

STA 314: Statistical Methods for Machine Learning I

Lecture 3 - Bias-Variance Decomposition

Chris J. Maddison

University of Toronto

- Expand a bit on Q3 and Q4 of the HW1.
- Today we will talk about the [bias-variance](#) decomposition, which is beginning to make more precise our discussion of overfitting and underfitting last class.

- Given any finite set $\{x_i\}_{i=1}^N$ of $x_i \in \mathbb{R}$, we can define the uniform random variable over $\{x_i\}_{i=1}^N$, which is any D such that

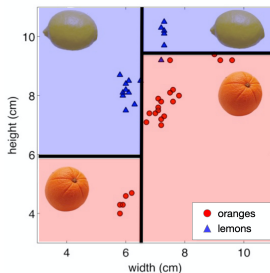
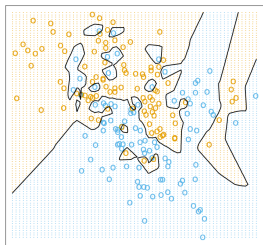
$$P(D = x_i) = \frac{1}{N}$$

- Sampling from this random variable is easy: sample an integer $J \in \{1, \dots, N\}$ uniformly at random and return x_J .
- For this distribution, we have

$$\mathbb{E}[D] = \sum_{i=1}^N P(D = x_i)x_i = \sum_{i=1}^N \frac{1}{N}x_i$$

Recall: supervised learning

- In supervised learning, our learning algorithms (k -NN, decision trees) produce predictions $\hat{y}^*(\mathbf{x}) \approx t$ for a query point \mathbf{x} .



Recall: supervised learning

- We can think of this as picking a predictor function $\hat{y}^* \in \mathcal{H}$ from a hypothesis class by minimizing the average loss on the training set

$$\hat{y}^* = \arg \min_{y \in \mathcal{H}} \hat{\mathcal{R}}[y, \mathcal{D}^{train}]$$

- Then, we measure the average loss on an **unseen test set** to approximate how well \hat{y}^* does on the true data generating distribution,

$$\hat{\mathcal{R}}[\hat{y}^*, \mathcal{D}_{test}] \approx \mathcal{R}[\hat{y}^*]$$

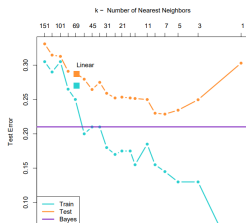
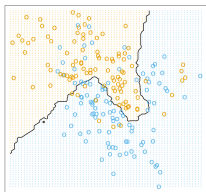
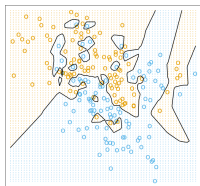
Recall: supervised learning

- This view of supervised learning is a very idealized view:
 - ▶ k -NN algorithm for $k > 1$ doesn't really select the predictor by minimizing a global loss.
 - ▶ Decision tree fitting does select \hat{y}^* based on training loss, but it is often greedy and sometimes does not find the global optimal \hat{y}^* .
- Still, it's a very useful general model for supervised learning.

- Let's consider Q4 in HW1 as a way to review this supervised framework.

Bias-Variance Decomposition

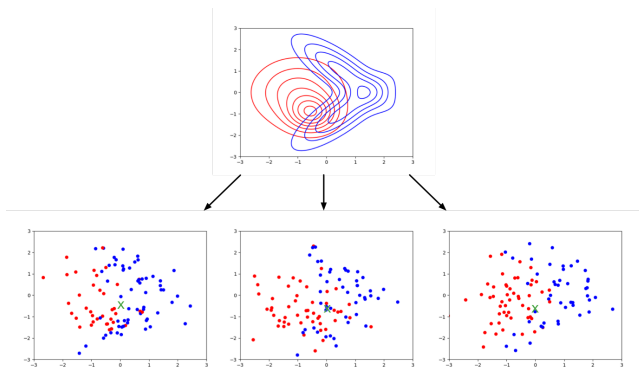
- Recall that overly simple hypothesis classes underfit the data, and overly complex ones overfit.



- Last lecture we talked about this intuitively.
- We can quantify this effect in terms of the **bias-variance decomposition**.
 - So far we've been talking about the training set as if it is fixed, but it makes more sense to think of it as random.
 - So, we'd like to understand how our learning algorithm is impacted by selecting a predictor on a finite, **random**, training set.

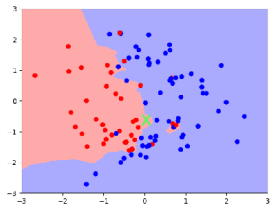
Bias-Variance Decomposition: Basic Setup

- Recall: the training set $\mathcal{D}^{train} = \{(\mathbf{x}^{(i)}, t^{(i)})\}_{i=1}^N$ contains N i.i.d. draws from a single **data generating distribution** p_{data} .
- Consider a **fixed query point** \mathbf{x} (green \mathbf{x} below).
- Consider **sampling many training sets** \mathcal{D}_n^{train} independently from p_{data} .

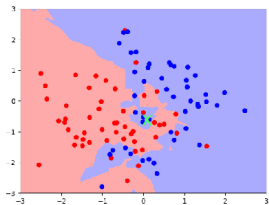


Bias-Variance Decomposition: Basic Setup

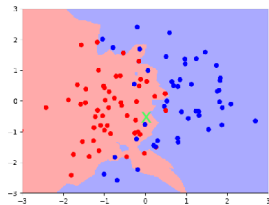
- For each training set \mathcal{D}_n^{train} , run learning alg. to get a predictor $\hat{y}_n^* \in \mathcal{H}$.
- Compute the prediction $\hat{y}_n^*(\mathbf{x})$ and compare it to a label t drawn from $p_{\text{data}}(t|\mathbf{x})$.
- We can view \hat{y}_n^* as a random variable, where the randomness comes from the choice of training set.



$y = \bullet$



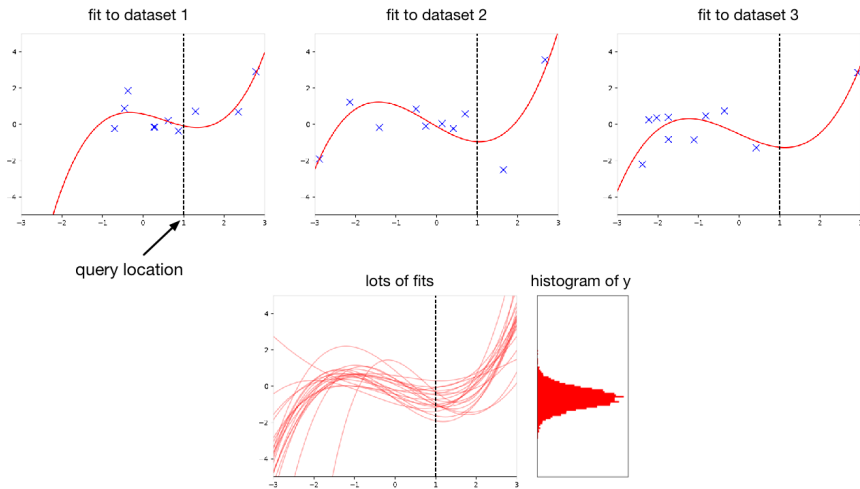
$y = \bullet$



$y = \bullet$

Bias-Variance Decomposition: Basic Setup

Here is the analogous setup for regression:



Bias-Variance Decomposition: Basic Setup

- Recap of basic setup:
 - ▶ Fix a query point \mathbf{x} .
 - ▶ Sample the (true) target t from the conditional distribution $p_{\text{data}}(t|\mathbf{x})$.
 - ▶ Repeat:
 - ▶ Sample a random training dataset $\mathcal{D}_n^{\text{train}}$ i.i.d. from the data generating distribution p_{data} .
 - ▶ Run the learning algorithm on $\mathcal{D}_n^{\text{train}}$ to get a prediction $\hat{y}_n^*(\mathbf{x})$ from \mathcal{H} at \mathbf{x} .
 - ▶ Compute the loss $L(\hat{y}_n^*(\mathbf{x}), t)$.
 - ▶ Average the losses.
- Notice: y is independent of t given \mathbf{x} .
- This gives a distribution over the loss at \mathbf{x} , with expectation $\mathbb{E}[L(\hat{y}^*(\mathbf{x}), t) | \mathbf{x}]$ taken over t and the *random* training set $\mathcal{D}^{\text{train}}$ where $\hat{y}^* = \arg \min_{y \in \mathcal{H}} \hat{\mathcal{R}}[y, \mathcal{D}^{\text{train}}]$.
- For each query point \mathbf{x} , the expected loss is different. We are interested in minimizing the expectation of this with respect to $\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$.

- For now, focus on squared error loss, $L(y, t) = \frac{1}{2}(y - t)^2$ with $y, t \in \mathbb{R}$.
- A first step: suppose we knew the conditional distribution $p_{\text{data}}(t | \mathbf{x})$. What is the best deterministic value $y(\mathbf{x}) \in \mathbb{R}$ should we predict?
 - ▶ Here, we are treating t as a random variable and choosing $y(\mathbf{x})$.
- **Claim:** $y^*(\mathbf{x}) = \mathbb{E}[t | \mathbf{x}]$ is the best possible prediction.
- **Proof:** Consider a fixed $y \in \mathbb{R}$

$$\begin{aligned}\mathbb{E}[(y - t)^2 | \mathbf{x}] &= \mathbb{E}[y^2 - 2yt + t^2 | \mathbf{x}] \\ &= y^2 - 2y\mathbb{E}[t | \mathbf{x}] + \mathbb{E}[t^2 | \mathbf{x}] \\ &= y^2 - 2y\mathbb{E}[t | \mathbf{x}] + \mathbb{E}[t | \mathbf{x}]^2 + \text{Var}[t | \mathbf{x}] \\ &= y^2 - 2yy^*(\mathbf{x}) + y^*(\mathbf{x})^2 + \text{Var}[t | \mathbf{x}] \\ &= (y - y^*(\mathbf{x}))^2 + \text{Var}[t | \mathbf{x}]\end{aligned}$$

$$\mathbb{E}[(y - t)^2 | \mathbf{x}] = (y - y^*(\mathbf{x}))^2 + \text{Var}[t | \mathbf{x}]$$

- The first term is nonnegative, and can be made 0 by setting $y = y^*(\mathbf{x})$.
- The second term corresponds to the inherent unpredictability, or **noise**, of the targets, and is called the **Bayes error**.
 - ▶ This is the best we can ever hope to do with any learning algorithm. An algorithm that achieves it is **Bayes optimal**.
 - ▶ Notice that this term doesn't depend on y .
- This process of choosing a single value $y^*(\mathbf{x})$ based on $p_{\text{data}}(t | \mathbf{x})$ is an example of **decision theory**.

Bayes Optimality

- But, in practice, our prediction $\hat{y}^*(\mathbf{x})$ is not $y^*(\mathbf{x})$. Instead, it is a random variable (where the randomness comes from randomness of the training set) taking values in \mathcal{H} .
- We can decompose out the expected loss.
- Suppressing the dependence on \mathbf{x} for clarity:

$$\begin{aligned}\mathbb{E}[(\hat{y}^* - t)^2] &= \mathbb{E}[(\hat{y}^* - y^*)^2] + \text{Var}(t) \\ &= \mathbb{E}[y^{*2} - 2y^*\hat{y}^* + \hat{y}^{*2}] + \text{Var}(t) \\ &= y^{*2} - 2y^* \mathbb{E}[\hat{y}^*] + \mathbb{E}[\hat{y}^{*2}] + \text{Var}(t) \\ &= y^{*2} - 2y^* \mathbb{E}[\hat{y}^*] + \mathbb{E}[\hat{y}^*]^2 + \text{Var}(\hat{y}^*) + \text{Var}(t) \\ &= \underbrace{(y^* - \mathbb{E}[\hat{y}^*])^2}_{\text{bias}} + \underbrace{\text{Var}(\hat{y}^*)}_{\text{variance}} + \underbrace{\text{Var}(t)}_{\text{Bayes error}}\end{aligned}$$

- Let's step back and consider what we just did. First, recall:
 - ▶ Picking a predictor by minimizing the average loss on the training set

$$\hat{y}^* = \arg \min_{y \in \mathcal{H}} \hat{\mathcal{R}}[y, \mathcal{D}^{train}]$$

returns a random predictor \hat{y}^* .

- ▶ But, we're interested in our performance in terms of expected loss:

$$\mathcal{R}[\hat{y}^*]$$

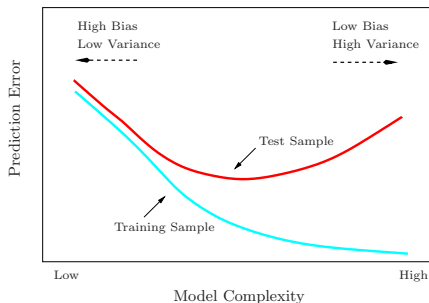
- In our case:

$$\mathcal{R}[\hat{y}^*] = \mathbb{E} \left[\mathbb{E} [(\hat{y}^*(\mathbf{x}) - t)^2 \mid \mathbf{x}] \right].$$

$$\mathbb{E} [\mathbb{E}[(\hat{y}^*(\mathbf{x}) - t)^2 | \mathbf{x}]] = \mathbb{E} \left[\underbrace{(y^*(\mathbf{x}) - \mathbb{E}[\hat{y}^*(\mathbf{x}) | \mathbf{x}])^2}_{\text{bias}} + \underbrace{\text{Var}[\hat{y}^*(\mathbf{x}) | \mathbf{x}]}_{\text{variance}} + \underbrace{\text{Var}[t | \mathbf{x}]}_{\text{Bayes error}} \right]$$

- So, we just split the expected loss $\mathcal{R}[\hat{y}^*]$ into three terms:
 - ▶ **bias**: how wrong the expected prediction is
 - ▶ **variance**: the amount of variability in the predictions
 - ▶ **Bayes error**: the inherent unpredictability of the targets
- How does our choice of \mathcal{H} interact with this analysis?

Bayes Optimality

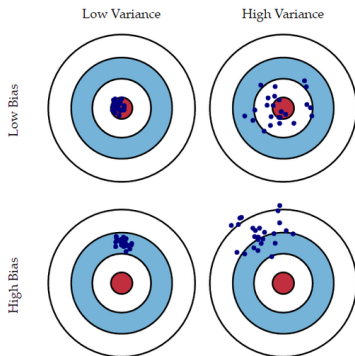


Source: ESL

- If \mathcal{H} is large, then \hat{y}^* can get close y^* , therefore reducing bias. It's also sensitive to the finite training set, therefore increasing variance.
- If \mathcal{H} is small, then \hat{y}^* is typically from y^* , therefore increasing bias. It's less sensitive to the finite training set, therefore reducing variance.
- Even though this analysis only applies to squared error, we often loosely use “bias” and “variance” as synonyms for “underfitting” and “overfitting”.

Bias and Variance

- Throwing darts = predictions for each draw of a dataset

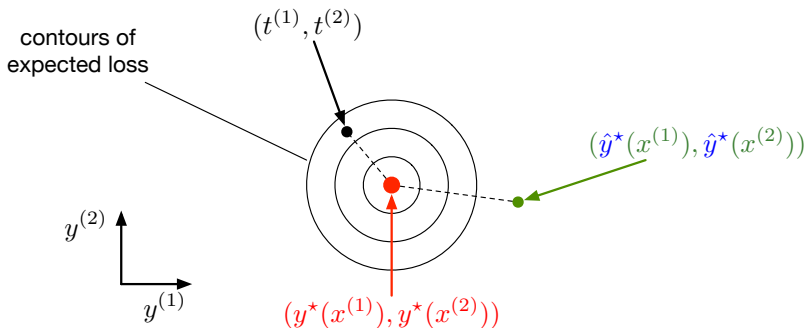


Source: ESL.

- Be careful, the expected loss averages over points \mathbf{x} from the data distribution, so this produces its own type of variance.

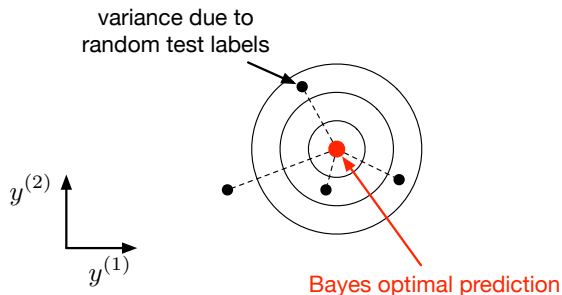
Bias/Variance Decomposition: Another Visualization

- In practice, measure the average loss $\hat{\mathcal{R}}[\hat{y}^*, \mathcal{D}_{test}]$ on the test set instead of $\mathcal{R}[\hat{y}^*]$.
- Let's visualize the bias-variance decomposition by plotting the space of predictions of the model, where each axis correspond to predictions on a two test examples $(\mathbf{x}^{(1)}, \mathbf{x}^{(2)})$.



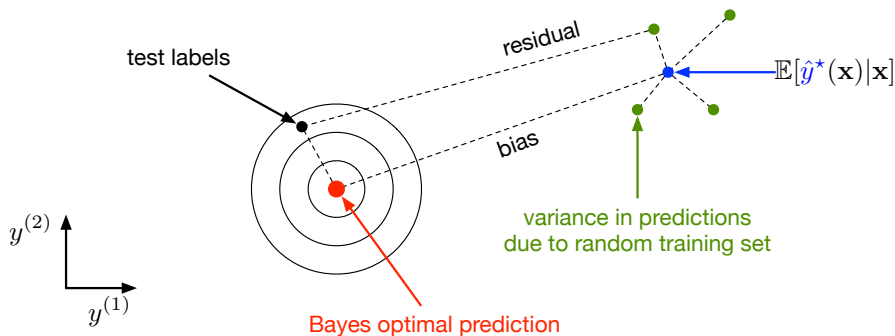
Bias/Variance Decomposition: Another Visualization

- The Bayes error is an irreducible error that comes from the randomness in $p_{\text{data}}(t | \mathbf{x})$.



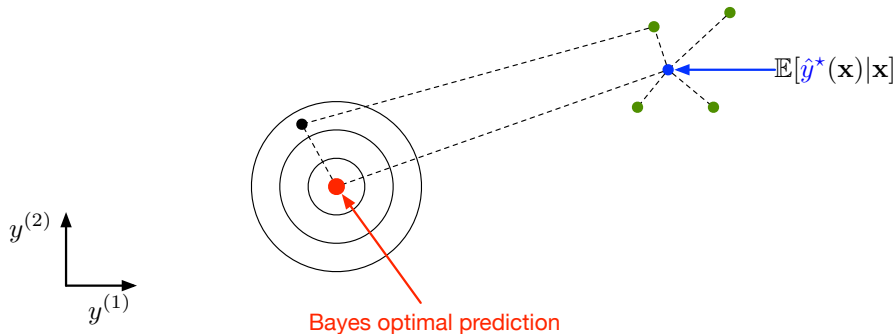
Bias/Variance Decomposition: Another Visualization

- Selecting a predictor $\hat{y}^* \in \mathcal{H}$ from a training set comes with bias and variance.



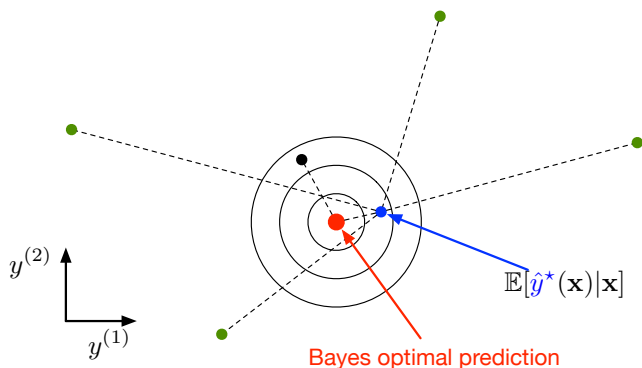
Bias/Variance Decomposition: Another Visualization

- An overly simple model (e.g. k -NN with large k) might have
 - ▶ high bias (too simplistic to capture the structure in the data)
 - ▶ low variance (there's enough data to get a stable estimate of the decision boundary)



Bias/Variance Decomposition: Another Visualization

- An overly complex model (e.g. KNN with $k = 1$) may have
 - ▶ low bias (since it learns all the relevant structure)
 - ▶ high variance (it fits the quirks of the data you happened to sample)

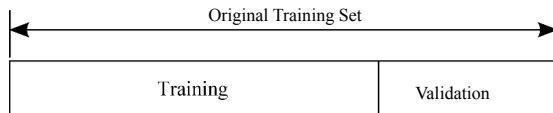


- Before we move on to bagging, it's a good time to mention **validation**.
- We may **want to assess how likely a learning algorithm is to generalize before picking one and reporting the final test error**.
- In other words, until now we've been picking predictors that optimize the training loss, but we want a technique for picking predictors that are likely to generalize as well.

- For example, we may want to assess the following types of choices:
 1. **Hyper-parameters of the learning algorithm that lead to better generalization.** Often there are parameters that cannot be fit on the training set, e.g., k in k -NN, because the training set would give meaningless answers about the best setting, i.e., $k = 1$ is always gives optimal training set loss for k -NN.
 2. **Picking predictors that generalize better.** E.g., should we use a decision tree or k -NN if we want to generalize?
- **We make these choices using validation** to avoid measuring test loss (then the test set would no longer be unseen data!).
- Suppose we are trying to estimate the generalization of two learning algorithms, e.g., a decision tree and a k -NN model.

Hold-out validation

- The most common method of validation is to hold-out a subset of the training set and use it to assess how likely we are to generalize to unseen data.



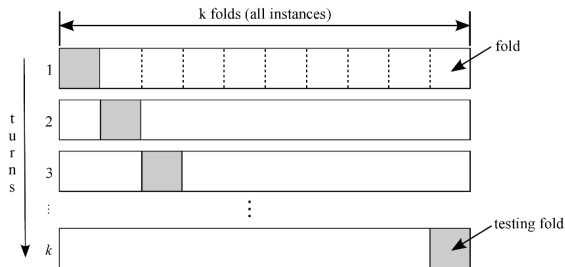
- In our example of deciding between a decision tree and k -NN in terms of generalization, we would fit \hat{y}_{kNN}^* and \hat{y}_{d-tree}^* on the training set and measure the average loss on the validation set

$$\hat{\mathcal{R}}[\hat{y}_{kNN}^*, \mathcal{D}^{valid}] \text{ vs. } \hat{\mathcal{R}}[\hat{y}_{d-tree}^*, \mathcal{D}^{valid}]$$

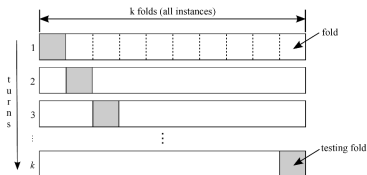
- We pick the predictor \hat{y}_{kNN}^* vs. \hat{y}_{d-tree}^* with lowest validation loss.
- Problem: this is usually a waste of data.

K -fold cross validation

- Second most common way: partition training data randomly into K equally sized subsets. For each “turn”, use the first $K - 1$ subsets (or “folds”) as training data and the last subset as validation



K-fold cross validation



- In our running example: fit a new predictor using each learning algorithm on $K - 1$ folds for each of the K turns, and measure the validation loss on the held-out fold, averaged over the turns:

$$\frac{1}{K} \sum_{i=1}^K \hat{\mathcal{R}}[\hat{y}_{kNN,i}^*, \mathcal{D}_i^{valid}] \text{ vs. } \frac{1}{K} \sum_{i=1}^K \hat{\mathcal{R}}[\hat{y}_{d-tree,i}^*, \mathcal{D}_i^{valid}]$$

where $\hat{y}_{A,i}^*$ is the predictor fit on the training subset of the i th turn using algorithm A and \mathcal{D}_i^{valid} is the validation subset of the i th turn.

- We pick the learning algorithm, e.g., k -NN v. decision tree, with lowest validation loss averaged across the K turns.