

1. Review: Non-linear Basis Functions
2. Overfitting and Regularization
3. Hyperparameter Tuning and Cross-Validation
4. Bias-Variance Trade-off

## **Review: Non-linear Basis Functions**

---

# General Nonlinear Basis Functions

We can use a nonlinear mapping to a new feature vector:

$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^D \rightarrow \mathbf{z} \in \mathbb{R}^M$$

- $M$  is dimensionality of new features  $\mathbf{z}$  (or  $\phi(\mathbf{x})$ )
- $M$  could be greater than, less than, or equal to  $D$

We can apply existing learning methods on the transformed data:

- linear methods: prediction is based on  $\mathbf{w}^\top \phi(\mathbf{x})$

## Residual sum of squares

$$\sum_n [\mathbf{w}^\top \phi(\mathbf{x}_n) - y_n]^2$$

where  $\mathbf{w} \in \mathbb{R}^M$ , the same dimensionality as the transformed features  $\phi(\mathbf{x})$ .

The LMS solution can be formulated with the new design matrix

$$\Phi = \begin{pmatrix} \phi(\mathbf{x}_1)^\top \\ \phi(\mathbf{x}_2)^\top \\ \vdots \\ \phi(\mathbf{x}_N)^\top \end{pmatrix} \in \mathbb{R}^{N \times M}, \quad \mathbf{w}^{\text{LMS}} = \left( \Phi^\top \Phi \right)^{-1} \Phi^\top \mathbf{y}$$

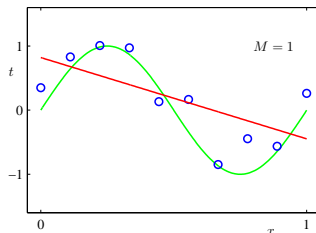
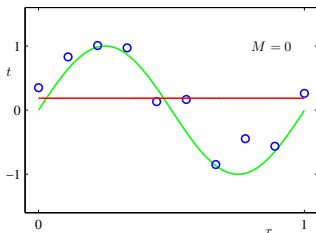
# Non-linear Basis Functions: Polynomial Regression

## Polynomial basis functions

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow f(x) = w_0 + \sum_{m=1}^M w_m x^m$$

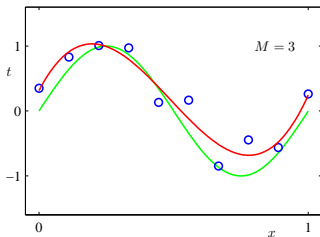
Fitting samples from a sine function:

underfitting since  $f(x)$  is too simple

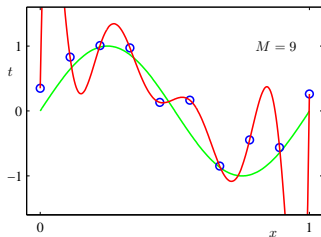


# Adding Higher-order Terms

$M=3$



$M=9$ : **overfitting**



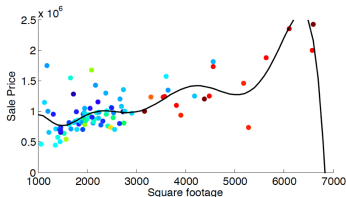
More complex features lead to better results on the training data, but potentially worse results on new data, e.g., test data!

# Overfitting and Regularization

---

# Overfitting Can Be Quite Disastrous

Fitting the housing price data with large  $M$ :



Predicted price goes to zero (and is ultimately negative) if you buy a big enough house!

This is called poor generalization/overfitting.

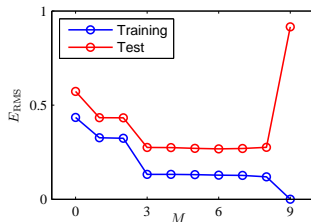


# Detecting Overfitting

## Plot model complexity versus objective function:

- X axis: model complexity, e.g.,  $M$
- Y axis: error, e.g., RSS, RMS  
(square root of RSS), 0-1 loss

Compute the objective on a **training** (used to train the model) and **test** (used to evaluate the model) dataset.

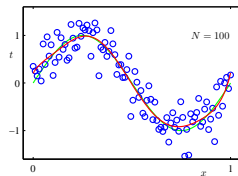
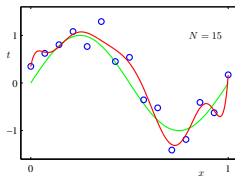
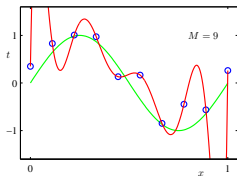


As a model increases in complexity:

- Training error keeps reducing
- Test error may first reduce but eventually increase

# Preventing Overfitting: Option 1

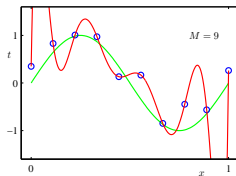
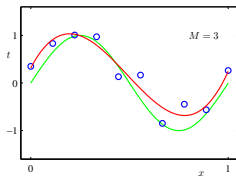
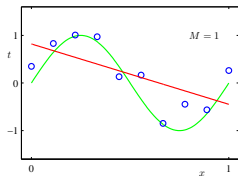
Try to use more training data



But getting a lot of data can be expensive and time-consuming

# Preventing Overfitting: Option 2

## Reduce the Number of Features



Q: how to select the right  $M$ ? - will talk about that later.

# Preventing Overfitting: Option 3

## Regularization Methods: Give preference to 'simpler' models

- How do we define a simple linear regression model —  $\mathbf{w}^\top \mathbf{x}$ ?

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
$w_0$	0.19	0.82	0.31	0.35
$w_1$		-1.27	7.99	232.37
$w_2$			-25.43	-5321.83
$w_3$			17.37	48568.31
$w_4$				-231639.30
$w_5$				640042.26
$w_6$				-1061800.52
$w_7$				1042400.18
$w_8$				-557682.99
$w_9$				125201.43

- Intuitively, the weights corresponding to higher order terms should not be “too large”

# Regularization Methods

## Add a term to the objective function.

Choose the parameters to not just minimize risk, but avoid being large.

$$\frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \frac{1}{2} \lambda \|\mathbf{w}\|_2^2$$

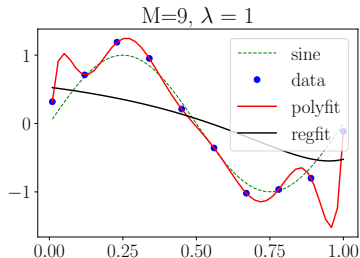
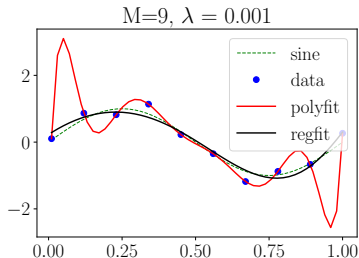
Ridge regression is just regularized linear regression.

## Advantages

- Forces the magnitude of  $\mathbf{w}$  to be small
- Alleviates overfitting and generalizes well to new data points

## Example: Effect of Regularization

As  $\lambda$  increases, the overfitting issue is alleviated.



**How to choose  $\lambda$ ? Up-next: hyperparameter tuning**

# Hyperparameter Tuning and Cross-Validation

---

# How Much Regularization Do We Need?

## Can we tune $\lambda$ on the training dataset?

i.e. try different  $\lambda$  and choose the one with the smallest training error.

**No:** In the overfitting example, we would pick  $\lambda = 0$  as it already has ZERO training error. This defeats our intention of alleviating overfitting.

To tune  $\lambda$ ,

- We can use a validation set or do cross validation.
- Pick the value of  $\lambda$  that yields lowest error on the **validation dataset**.

Similar idea applies to tuning degree of polynomial  $M$ , learning rate  $\eta$  (or any other hyperparameter) as well.



# Tuning with a Validation Dataset

**Training data are used to learn  $f(\cdot)$  (our model).**

N samples/instances:  $\mathcal{D}^{\text{TRAIN}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$

**Test data are used to assess the prediction error.**

- M samples/instances:  $\mathcal{D}^{\text{TEST}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_M, y_M)\}$
- They are used for assessing how well  $f(\cdot)$  will do in predicting an unseen  $\mathbf{x} \notin \mathcal{D}^{\text{TRAIN}}$

**Validation data are used to optimize hyperparameter(s).**

L samples/instances:  $\mathcal{D}^{\text{VAL}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_L, y_L)\}$

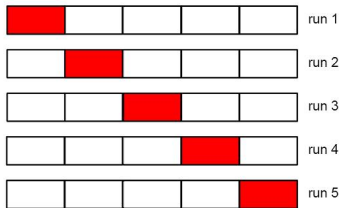
Training, validation and test data **should not overlap!**

- For each possible value of the hyperparameter (say  $\lambda = 1, 3, \dots, 100$ )
  - Train a model using  $\mathcal{D}^{\text{TRAIN}}$
  - Evaluate the performance of the model on  $\mathcal{D}^{\text{VAL}}$
- Choose the  $\lambda$  with the best performance on  $\mathcal{D}^{\text{VAL}}$
- Evaluate this model on  $\mathcal{D}^{\text{TEST}}$  to get the final prediction error

## What if we do not have validation data?

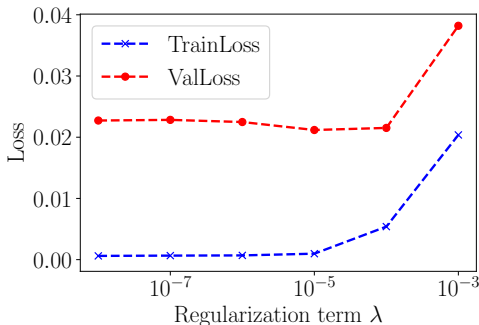
- We split the training data into  $S$  equal parts.
- We use **each part in turn** as a validation dataset and use the others as a training dataset.
- We choose the hyperparameter such that the model performs the best over  $S$  trials (based on average, variance, etc.)
- Re-train with this hyperparameter on the entire training dataset.

**Figure 1:**  $S = 5$ : 5-fold cross validation



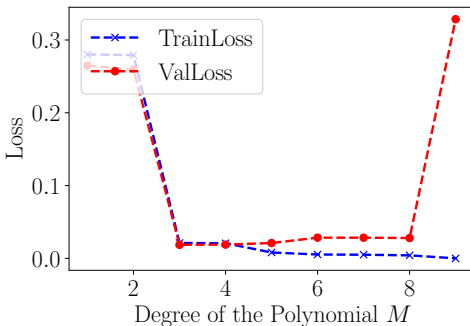
## Example: Hyper-parameter Tuning $\lambda$

- $\lambda = 10^{-4}$  gives the smallest validation loss
- Strikes a balance between over- and under-fitting



## Example: Hyper-parameter Tuning $M$

- Considering polynomial regression without regularization
- $M = 3$  or  $M = 4$  gives the smallest validation loss
- Strikes a balance between over- and under-fitting



## **Bias-Variance Trade-off**

---

# Empirical Risk Minimization

## Supervised learning

We aim to build a function  $h(\mathbf{x})$  to predict the true value  $y$  associated with  $\mathbf{x}$ . If we make a mistake, we incur a **loss**

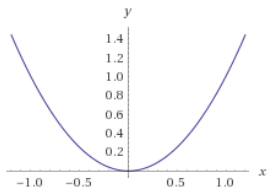
$$\ell(h(\mathbf{x}), y)$$

## Example:

Quadratic loss function for regression when  $y$  is continuous:

$$\ell(h(\mathbf{x}), y) = [h(\mathbf{x}) - y]^2$$

Ex: linear regression



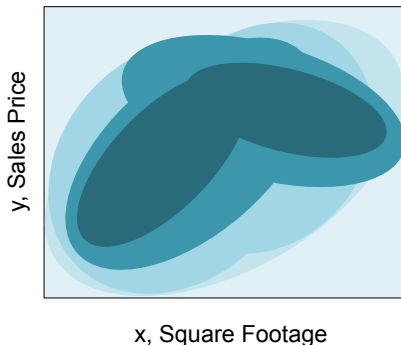
# Bias-Variance Trade-off: Illustration

How to evaluate  $\ell(h(\mathbf{x}), y)$  “generally” over  $(\mathbf{x}, y)$ ?

Using “true” data joint distribution  $p(\mathbf{x}, y)$ .

## Example:

- Joint distribution of square footage  $x$  and house sales price  $y$
- Darker color indicates higher probability regions





# How Good Is Our Predictor?

## Risk:

Given the true distribution of data  $p(\mathbf{x}, y)$ , the **risk** of a given predictor  $h(\mathbf{x})$  is its expected loss  $\ell$ :

$$R[h(\mathbf{x})] = \int_{\mathbf{x}, y} \ell(h(\mathbf{x}), y) p(\mathbf{x}, y) d\mathbf{x} dy$$

However, we cannot compute  $R[h(\mathbf{x})]$  (we do not know  $p$ ), so we use the **empirical risk**, given a training dataset  $\mathcal{D}$ :

$$R^{\text{EMP}}[h(\mathbf{x})] = \frac{1}{N} \sum_n \ell(h(\mathbf{x}_n), y_n)$$

Intuitively, as  $N \rightarrow +\infty$ ,

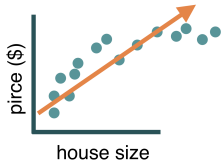
$$R^{\text{EMP}}[h(\mathbf{x})] \rightarrow R[h(\mathbf{x})]$$

# Connection to Underfitting/Overfitting

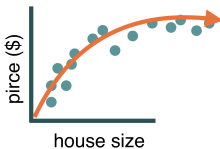
## Error decomposes into 3 terms

$$\text{Expected Risk} = \text{VARIANCE} + \text{BIAS}^2 + \text{NOISE}$$

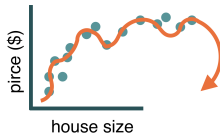
And underfitting/overfitting can be interpreted as **bias/variance** tradeoff.



**Figure 2:** High Bias



**Figure 3:** Just Right



**Figure 4:** High Variance

**Next:** Set up the notations needed to state and derive the decomposition.

# Bias-Variance Tradeoff for Regression

- $\mathcal{D}$ : our training data
- $\ell(h(\mathbf{x}), y)$ : our square loss function for regression

$$\ell(h(\mathbf{x}), y) = [h(\mathbf{x}) - y]^2$$

- $h_{\mathcal{D}}(\mathbf{x})$ : our prediction function learned from training data  $\mathcal{D}$ 
  - We are using the subscript  $\mathcal{D}$  to indicate that the prediction function is learned on the specific set of training data  $\mathcal{D}$
  - $h_{\mathcal{D}}(\mathbf{x})$  is the model within a model class (e.g. linear, degree  $M$  polynomial) that minimizes the training loss, e.g. the (regularized) empirical risk based on  $\mathcal{D}$ .
- Unknown joint distribution  $p(\mathbf{x}, y)$
- Risk of  $h_{\mathcal{D}}$ :  $R[h_{\mathcal{D}}(\mathbf{x})] = \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy$

# The Effect of Finite Training Samples

Every training sample  $\mathcal{D}$  is a sample from the following joint distribution of all possible training datasets

$$\mathcal{D} \sim P(\mathcal{D}) = \prod_{n=1}^N p(\mathbf{x}_n, y_n)$$

Thus, the prediction function  $h_{\mathcal{D}}(\mathbf{x})$  is a random function with respect to this distribution of possible training datasets. So is also its risk

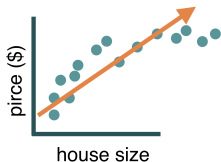
$$R[h_{\mathcal{D}}(\mathbf{x})] = \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

We will now evaluate the **expected risk**  $\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})]$ : the average risk over the distribution of possible training datasets,  $P(\mathcal{D})$ .

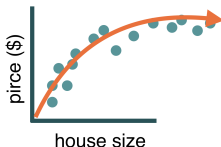
# Bias-Variance Trade-off: Intuition

Error decomposes into 3 terms

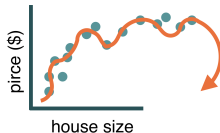
$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \text{VARIANCE} + \text{BIAS}^2 + \text{NOISE}$$



**Figure 5:** High Bias



**Figure 6:** Just Right



**Figure 7:** High Variance

We will prove this result, and interpret what it means...

**Warning:** The next few slides are somewhat mathematical – please review them carefully after class.

# Average over the Distribution of the Training Data

## Expected risk

$$\mathbb{E}_{\mathcal{D}} [R[h_{\mathcal{D}}(\mathbf{x})]] = \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}$$

Namely, the randomness with respect to  $\mathcal{D}$  is marginalized out.

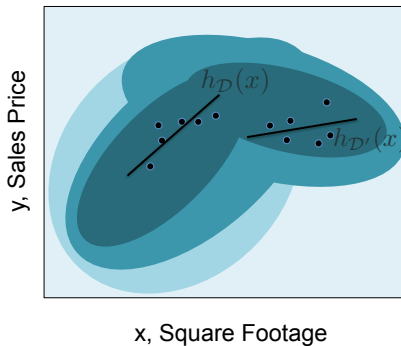
## Averaged prediction

$$\mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}(\mathbf{x})] = \int_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) P(\mathcal{D}) d\mathcal{D}$$

Namely, if we have seen many training datasets, we predict with the average of the prediction functions learned on each training dataset.

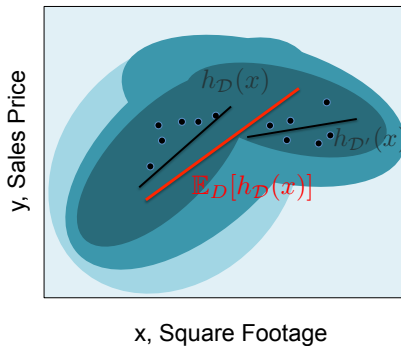
# Bias-Variance Trade-off: Illustration

- Learning the model for a different dataset  $\mathcal{D}'$  yields a new linear model  $h_{\mathcal{D}'}(x)$
- Average of such models over infinitely many datasets sampled from the joint distribution is  $\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})$



# Bias-Variance Trade-off: Illustration

- Learning the model for a different dataset  $\mathcal{D}'$  yields a new linear model  $h_{\mathcal{D}'}(x)$
- Average of such models over infinitely many datasets sampled from the joint distribution is  $\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(x)$





We will subtract the averaged prediction from the averaged risk

$$\begin{aligned}\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] &= \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D} \\&= \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) \\&\quad + \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D} \\&= \underbrace{\int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}}_{\text{VARIANCE}} \\&\quad + \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}\end{aligned}$$

# Where Does the Cross-term Go?

**It is zero**

$$\begin{aligned} & \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})][\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y] p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D} \\ &= \int_{\mathbf{x}} \int_y \left\{ \int_{\mathcal{D}} [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})] P(\mathcal{D}) d\mathcal{D} \right\} [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y] p(\mathbf{x}, y) d\mathbf{x} dy \\ &= \int_{\mathbf{x}} \int_y \left\{ \int_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) P(\mathcal{D}) d\mathcal{D} - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) \right\} [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y] p(\mathbf{x}, y) d\mathbf{x} dy \\ &= 0 \leftarrow \text{(the term within the braces vanishes, by definition)} \end{aligned}$$

We will subtract the averaged prediction from the averaged risk

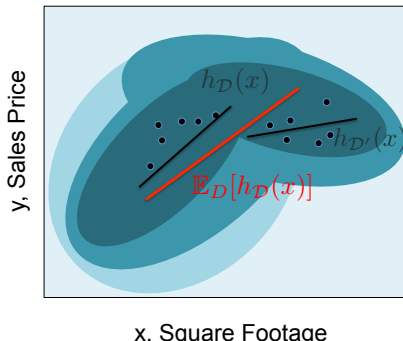
$$\begin{aligned}\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] &= \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D} \\&= \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) \\&\quad + \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D} \\&= \underbrace{\int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}}_{\text{VARIANCE}} \\&\quad + \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}\end{aligned}$$

Now let's analyze the Variance term!

# Bias-Variance Trade-off: Illustration

$$\underbrace{\int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}}_{\text{VARIANCE: error due to training dataset}}$$

- $\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})$ : Average of such models over infinitely many datasets sampled from the joint distribution.
- Variance term captures how much individual models differ from the average



# Analyzing the Variance

$$\underbrace{\int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}}_{\text{VARIANCE: error due to training dataset}}$$

- For each  $(\mathbf{x}, y)$  pair, we compute the squared difference of  $h_{\mathcal{D}}(\mathbf{x})$  (the prediction with training dataset  $\mathcal{D}$ ) and the averaged prediction  $\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})$ .
- High variances correspond to **Overfitting!**

## How can we reduce the variance?

- Use a lot of data (ie, increase the size of  $\mathcal{D}$ )
- Use a simple  $h(\cdot)$  so that  $h_{\mathcal{D}}(\mathbf{x})$  does not vary much across different training datasets. An extreme example is  $h(\mathbf{x}) = 0$  (not sensitive to  $\mathcal{D}$  at all).

## The Remaining Item

$$\begin{aligned}\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] &= \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D} \\ &\quad + \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}\end{aligned}$$

The integrand has no dependency on  $\mathcal{D}$  anymore and simplifies to

$$\int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

We will apply a similar add-and-subtract trick, by using an averaged target  $y$  (what we want to predict from  $\mathbf{x}$ ):

$$\mathbb{E}_y[y|\mathbf{x}] = \int_y y p(y|\mathbf{x}) dy$$

Decompose again

$$\begin{aligned} & \int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy \\ &= \int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y|\mathbf{x}] + \mathbb{E}_y[y|\mathbf{x}] - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy \\ &= \underbrace{\int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y|\mathbf{x}]]^2 p(\mathbf{x}, y) d\mathbf{x} dy}_{\text{BIAS}^2} \\ &\quad + \underbrace{\int_{\mathbf{x}} \int_y [\mathbb{E}_y[y|\mathbf{x}] - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy}_{\text{NOISE}} \end{aligned}$$

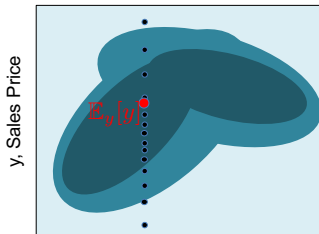
Where is the cross-term?

Take-home exercise: Show that it is zero

# Bias-Variance Trade-off: Illustration

$$\underbrace{\int_{\mathbf{x}} \int_y [\mathbb{E}_y[y|\mathbf{x}] - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy}_{\text{NOISE: error due to randomness of } y}$$

- For a given  $\mathbf{x}$ , we have a conditional distribution  $p(y|\mathbf{x})$ : the **Bayesian optimal prediction** of the label value for a given  $\mathbf{x}$  is  $\mathbb{E}_y[y|\mathbf{x}]$ ;
- The noise term measures the inherent variance in labels  $y$
- This term has nothing to do with our prediction  $h_{\mathcal{D}}(\mathbf{x})$





$$\underbrace{\int_{\mathbf{x}} \int_y [\mathbb{E}_y[y|\mathbf{x}] - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy}_{\text{NOISE: error due to randomness of } y}$$

## How can we reduce noise?

There is **nothing** we can do. This quantity depends on  $p(\mathbf{x}, y)$  only; choosing  $h(\cdot)$  or the training dataset  $\mathcal{D}$  will not affect it.

# Analyzing the Bias Term

$$\underbrace{\int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y|\mathbf{x}]]^2 p(\mathbf{x}, y) d\mathbf{x} dy}_{\text{BIAS}^2: \text{error due to the model approximation}}$$

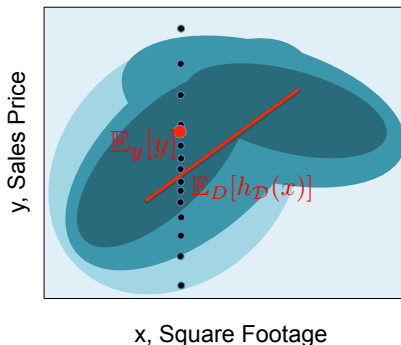
## Understanding the bias

- For each  $(\mathbf{x}, y)$  pair, we compute the loss of our averaged prediction  $\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})$  compared to the  $\mathbb{E}[y|\mathbf{x}]$ , the optimal prediction.
- If our model class was rich enough (eg. a high-degree polynomial), then  $\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})$  should perfectly match  $\mathbb{E}_y[y|\mathbf{x}]$
- Restricting to simpler models (eg. linear) results in a large bias, aka **underfitting!**

# Analyzing the Bias Term

$$\underbrace{\int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y|\mathbf{x}]]^2 p(\mathbf{x}, y) d\mathbf{x} dy}_{\text{BIAS}^2: \text{error due to the model approximation}}$$

Restricting to simpler models (eg. linear) results in a large bias



## How can we reduce the bias?

- It can be reduced by using more complex models. We shall choose  $h(\cdot)$  to be as flexible as possible: the better  $h(\cdot)$  approximates  $\mathbb{E}_y[y|\mathbf{x}]$ , the smaller the bias.
- However, this will increase the VARIANCE term.

# Summary of Risk Components

The average risk (with quadratic loss) can be decomposed as:

$$\begin{aligned}\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] &= \underbrace{\int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}}_{\text{VARIANCE: error due to training dataset}} \\ &+ \underbrace{\int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y|\mathbf{x}]]^2 p(\mathbf{x}, y) d\mathbf{x} dy}_{\text{BIAS}^2: \text{error due to the model approximation}} \\ &+ \underbrace{\int_{\mathbf{x}} \int_y [\mathbb{E}_y[y|\mathbf{x}] - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy}_{\text{NOISE: error due to randomness of } y}\end{aligned}$$

Here we define:  $h_{\mathcal{D}}(\mathbf{x})$  as the output of the model trained on  $\mathcal{D}$ ,  $\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})$  as the expectation of the model over all datasets  $\mathcal{D}$ , and  $\mathbb{E}_y[y|\mathbf{x}]$  as the expected value of  $y$ .

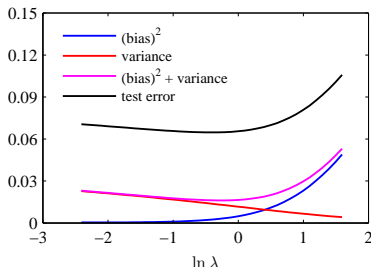
# Bias-Variance Tradeoff

Written compactly, the risk decomposition is

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \text{VARIANCE} + \text{BIAS}^2 + \text{NOISE}$$

where the first and the second term are inherently in conflict in terms of choosing what kind of  $h(\mathbf{x})$  we should use (unless we have an infinite amount of data).

If we can compute all terms analytically, they will look like this



# Summary for Today

- Overfitting and how to cope with it (more data, simpler model, regularization)
- Validation datasets (or cross-validation) are used to tune model hyperparameters.
- “Risk” framework for analyzing machine learning, and error decomposition into bias, variance, and noise terms.
  - Variance: Due to only optimizing over an empirical sample of the complete  $(x, y)$  distribution. High variance is responsible for **overfitting**.
  - Bias: Due to our choosing a model that does not fit the exact  $(x, y)$  relationship. High bias is responsible for **underfitting**.
  - Noise: Due to the output  $y$ 's randomness with respect to the input  $x$ .
- Choosing a more complex model improves the bias, but increases the variance (and vice versa for less complex models).
- The noise is independent of the model that we choose.