

CHEAT SHEET

Model Debugging

If a classifier is not performing well, there are several ways to improve its performance. To determine which of the many techniques to use, the first step is to determine the root of the problem.

	High Variance	High Bias
Description	Models with high variance are capable of memorizing many more properties of the training data and do not do well on unseen data, resulting in low training error but high test error.	Models with high bias are too simplistic or make ill-suited assumptions and therefore cannot even achieve low error on the training data set.
Symptoms	Training error is much lower than test error.	Training error is higher than ϵ .
Remedies	<ul style="list-style-type: none"> • Add more training data • Reduce model complexity (complex models are prone to high variance) • Bagging 	<ul style="list-style-type: none"> • Use more complex model (e.g., Kernelize or use non-linear models) • Add features • Boosting

Visualize Variance and Bias:

The graph below exemplifies data with high/low variance and high/low bias as “darts” thrown at a target. The bullseye at the center of the target is the location of the perfect classifier on the testing data. The blue dots illustrate the darts, which represent classifiers trained on different training data sets.

High variance / low bias settings lead to classifiers that are unbiased; i.e., in expectation they are very close to the bullseye. However, different training sets lead to vastly different outcomes as the classifiers overspecialize to the data they were trained on.

High bias / low variance settings lead to classifiers that are very similar across different training data sets (the blue dots are close together); however, they are systematically off-target (i.e., they make wrong assumptions).

The worst case is **high bias and high variance**. The goal is to achieve a classifier with low bias and low variance.

