

Mismatched training and dev/test set

Tuesday, September 1, 2020 9:24 PM

① Training and testing on different distributions

Cat app example ↴

Data from webpages



$\rightarrow \approx 200,000$

care about this

Data from mobile app

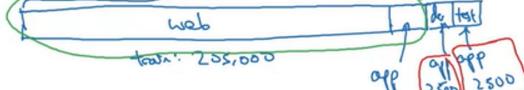


$\rightarrow \approx 10,000$

X Option 1:



Option 2:



Andrew Ng

\rightarrow Prioritize the data you care for in the dev/test set

② Bias and Variance with mismatched data distributions

\rightarrow Since Training and dev set are different, we can't evaluate bias and variance

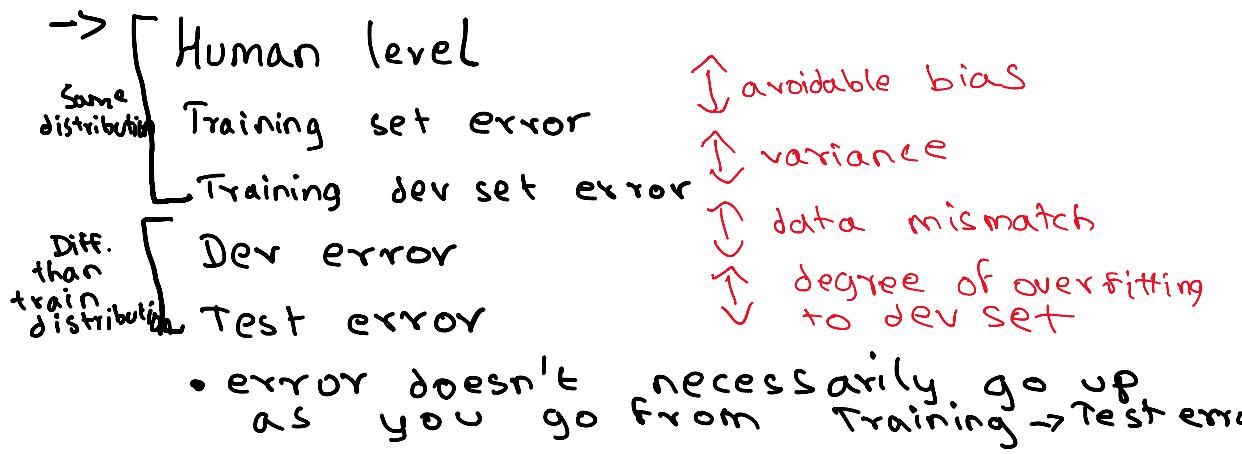
\rightarrow We define Training dev-set: Same distribution as training set, but not used for training

\rightarrow e.g.:

Training error	1%	\uparrow variance problem	1%
Training-dev error	9%		1.5%
Dev error	10%		10% \uparrow data mismatch problem

e.g.:

Human error	0%	\uparrow avoidable bias problem	0% \uparrow avoidable bias problem
Training error	10%		10% \uparrow data mismatch problem
Training-dev error	11%		11% \uparrow data mismatch problem
Dev error	12%		20% \uparrow data mismatch problem



③ Addressing Data Mismatch

- Carry out manual error analysis to try to understand difference between training and dev/test sets.
- Make training data more similar; or collect more data similar to dev/test sets (e.g. Artificial data synthesis (adding noise))
 - Artificial data synthesis might lead to overfitting of model to tiny subset of the data that we can simulate