets

Human level

Training set error

Training - dev set error

→ Dev error

→ Test error

4% ↑ avoidable bias

7% ↑ variance

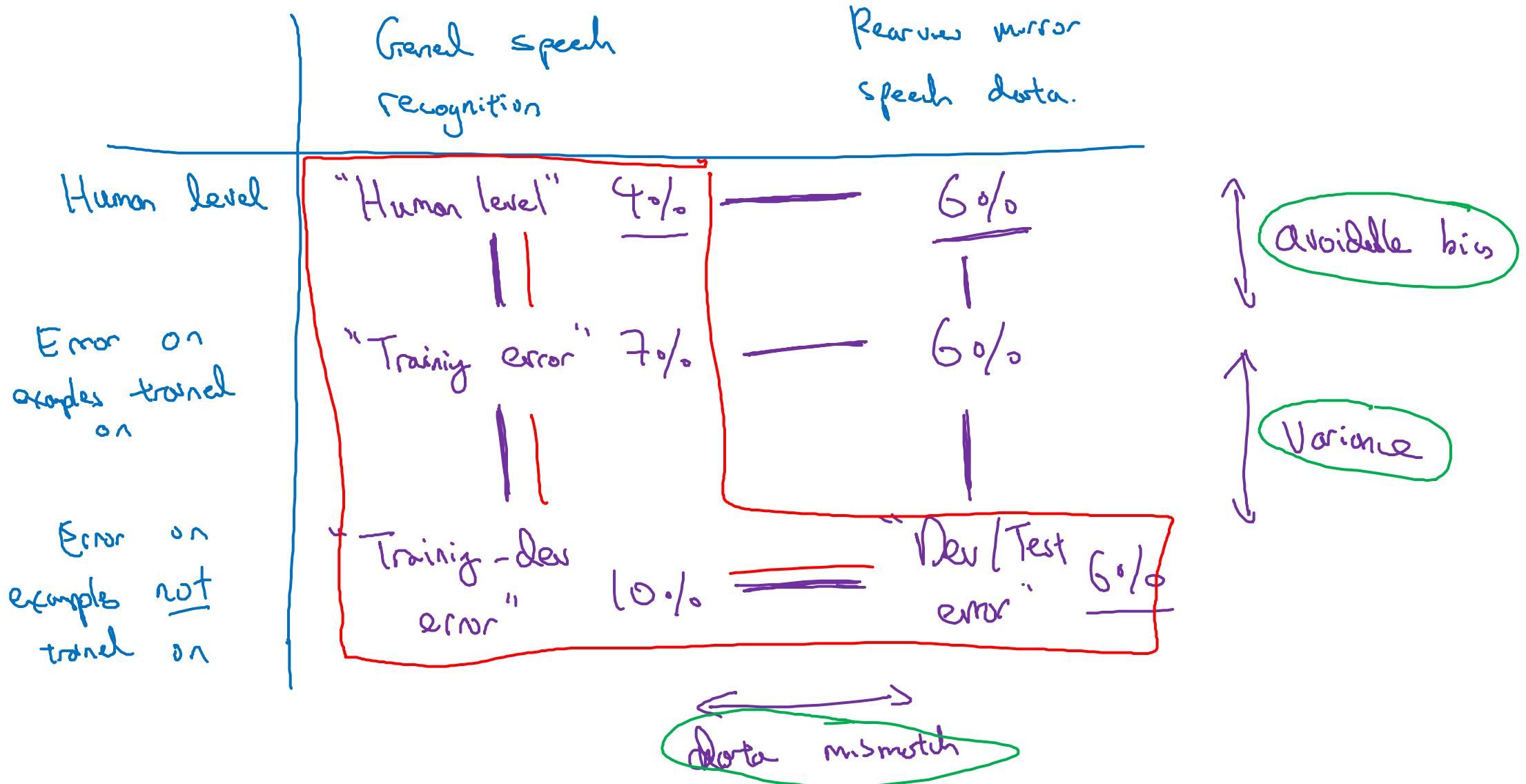
10% ↓ data mismatch

12% ↓ degree of overfitting
to dev set.

4%

7% }
10% }
6% }
6% }

More general formulation





deeplearning.ai

Mismatched training
and dev/test data

Addressing data
mismatch

Addressing data mismatch

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

E.g. noisy - car noise

street numbers

- • Make training data more similar; or collect more data similar to dev/test sets

E.g. Simulate noisy in-car data

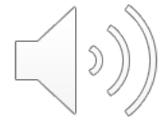
Artificial data synthesis



+

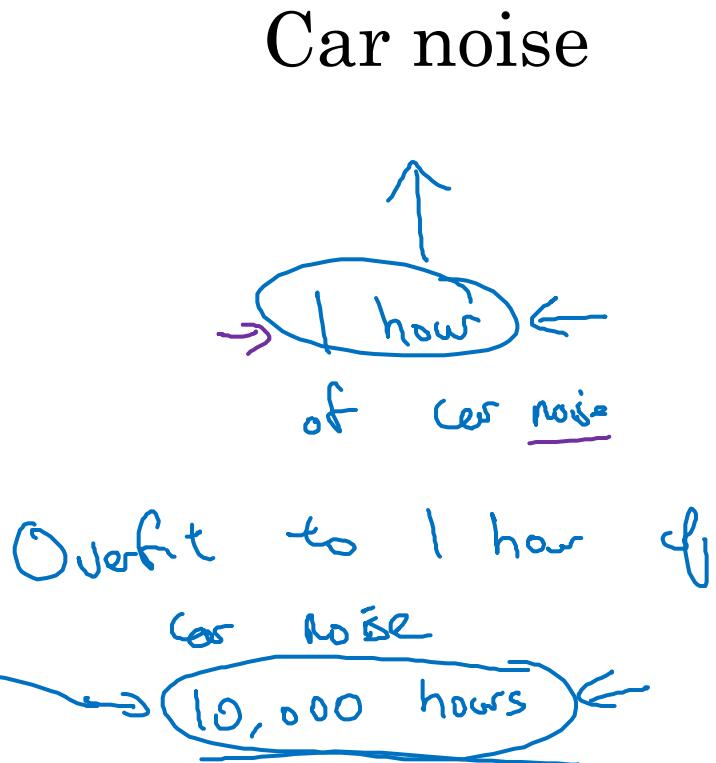


=

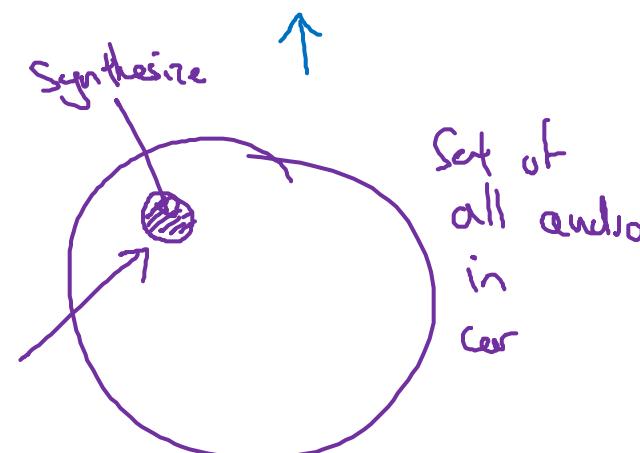


“The quick brown  fox jumps over the lazy dog.”

10,000 hours



Synthesized
in-car audio



Artificial data synthesis

Car recognition:



≈ 20 cars

