# Bias and variance with mismatched data distributions

Example: Cat classifier with mismatch data distribution

When the training set is from a different distribution than the development and test sets, the method to analyze bias and variance changes.

| | Classification error (%) | | | | | |
|---|---|---|---|---|---|---|
| | Scenario A | Scenario B | Scenario C | Scenario D | Scenario E | Scenario F |
| Human (proxy for Bayes error) | 0 | 0 | 0 | 0 | 0 | 4 |
| Training error | 1 | 1 | 1 | 10 | 10 | 7 |
| Training-development error | - | 9 | 1.5 | 11 | 11 | 10 |
| Development error | 10 | 10 | 10 | 12 | 20 | 6 |
| Test error | - | - | - | - | - | 6 |

Scenario A

If the development data comes from the same distribution as the training set, then there is a large variance problem and the algorithm is not generalizing well from the training set.

However, since the training data and the development data come from a different distribution, this conclusion cannot be drawn. There isn't necessarily a variance problem. The problem might be that the development set contains images that are more difficult to classify accurately.

When the training set, development and test sets distributions are different, two things change at the same time. First of all, the algorithm trained in the training set but not in the development set. Second of all, the distribution of data in the development set is different.

It's difficult to know which of these two changes what produces this 9% increase in error between the training set and the development set. To resolve this issue, we define a new subset called training-development set. This new subset has the same distribution as the training set, but it is not used for training the neural network. (Note: This training-development set is used to test whether there is a variance problem. )

Scenario B

The error between the training set and the training- development set is 8%. In this case, since the training set and training-development set come from the same distribution, the only difference between them is the neural network sorted the data in the training and not in the training development. The neural network is not generalizing well to data from the same distribution that it hadn't seen before

Therefore, we have really a variance problem.

Scenario C

In this case, we have a mismatch data problem since the 2 data sets come from different distribution.

Scenario D

In this case, the avoidable bias is high since the difference between Bayes error and training error is 10 %.

Scenario E

In this case, there are 2 problems. The first one is that the avoidable bias is high since the difference between Bayes error and training error is 10 % and the second one is a data mismatched problem.

Scenario F

Development should never be done on the test set. However, the difference between the development set and the test set gives the degree of overfitting to the development set.

Bayes error

↕ Avoidable Bias

Training set error

↕ Variance

Development - Training set error

↕ Data mismatch

Development set error

↕ Degree of overfitting to the development set

Test set error