

# **18-661 Introduction to Machine Learning**

## Overfitting and the Bias-Variance Trade-off

---

Spring 2020

ECE – Carnegie Mellon University

# Announcements

- HW 2 will be released today, due February 10.
- Recitation on Friday will cover linear regression, gradient descent, and the bias-variance tradeoff (today's lecture). These problems will be helpful for HW 2.
- There has been some student demand for a Python and Jupyter tutorial.
  - SV will include this material at the end of Friday's recitation.
  - Pittsburgh will hold the tutorial tomorrow and broadcast to Rwanda around 1pm ET (rooms TBD).
- HW 1 solutions will be posted on Canvas later this week.

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

# Review of Ridge Regression

---

# What can go wrong with the LMS solution?

$$\mathbf{w}^{LMS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Why might  $\mathbf{X}^T \mathbf{X}$  be non-invertible?

- **Answer 1:**  $N < D$ . Not enough data to estimate all parameters.
- **Answer 2:** Columns of  $\mathbf{X}$  are not linearly independent, e.g., some features are linear functions of other features. In this case, solution is not unique. Examples:
  - A feature is a re-scaled version of another, for example, having two features correspond to length in meters and feet respectively
  - Same feature is repeated twice – could happen when there are many features
  - A feature has the same value for all data points
  - Sum of two features is equal to a third feature

## Example: Matrix $X^T X$ is not invertible

sqft (1000's)	bathrooms	sale price (100k)
1	2	2
2	2	3.5
1.5	2	3
2.5	2	4.5

Design matrix and target vector:

$$\mathbf{X} = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 2 & 2 \\ 1 & 1.5 & 2 \\ 1 & 2.5 & 2 \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 2 \\ 3.5 \\ 3 \\ 4.5 \end{bmatrix}$$

The 'bathrooms' feature is redundant, so we don't need  $w_2$

$$\begin{aligned} y &= w_0 + w_1 x_1 + w_2 x_2 \\ &= w_0 + w_1 x_1 + w_2 \times 2, \quad \text{since } x_2 \text{ is always } 2! \\ &= w_{0,eff} + w_1 x_1, \quad \text{where } w_{0,eff} = (w_0 + 2w_2) \end{aligned}$$

## Ridge regression

**Intuition:** what does a non-invertible  $\mathbf{X}^T \mathbf{X}$  mean?

Consider the SVD of this matrix:

$$\mathbf{X}^T \mathbf{X} = \mathbf{V} \begin{bmatrix} \lambda_1 & 0 & 0 & \cdots & 0 \\ 0 & \lambda_2 & 0 & \cdots & 0 \\ 0 & \cdots & \cdots & \cdots & 0 \\ 0 & \cdots & \cdots & \lambda_r & 0 \\ 0 & \cdots & \cdots & 0 & 0 \end{bmatrix} \mathbf{V}^T$$

where  $\lambda_1 \geq \lambda_2 \geq \cdots \lambda_r > 0$ ,  $r < D$ , and  $\mathbf{V}$  is a unitary matrix (its transpose is its inverse). We will need to divide by zero to compute  $(\mathbf{X}^T \mathbf{X})^{-1} \dots$

**Fix the problem:** ensure all singular values are non-zero:

$$\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} = \mathbf{V} \text{diag}(\lambda_1 + \lambda, \lambda_2 + \lambda, \dots, \lambda) \mathbf{V}^T$$

where  $\lambda > 0$  and  $\mathbf{I}$  is the identity matrix.

# Regularized least squares (ridge regression)

## Solution

$$\mathbf{w} = \left( \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

This is equivalent to adding an extra term to  $RSS(\mathbf{w})$

$$\overbrace{\frac{1}{2} \left\{ \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2 \left( \mathbf{X}^T \mathbf{y} \right)^T \mathbf{w} \right\}}^{RSS(\mathbf{w})} + \underbrace{\frac{1}{2} \lambda \|\mathbf{w}\|_2^2}_{\text{regularization}}$$

## Benefits

- Numerically more stable, invertible matrix
- Force  $\mathbf{w}$  to be small
- Prevent overfitting — more on this later



## Applying this to our example

sqft (1000's)	bathrooms	sale price (100k)
1	2	2
2	2	3.5
1.5	2	3
2.5	2	4.5

The 'bathrooms' feature is redundant, so we don't need  $w_2$

$$y = w_0 + w_1x_1 + w_2x_2$$

$$= w_0 + w_1x_1 + w_2 \times 2,$$

$$= w_{0,eff} + w_1x_1,$$

$$= 0.45 + 1.6x_1$$

since  $x_2$  is always 2!

where  $w_{0,eff} = (w_0 + 2w_2)$

Should get this

## Applying this to our example

The 'bathrooms' feature is redundant, so we don't need  $w_2$

$$\begin{aligned}y &= w_0 + w_1x_1 + w_2x_2 \\ &= w_0 + w_1x_1 + w_2 \times 2, \quad \text{since } x_2 \text{ is always } 2! \\ &= w_{0,eff} + w_1x_1, \quad \text{where } w_{0,eff} = (w_0 + 2w_2) \\ &= 0.45 + 1.6x_1 \quad \text{Should get this}\end{aligned}$$

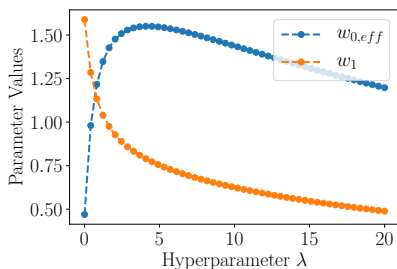
Compute the solution for  $\lambda = 0.5$

$$\begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = (\mathbf{X}^T \mathbf{X} + \lambda I)^{-1} \mathbf{X}^T \mathbf{y}$$
$$\begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 0.208 \\ 1.247 \\ 0.4166 \end{bmatrix}$$

## How does $\lambda$ affect the solution?

$$\begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

Let us plot  $w'_o = w_0 + 2w_2$  and  $w_1$  for different  $\lambda \in [0.01, 20]$

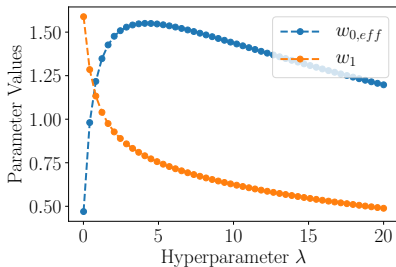


Setting small  $\lambda$  gives almost the least-squares solution ( $w'_o = 0.45$  and  $w_1 = 1.6$ ), but it can cause numerical instability in the inversion.

# How to choose $\lambda$ ?

$\lambda$  is referred as *hyperparameter*

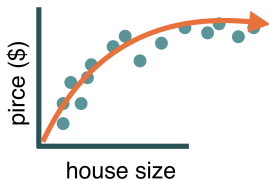
- Associated with the estimation method, not the dataset
- In contrast  $\mathbf{w}$  is the parameter vector
- Use validation set or cross-validation to find good choice of  $\lambda$



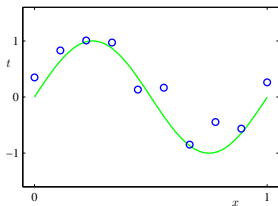
# Review of Non-linear Basis Functions

---

# Is a linear modeling assumption always a good idea?



**Figure 1:** Sale price can saturate as sq.footage increases



**Figure 2:** Temperature has cyclic variations over each year

# General nonlinear basis functions

We can use a nonlinear mapping:

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

- $M$  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  $\mathbf{w}^\top \phi(\mathbf{x})$
- other methods: nearest neighbors, decision trees, etc

## 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

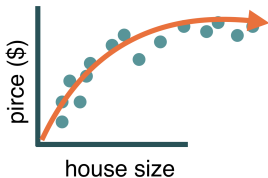
$$\Phi = \begin{pmatrix} \phi(\mathbf{x}_1)^\top \\ \phi(\mathbf{x}_2)^\top \\ \vdots \\ \phi(\mathbf{x}_N)^\top \end{pmatrix} \in \mathbb{R}^{N \times M}, \quad \mathbf{w}^{\text{LMS}} = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$$



## Example: Lots of flexibility in designing features!

$x_1$ , Area (1k sqft)	$\sqrt{x_1}$	Price (100k)
1	1	1
2.25	1.5	2
4	2	2.2
6.25	2.5	2.5

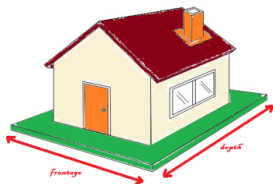
Price =  $\sqrt{x_1}$  is more accurate than Price =  $x_1$ .



**Figure 3:** Add  $\sqrt{x_1}$  as a feature to allow us to fit square-root, instead of linear, functions of the house area  $x_1$ .

## Example: Lots of flexibility in designing features!

$x_1$ , front (100ft)	$x_2$ depth (100ft)	$10x_1x_2$ , Lot (1k sqft)	Price (100k)
0.5	0.5	2.5	2
0.5	1	5	3.5
0.8	1.5	12	3
1.0	1.5	15	4.5



**Figure 4:** Instead of having frontage and depth as two separate features, it may be better to consider the lot-area, which is equal to frontage  $\times$  depth

# Overfitting and Regularization

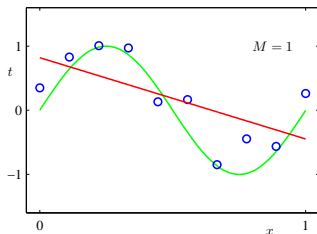
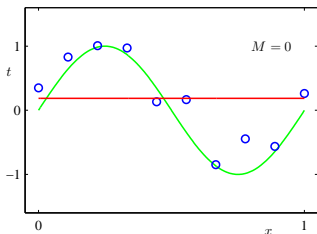
---

# Non-linear basis functions: Polynomial regression

## Polynomial basis functions

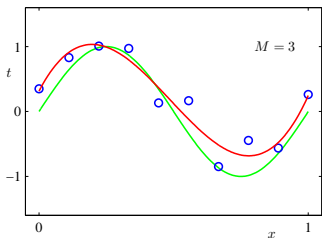
$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow f(x) = w_0 + \sum_{m=1}^M w_m x^m$$

Fitting samples from a sine function:

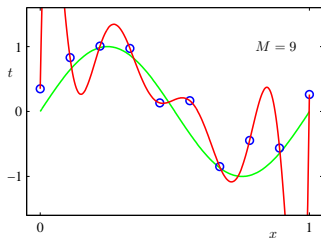


# Adding high-order terms

$M=3$



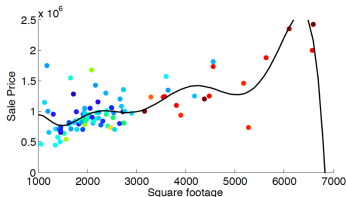
$M=9$ : **overfitting**



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

# Overfitting 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!

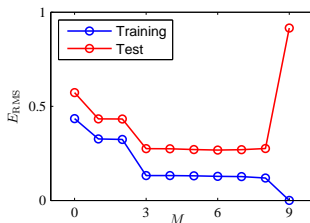
This is called **poor generalization/overfitting**.

# 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 and test dataset.

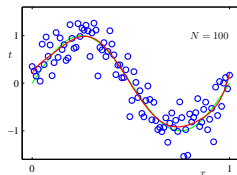
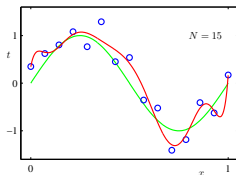
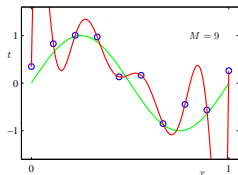


As a model increases in complexity:

- Training error keeps improving
- Test error may first improve but eventually will deteriorate

# Dealing with overfitting: Option 1

Try to use more training data

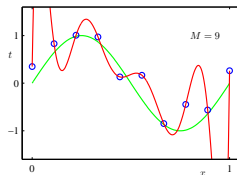
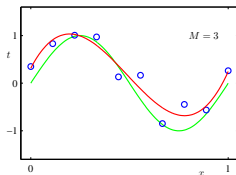
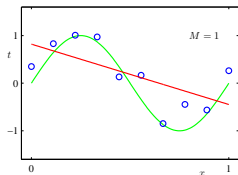


But getting a lot of data can be expensive and time-consuming



# Dealing with overfitting: Option 2

## Reduce the Number of Features



May not know which and how many features to remove

## Dealing with overfitting: Option 3

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

- How do we define a simple linear regression model —  $\mathbf{w}^\top \mathbf{x}$ ?
- Intuitively, the weights corresponding to higher order terms should not be “too large”

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
$w_0$	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

## Add a term to the objective function.

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

$$\frac{1}{2} \left\{ \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2 \left( \mathbf{X}^\top \mathbf{y} \right)^\top \mathbf{w} \right\} + \frac{1}{2} \lambda \|\mathbf{w}\|_2^2$$

Ridge regression is just regularized linear regression.

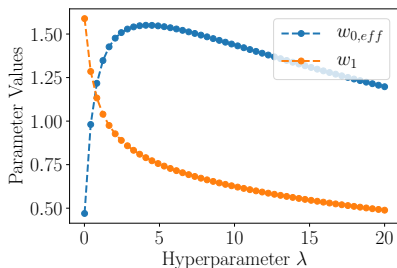
## Advantages

- Forces the magnitude of  $\mathbf{w}$  to be small
- Tries to find a simple model with few parameters
- Generalizes well to new data points

## Ridge regression as regularization

$$\begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = (\mathbf{X}^\top \mathbf{X} + \lambda I)^{-1} \mathbf{X}^\top \mathbf{y}$$

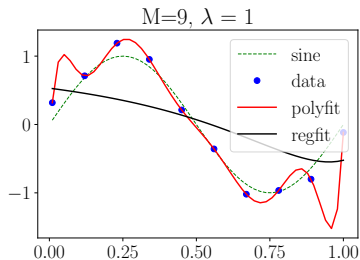
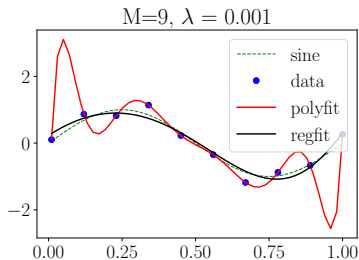
Let us plot  $w'_o = w_0 + 2w_2$  and  $w_1$  for different  $\lambda \in [0.01, 20]$



Setting small  $\lambda$  gives almost the least-squares solution, but it can cause numerical instability in the inversion

## Example: Effect of regularization

- Regularization makes the higher order  $w_i$ 's smaller
- Regularized polynomial fit will generalize much better
- As  $\lambda$  increases, the model becomes simpler



# Probabilistic interpretation of regularization

**Regularized Regression model:**  $Y = \mathbf{w}^\top \mathbf{X} + \eta$

- $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$  is a Gaussian random variable (as before)
- $w_d \sim N(0, \sigma^2)$  are i.i.d. Gaussian random variables (**unlike before**)
- We first choose the weight for each feature  $d$ ,  $w_d$ , from  $N(0, \sigma^2)$ . Then for each input vector  $\mathbf{x}_n$ , draw  $y_n$  from the distribution  $N(\mathbf{w}^\top \mathbf{x}_n, \sigma_0^2)$ .

**How do we estimate  $\mathbf{w}$  for this model?**

**Maximum a posterior (MAP)** estimate:

$$\begin{aligned}\mathbf{w}^{\text{MAP}} &= \arg \max_{\mathbf{w}} p(\mathbf{w} | \mathcal{D}) = \arg \max_{\mathbf{w}} \frac{p(\mathcal{D} | \mathbf{w}) p(\mathbf{w})}{p(\mathcal{D})} \\ &= \arg \max_{\mathbf{w}} p(\mathcal{D} | \mathbf{w}) p(\mathbf{w})\end{aligned}$$

## Estimating $\mathbf{w}$

Let  $\mathbf{x}_1, \dots, \mathbf{x}_N$  be i.i.d. with  $y|\mathbf{w}, \mathbf{x} \sim N(\mathbf{w}^\top \mathbf{x}, \sigma_0^2)$ ;  $w_d \sim N(0, \sigma^2)$ .

Given  $\sigma_0, \sigma$ , we choose  $\mathbf{w}$  so as to maximize:

$$p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) = \prod_n p(y_n|\mathbf{x}_n, \mathbf{w}) \prod_d p(w_d)$$

Now we know  $p(w_d) \propto \exp\left(\frac{-w_d^2}{2\sigma^2}\right)$  and  $p(y_n|\mathbf{x}_n, \mathbf{w}) \propto \exp\left(\frac{-(\mathbf{w}^\top \mathbf{x}_n - y_n)^2}{2\sigma_0^2}\right)$ :

$$\begin{aligned} \log p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) &= \sum_n \log p(y_n|\mathbf{x}_n, \mathbf{w}) + \sum_d \log p(w_d) \\ &= -\frac{\sum_n (\mathbf{w}^\top \mathbf{x}_n - y_n)^2}{2\sigma_0^2} - \sum_d \frac{1}{2\sigma^2} w_d^2 + \text{const} \end{aligned}$$

MAP estimate:  $\mathbf{w}^{\text{MAP}} = \arg \max_{\mathbf{w}} \log p(\mathcal{D}|\mathbf{w})p(\mathbf{w})$

$$\mathbf{w}^{\text{MAP}} = \operatorname{argmin}_{\mathbf{w}} \frac{\sum_n (\mathbf{w}^\top \mathbf{x}_n - y_n)^2}{2\sigma_0^2} + \frac{1}{2\sigma^2} \|\mathbf{w}\|_2^2$$

## Maximum a posterior (MAP) estimate

$$\mathcal{E}(\mathbf{w}) = \sum_n (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \lambda \|\mathbf{w}\|_2^2$$

where  $\lambda > 0$  is used to denote  $\sigma_0^2/\sigma^2$ . This extra term  $\|\mathbf{w}\|_2^2$  is called regularization/regularizer and controls the magnitude of  $\mathbf{w}$ .

- If  $\lambda \rightarrow +\infty$ , then  $\sigma_0^2 \gg \sigma^2$ : the variance of noise is far greater than what our prior model can allow for  $\mathbf{w}$ . In this case, our prior model on  $\mathbf{w}$  will give a simpler model. Numerically,

$$\mathbf{w}^{\text{MAP}} \rightarrow \mathbf{0}$$

- If  $\lambda \rightarrow 0$ , then we trust our data more. Numerically,

$$\mathbf{w}^{\text{MAP}} \rightarrow \mathbf{w}^{\text{LMS}} = \operatorname{argmin} \sum_n (\mathbf{w}^\top \mathbf{x}_n - y_n)^2$$



# Hyperparameter Tuning and Cross-Validation

---

# How should we choose the right amount of regularization?

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

**No:** as this will always set  $\lambda$  to zero, i.e., no regularization, defeating our intention of controlling model complexity

$\lambda$  is thus a hyperparameter. To tune it,

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

Similar idea applies to tuning learning rate  $\eta$  (or any other hyperparameter) as well.

# Tuning by using a validation dataset

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

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

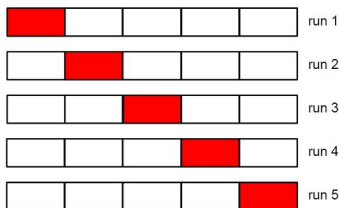
- For each possible value of the hyperparameter (say  $\lambda = 1, 3, \dots, 100$ )
  - Train a model using  $\mathcal{D}^{\text{TRAIN}}$
  - Evaluate the performance of the model on  $\mathcal{D}^{\text{VAL}}$
- Choose the model with the best performance on  $\mathcal{D}^{\text{VAL}}$
- Evaluate this 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 (based on average, variance, etc.)

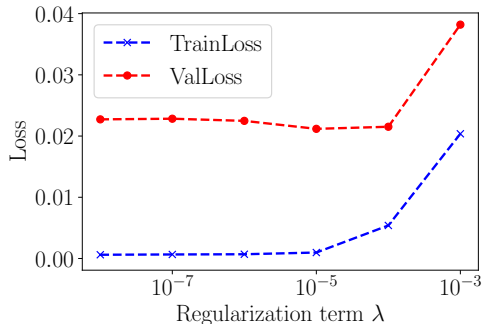
**Special case:** when  $S = N$ , this will be leave-one-out.

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



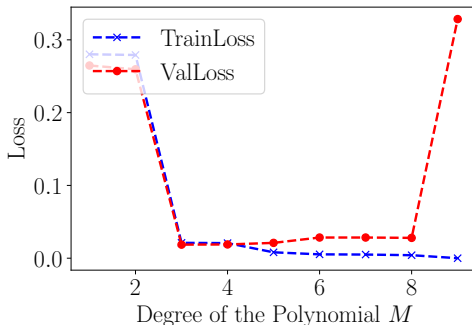
## Example: Hyper-parameter Tuning $\lambda$

- $\lambda = 10^{-4}$  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**

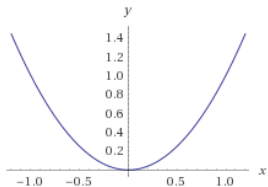
$$\ell(h(\mathbf{x}), y)$$

### Example:

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

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

Ex: when  $y = 0$



# 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  $N \rightarrow +\infty$ ,

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

# How could this go wrong?

So far, we have been doing **empirical risk minimization (ERM)**

For linear regression,  $h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$ , and we use squared loss  $\ell$ .

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

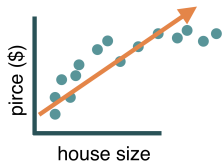
$$R[h(\mathbf{x})] = \int_{\mathbf{x}, y} \ell(h(\mathbf{x}), y) p(\mathbf{x}, y) d\mathbf{x} dy$$

**What could go wrong with ERM?**

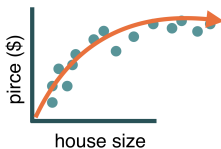
- **Limited Function Class:** The function  $h(\mathbf{x})$  is restricted to a limited class (e.g. linear functions), which does not allow us to perfectly fit  $y$ , even if we had infinitely many training data points.
- **Limited Data:** We don't know  $p(\mathbf{x}, y)$ , so we must hope that we have enough training data that the empirical risk approximates the real risk. Otherwise, we will **overfit** to the training data.

# Bias-Variance Trade-off: Intuition

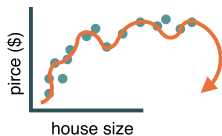
- **High Bias:** Model is not rich enough to fit the training dataset and achieve low training loss
- **High Variance:** If the training dataset changes slightly, the model changes a lot
- Regularization helps find a middle ground



**Figure 6:** High Bias



**Figure 7:** Just Right



**Figure 8:** High Variance

## Goal: to understand the sources of prediction errors

- $\mathcal{D}$ : our training data
- $h_{\mathcal{D}}(\mathbf{x})$ : our prediction function  
We are using the subscript  $\mathcal{D}$  to indicate that the prediction function is learned on the specific set of training data  $\mathcal{D}$
- $\ell(h(\mathbf{x}), y)$ : our square loss function for regression

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

- Unknown joint distribution  $p(\mathbf{x}, y)$

## 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

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

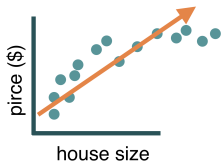
We will now evaluate the **expected risk**  $\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})]$ : the average risk over the distribution of possible training datasets,  $P(\mathcal{D})$ .

# Bias-Variance Trade-off: Intuition

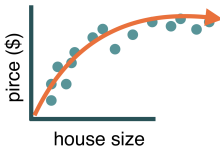
Error decomposes into 3 terms

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

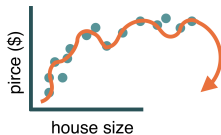
We will prove this result, and interpret what it means...



**Figure 9:** High Bias



**Figure 10:** Just Right



**Figure 11:** High Variance

# 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.



We will subtract the averaged prediction from the averaged risk

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

## Where does the cross-term go?

It is zero

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

## Understanding the variance

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

For each  $(\mathbf{x}, \mathbf{y})$  pair, we compute the squared difference of  $h_{\mathcal{D}}(\mathbf{x})$  (the prediction with training dataset  $\mathcal{D}$ ) and the averaged prediction  $\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})$ : the average (over all  $(\mathbf{x}, \mathbf{y}) \sim p$ ) **variance of the prediction** over  $\mathcal{D}$ .

## 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  $h(\mathbf{x}) = \text{const}$ .

## The remaining item

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

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

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

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

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

## 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

# Analyzing the noise

How can we reduce noise?

$$\int_{\mathbf{x}} \int_y [\mathbb{E}_y[y|\mathbf{x}] - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy = \int_{\mathbf{x}} \left( \int_y [\mathbb{E}_y[y|\mathbf{x}] - y]^2 p(y|\mathbf{x}) dy \right) p(\mathbf{x}) d\mathbf{x}$$

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. Note that the integral inside the parentheses is the *variance* (noise) of the posterior distribution  $p(y|\mathbf{x})$  at the given  $\mathbf{x}$ .

Figure 12: Somewhat difficult posterior

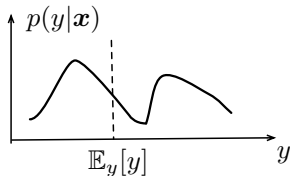
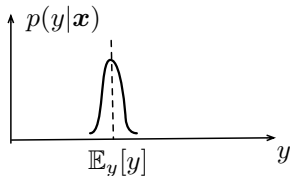


Figure 13: Somewhat easy posterior



## Understanding the bias

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

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 expected value of  $y$  given  $\mathbf{x}$ , which we compute as  $\mathbb{E}_y[y|\mathbf{x}] = \int_y yp(y|\mathbf{x})dy$ . Then we take the average over all pairs  $(\mathbf{x}, y) \sim p(\mathbf{x}, y)$ .

### 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.

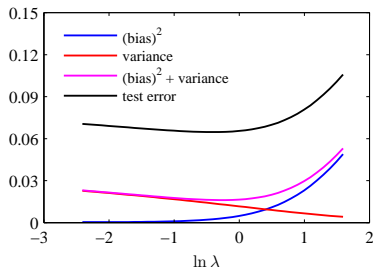
# Bias/variance tradeoff

## Error decomposes into 3 terms

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

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 of risk components

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

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

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  $y$ .

## Example: Why regularized linear regression could be helpful?

### Model

$$h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$$

Consider the best possible (linear)  $h^*(\mathbf{x})$

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \int_{\mathbf{x}} [\mathbb{E}_y[y|\mathbf{x}] - \mathbf{w}^\top \mathbf{x}]^2 p(\mathbf{x}) d\mathbf{x}$$

Note that this linear model assumes the knowledge of joint distribution, thus, not achievable. Intuitively, it is the *best* linear model that can predict the data most accurately.

## More refined decomposition of the bias

$$\int_{\mathbf{x}} [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y|\mathbf{x}]]^2 p(\mathbf{x}) d\mathbf{x} = \int_{\mathbf{x}} [h^*(\mathbf{x}) - \mathbb{E}_y[y|\mathbf{x}]]^2 p(\mathbf{x}) d\mathbf{x} \\ + \int_{\mathbf{x}} [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - h^*(\mathbf{x})]^2 p(\mathbf{x}) d\mathbf{x}$$

- **Model bias:** the price we pay for choosing linear functions to model data. This is the difference between the prediction of the best possible linear model and the actual target.
- **Estimation bias:** the difference between the optimal model and the estimated model.

*Normally, the estimation bias is zero if we do not regularize.*

# Bias/variance tradeoff for regularized linear regression

We can only adjust estimation bias

$$\int_{\mathbf{x}} [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}; \lambda) - h^*(\mathbf{x})]^2 p(\mathbf{x}) d\mathbf{x}$$

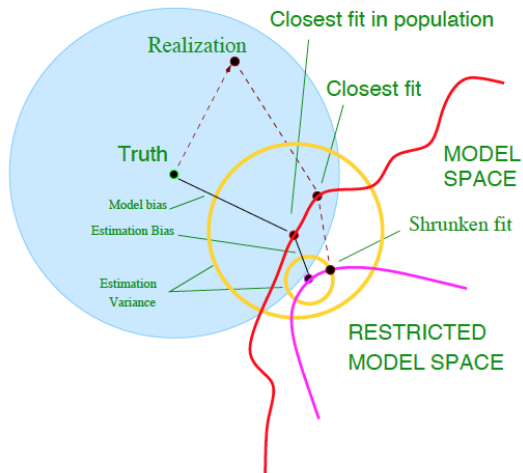
where  $h(\mathbf{x}; \lambda)$  is the estimated model with regularized linear regression (parameterized with  $\lambda$ ).

This term will not be zero anymore!

Thus, bias goes up.

But, as long as this is balanced with a decrease in variance, we are willing to do so.

# Visualizing the tradeoff



# Lecture Summary

- Validation datasets (or cross-validation) are used to determine model hyperparameters.
- Many ML models use empirical risk minimization to find the optimal parameters.
- ERM leads to an error consisting of bias, variance, and noise terms.
  - Variance: Due to only optimizing over an empirical sample of the complete  $(x, y)$  distribution.
  - Bias: Due to our choosing a model that does not fit the exact  $(x, y)$  relationship.
  - 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.