# Mismatched training and dev/test data

---

# Bias and Variance with mismatched data distributions

deeplearning.ai

# Cat classifier example

Assume humans get ≈ 0% error.
(Bayes error)

Training error    1%
Dev error        10%

— Estimating Absolute bias & variance helps when you want to find what to focus on ↓, bias or variance?
— But, the way we calculate bias & variance is diff when training dist is diff compared to dev/test dist
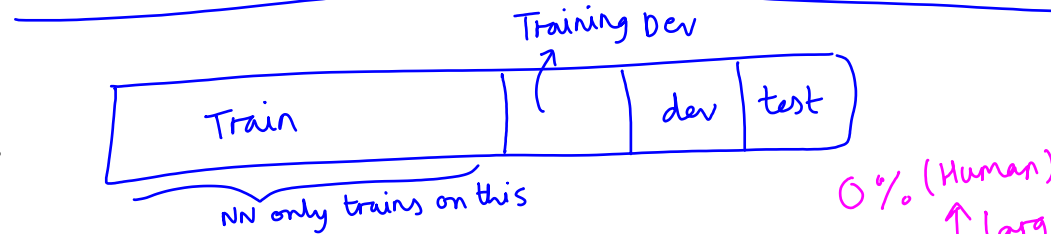
→ If dev data came from the same dist as training data, we would say we have large variance problem & work to ↓ variance

✗ — But what if it didn't come from same dist?

eg training data was very easy → high Res Images
= ∴ error = 1%, dev set is blurry Images ∴ 10%
⇒ its not a variance problem, its a data mismatch problem

New data set
(dont explicitly
train the NN on this)
→

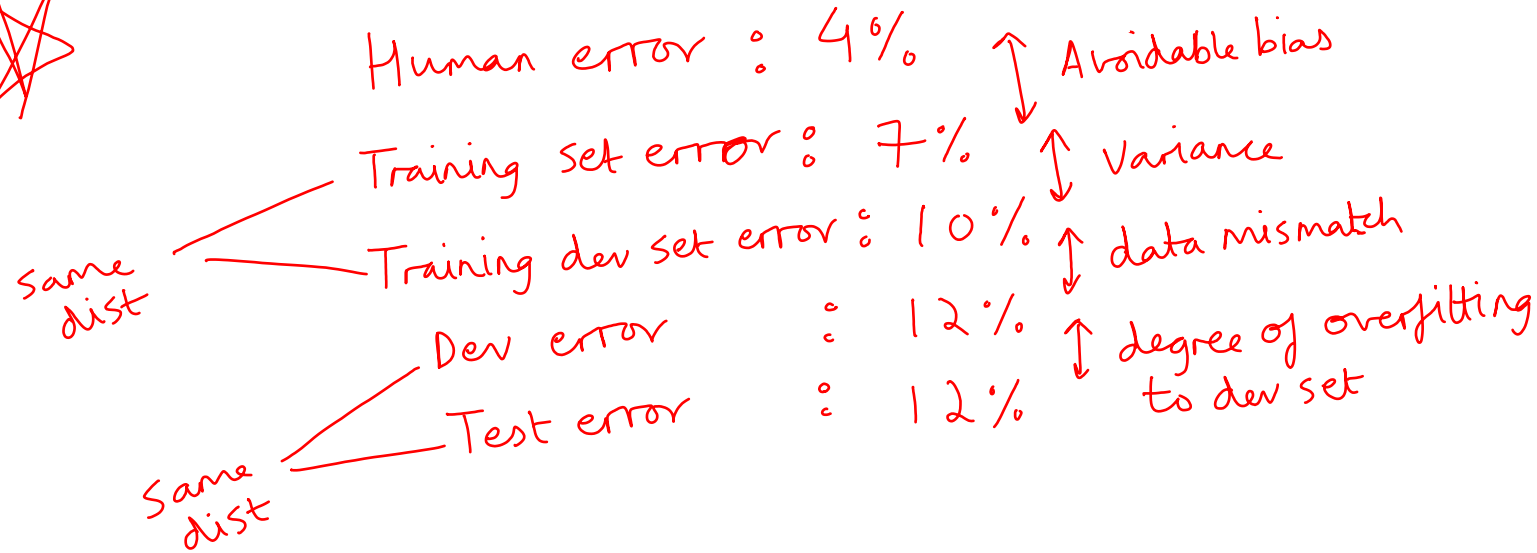Training-dev set: Same distribution as training set, but not used for training

Training Dev
↑

| Train | | dev | test |
|---|---|---|---|

NN only trains on this

Same dist
Training error   :   1%   ↑large
Training Dev error: 9%
Dev error    :   10%
Variance problem

1%
1.5% ↑large
10% ↑large
Data mismatch Problem

0% (Human) ↑large
10% ↑large
11%
12%
Avoidable bias

0% ↑large
10%
11% ↑large
20%
2 problems
① Avoidable bias
② Data mismatch

Andrew Ng

# Bias/variance on mismatched training and dev/test sets

Human error : 4%          ↕ Avoidable bias

Training set error : 7%    ↕ Variance

Training dev set error : 10%   ↕ data mismatch

Dev error    : 12%    ↕ degree of overfitting
                              to dev set

Test error   : 12%

same dist

same dist

4%

7% }

10% }

6% }   Dev/Test set
6% }   can have lesser
        error compared
        to Train sets
    ⟹ Training data
    was more difficult
    than dev/test set

Andrew Ng

# More general formulation

Speech Activated Rear view mirror

Purchased data, Samples from Amzn Echo etc.

General Speech Recognition data Samples

Rear view mirror speech data

Human level

Error on examples the NN has trained on

Error on examples the NN has NoT trained on

"Human Level" 4%

↕ Avoidable bias

"Training error" 7%

↕ Variance

"Training Dev Set error" 10%

"Dev/Test set error" 6%

↔ data mismatch

6% → How to get this? Make your Audience use the App

6% → How to get this? Train NN on samples from the Rear view App

⇒ If we get 6% Human & 6% training from NN, we're already at human level performance

We can handle avoidable bias & variance, How to handle data mis match?

Next slide

Andrew Ng