

# Model Flexibility - Takeaways

- If you **don't have too many samples**, you should prefer a **simpler model**
- If you have **many samples**, you can afford a **more complex model**
- We'll need **some sort of mechanism** to tell which regime we're in
- Examples of different model families
  - Useful tutorial: pyGAM

# Model Flexibility - Takeaways

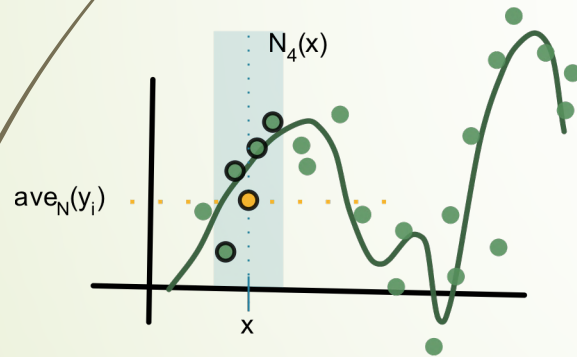
- If you **don't have too many samples**, you should prefer a **simpler model**
- If you have **many samples**, you can afford a **more complex model**
- We'll need **some sort of mechanism** to tell which regime we're in
- Examples of different model families
  - Useful tutorial: pyGAM
- **Reading:** ISLR 2.1, 3.2.1, 3.5

# Algorithms: K-Nearest Neighbors (KNN) Regression

- $K$  fixed and given
- **Samples:**  $(x_i, y_i)_{i=1}^N$
- **Estimate data generating function:**  $\hat{f}(x) = \frac{1}{K} \sum_{i \in N_K(x)} y_i$
- $N_K$ :  $K$  nearest neighbors of  $x$  within the **training set**

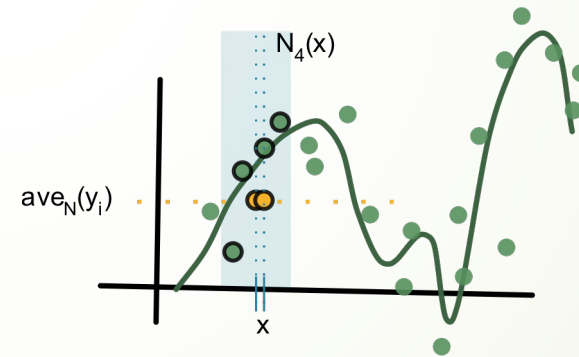
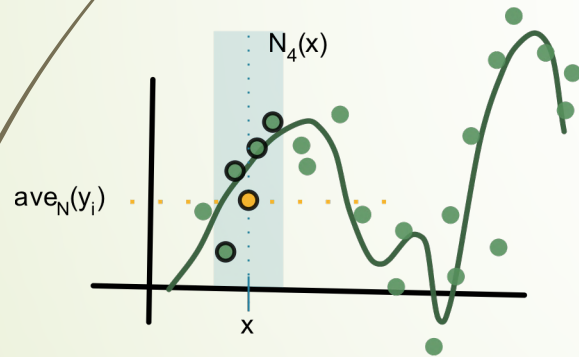
# Algorithms: K-Nearest Neighbors (KNN) Regression

- $K$  fixed and given
- **Samples:**  $(x_i, y_i)_{i=1}^N$
- **Estimate data generating function:**  $\hat{f}(x) = \frac{1}{K} \sum_{i \in N_K(x)} y_i$
- $N_K$ :  $K$  nearest neighbors of  $x$  within the **training set**



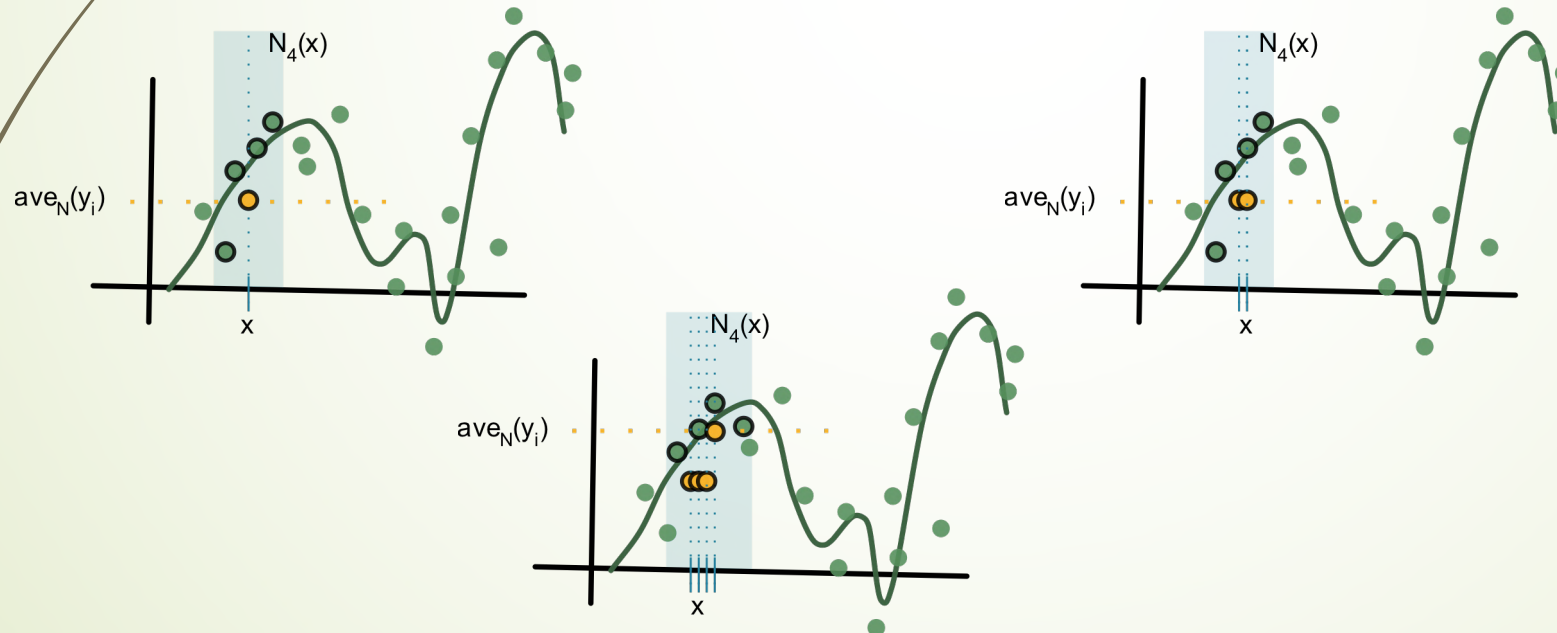
# Algorithms: K-Nearest Neighbors (KNN) Regression

- $K$  fixed and given
- **Samples:**  $(x_i, y_i)_{i=1}^N$
- **Estimate data generating function:**  $\hat{f}(x) = \frac{1}{K} \sum_{i \in N_K(x)} y_i$
- $N_K$ :  $K$  nearest neighbors of  $x$  within the **training set**



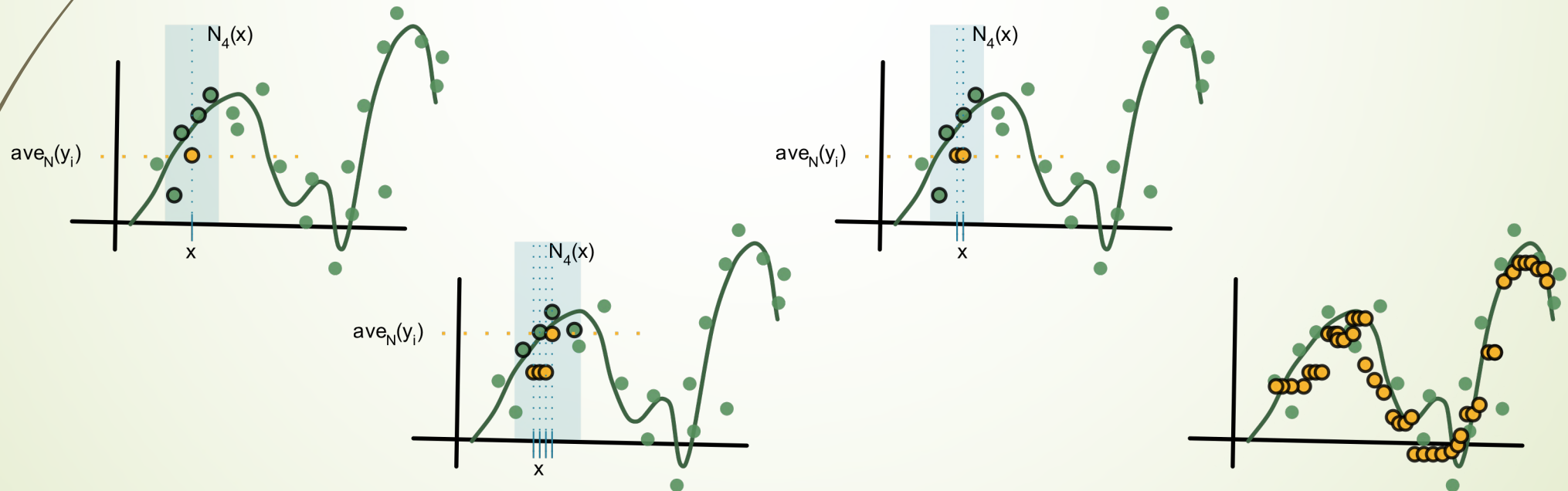
# Algorithms: K-Nearest Neighbors (KNN) Regression

- $K$  fixed and given
- **Samples:**  $(x_i, y_i)_{i=1}^N$
- **Estimate data generating function:**  $\hat{f}(x) = \frac{1}{K} \sum_{i \in N_K(x)} y_i$
- $N_K$ :  $K$  nearest neighbors of  $x$  within the **training set**



# Algorithms: K-Nearest Neighbors (KNN) Regression

- $K$  fixed and given
- **Samples:**  $(x_i, y_i)_{i=1}^N$
- **Estimate data generating function:**  $\hat{f}(x) = \frac{1}{K} \sum_{i \in N_K(x)} y_i$
- $N_K$ :  $K$  nearest neighbors of  $x$  within the **training set**







## Impact of $K$ : Bias-Variance Trade-off







# Impact of $K$ : Bias-Variance Trade-off

- **Model complexity** is controlled by the **size of the neighborhood**

# Impact of $K$ : Bias-Variance Trade-off

- ▶ **Model complexity** is controlled by the **size of the neighborhood**
  - ▶ Large  $K$  → Lower variance, larger bias
  - ▶ Small  $K$  → Higher variance, smaller bias

# Impact of $K$ : Bias-Variance Trade-off

- ▶ **Model complexity** is controlled by the size of the neighborhood
  - ▶ Large  $K$  → Lower variance, larger bias
  - ▶ Small  $K$  → Higher variance, smaller bias
- ▶ Larger  $K$  learns smoother functions; smaller  $K$  can match more complex functions when the sampling density is high enough



# Curse of dimensionality





# Curse of dimensionality

- The density of samples decreases with increase in dimension
- 

# Curse of dimensionality

- The density of samples decreases with increase in dimension

