

# Mathematical Foundations of ML

Philipp Grohs



Bedlewo, Nov. 2019

## Short Reading List

- 1 Felipe Cucker and Ding Yuan Zhou: Learning Theory: An Approximation Theory Viewpoint, 2001
- 2 Luc Devroye, Laszlo Györfi, Gabor Lugosi: A Probabilistic Theory of Pattern Recognition; Springer, 2013.
- 3 Aurelien Geron: Hands-On Machine Learning with Scikit-Learn and TensorFlow; O'Reilly, 2017
- 4 Brian Steele and John Chandler and Swarna Reddy: Algorithms for Data Science; Springer, 2017

# Syllabus

- 1 Basic Concepts
- 2 Mathematical Foundations of General Regression Problems

# 1. Mathematical Foundations of Machine Learning

## 1.1 Basic Concepts

# Definition of Learning

## Definition [Mitchell (1997)]

“A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ ”

# The Task $T$

## Classification

Compute  $f : \mathbb{R}^n \rightarrow \{1, \dots, k\}$  which maps data  $x \in \mathbb{R}^n$  to a category in  $\{1, \dots, k\}$ . Alternative: Compute  $f : \mathbb{R}^n \rightarrow \mathbb{R}^k$  which maps data  $x \in \mathbb{R}^n$  to a histogram with respect to  $k$  categories.

# The Task $T$

## Classification

Compute  $f : \mathbb{R}^n \rightarrow \{1, \dots, k\}$  which maps data  $x \in \mathbb{R}^n$  to a category in  $\{1, \dots, k\}$ . Alternative: Compute  $f : \mathbb{R}^n \rightarrow \mathbb{R}^k$  which maps data  $x \in \mathbb{R}^n$  to a histogram with respect to  $k$  categories.





# The Task $T$

## Classification

Compute  $f : \mathbb{R}^n \rightarrow \{1, \dots, k\}$  which maps data  $x \in \mathbb{R}^n$  to a category in  $\{1, \dots, k\}$ . Alternative: Compute  $f : \mathbb{R}^n \rightarrow \mathbb{R}^k$  which maps data  $x \in \mathbb{R}^n$  to a histogram with respect to  $k$  categories.



$$x = \text{[image of handwritten digit 5]} \mapsto f(x) = 5.$$

# The Task $T$

## Regression

Predict a numerical value  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ .

# The Task $T$

## Regression

Predict a numerical value  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ .

- Expected claim of insured person

# The Task $T$

## Regression

Predict a numerical value  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ .

- Expected claim of insured person
- Algorithmic trading

# The Task $T$

## Density Estimation

Estimate a probability density  $p : \mathbb{R}^n \rightarrow \mathbb{R}_+$  which can be interpreted as a probability distribution on the space that the examples were drawn from.

# The Task $T$

## Density Estimation

Estimate a probability density  $p : \mathbb{R}^n \rightarrow \mathbb{R}_+$  which can be interpreted as a probability distribution on the space that the examples were drawn from.

- Useful for many tasks in data processing, for example if we observe corrupted data  $\tilde{x}$  we may estimate the original  $x$  as the argmax of  $p(\tilde{x}|x)$ .

## The Experience $E$

The experience typically consists of a dataset which consists of many examples (aka data points).

# The Experience $E$

The experience typically consists of a dataset which consists of many examples (aka data points).

- If these data points are labeled (for example in the classification problem, if we know the classifier of our given data points) we speak of *supervised learning*.



# The Experience $E$

The experience typically consists of a dataset which consists of many examples (aka data points).

- If these data points are labeled (for example in the classification problem, if we know the classifier of our given data points) we speak of *supervised learning*.
- If these data points are not labeled (for example in the classification problem, the algorithm would have to find the clusters itself from the given dataset) we speak of *unsupervised learning*.

# The Performance Measure $P$

In classification problems this is typically the *accuracy*, i.e., the proportion of examples for which the model produces the correct output.

# The Performance Measure $P$

In classification problems this is typically the *accuracy*, i.e., the proportion of examples for which the model produces the correct output.

- Often the given dataset is split into a *training set* on which the algorithm operates and a *test set* on which its performance is measured.

## An Example: Linear Regression

# An Example: Linear Regression

## The Task

Regression: Predict  $\hat{f} : \mathbb{R}^d \rightarrow \mathbb{R}$ .

# An Example: Linear Regression

## The Task

Regression: Predict  $\hat{f} : \mathbb{R}^d \rightarrow \mathbb{R}$ .

## The Experience

Training data  $((x_i^{train}, y_i^{train}))_{i=1}^m$  with  $y_i^{train} \sim \hat{f}(x_i^{train})$

# An Example: Linear Regression

## The Task

Regression: Predict  $\hat{f} : \mathbb{R}^d \rightarrow \mathbb{R}$ .

## The Experience

Training data  $((x_i^{train}, y_i^{train}))_{i=1}^m$  with  $y_i^{train} \sim \hat{f}(x_i^{train})$

## The Performance Measure

Given test data  $((x_i^{test}, y_i^{test}))_{i=1}^n$  we evaluate the performance of an estimator  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  as the *mean squared error*

$$\frac{1}{n} \sum_{i=1}^n |f(x_i^{test}) - y_i^{test}|^2.$$

# An Example: Linear Regression

## The Computer Program

Define a *Hypothesis Space*

$$\mathcal{H} = \text{span}\{\varphi_1, \dots, \varphi_l\} \subset C(\mathbb{R}^d)$$



# An Example: Linear Regression

## The Computer Program

Define a *Hypothesis Space*

$$\mathcal{H} = \text{span}\{\varphi_1, \dots, \varphi_l\} \subset C(\mathbb{R}^d)$$

and, given training data

$$\mathbf{z} = (x_i, y_i)_{i=1}^m,$$

# An Example: Linear Regression

## The Computer Program

Define a *Hypothesis Space*

$$\mathcal{H} = \text{span}\{\varphi_1, \dots, \varphi_l\} \subset C(\mathbb{R}^d)$$

and, given training data

$$\mathbf{z} = (x_i, y_i)_{i=1}^m,$$

define the *empirical risk*

$$\mathcal{E}_{\mathbf{z}}(f) := \frac{1}{m} \sum_{i=1}^m (f(x_i) - y_i)^2.$$

# An Example: Linear Regression

## The Computer Program

Define a *Hypothesis Space*

$$\mathcal{H} = \text{span}\{\varphi_1, \dots, \varphi_l\} \subset C(\mathbb{R}^d)$$

and, given training data

$$\mathbf{z} = (x_i, y_i)_{i=1}^m,$$

define the *empirical risk*

$$\mathcal{E}_{\mathbf{z}}(f) := \frac{1}{m} \sum_{i=1}^m (f(x_i) - y_i)^2.$$

We let our algorithm find the minimizer (a.k.a. *empirical regression function*)

$$\hat{f}_{\mathcal{H}, \mathbf{z}} := \operatorname{argmin}_{f \in \mathcal{H}} \mathcal{E}_{\mathbf{z}}(f).$$

# An Example: Linear Regression

Computing the Empirical Target Function

# An Example: Linear Regression

## Computing the Empirical Target Function

■ Let

$$\mathbf{A} = (\varphi_j(x_i))_{i,j} \in \mathbb{R}^{m \times l}.$$

# An Example: Linear Regression

## Computing the Empirical Target Function

- Let

$$\mathbf{A} = (\varphi_j(x_i))_{i,j} \in \mathbb{R}^{m \times l}.$$

- Every  $f \in \mathcal{H}$  can be written as  $\sum_{i=1}^m w_i \varphi_i$  and we denote  $\mathbf{w} := (w_i)_{i=1}^l$ .

# An Example: Linear Regression

## Computing the Empirical Target Function

- Let

$$\mathbf{A} = (\varphi_j(x_i))_{i,j} \in \mathbb{R}^{m \times l}.$$

- Every  $f \in \mathcal{H}$  can be written as  $\sum_{i=1}^m w_i \varphi_i$  and we denote  $\mathbf{w} := (w_i)_{i=1}^l$ .
- We let  $\mathbf{y} := (y_i)_{i=1}^m$ .

# An Example: Linear Regression

## Computing the Empirical Target Function

- Let

$$\mathbf{A} = (\varphi_j(x_i))_{i,j} \in \mathbb{R}^{m \times l}.$$

- Every  $f \in \mathcal{H}$  can be written as  $\sum_{i=1}^m w_i \varphi_i$  and we denote  $\mathbf{w} := (w_i)_{i=1}^l$ .
- We let  $\mathbf{y} := (y_i)_{i=1}^m$ .
- We get that

$$\mathcal{E}_{\mathbf{z}}(f) = \|\mathbf{A}\mathbf{w} - \mathbf{y}\|^2.$$



# An Example: Linear Regression

## Computing the Empirical Target Function

- Let

$$\mathbf{A} = (\varphi_j(x_i))_{i,j} \in \mathbb{R}^{m \times l}.$$

- Every  $f \in \mathcal{H}$  can be written as  $\sum_{i=1}^m w_i \varphi_i$  and we denote  $\mathbf{w} := (w_i)_{i=1}^l$ .
- We let  $\mathbf{y} := (y_i)_{i=1}^m$ .
- We get that

$$\mathcal{E}_{\mathbf{z}}(f) = \|\mathbf{A}\mathbf{w} - \mathbf{y}\|^2.$$

- A minimizer is given by  $\mathbf{w}_* := \mathbf{A}^\dagger \mathbf{y}$ , and we get our estimate

$$f_* := \sum_{i=1}^l (\mathbf{w}_*)_i \varphi_i.$$

Proof.

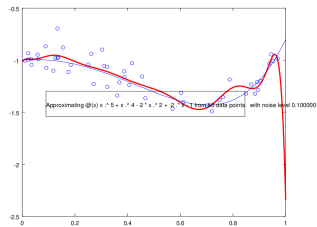
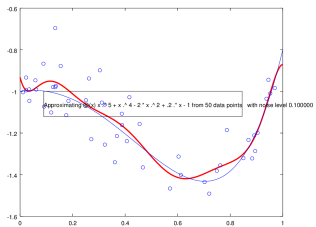
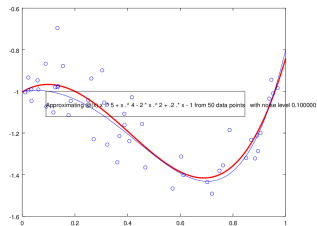
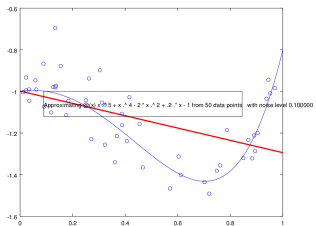
We want to minimize the function

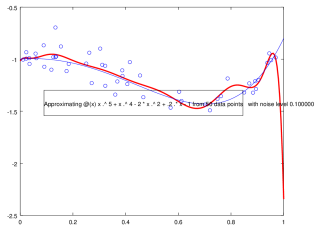
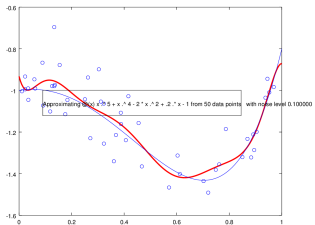
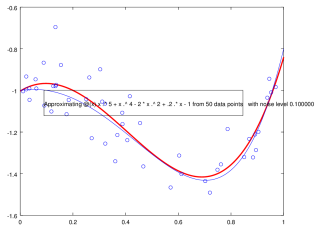
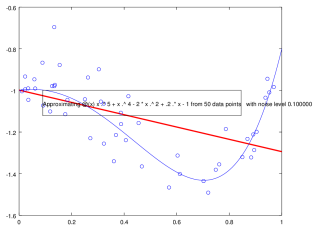
$$\mathcal{X}(\mathbf{w}) := \mathbf{w} \mapsto \|\mathbf{A}\mathbf{w} - \mathbf{y}\|^2,$$

which is (more or less...) equivalent to setting its first derivative to zero. It holds that

$$\frac{d\mathcal{X}(\mathbf{w})}{d\mathbf{w}} = 2\mathbf{A}^\dagger(\mathbf{A}\mathbf{w} - \mathbf{y}),$$

which, if set to zero, are precisely the normal equations. □





Degree too low: underfitting. Degree to high: overfitting!

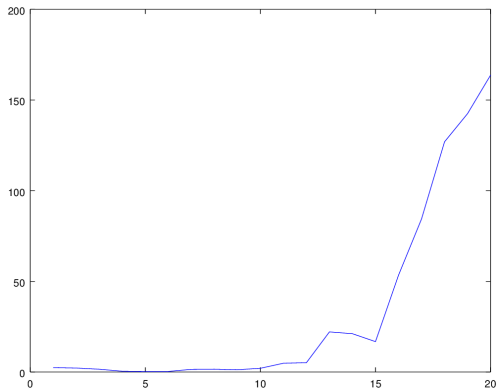


Figure: Error with Polynomial Degree

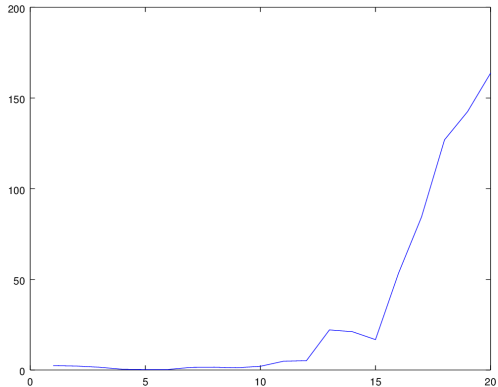


Figure: Error with Polynomial Degree

### Bias-Variance Problem

“Capacity” of the hypothesis space has to be adapted to the complexity of the target function and the sample size!

## 1.2 Mathematical Foundations of General Regression Problems

## 1.2.1 Basic Definitions



# The Mathematical Learning Problem

## The Mathematical Learning Problem

Let  $(\Sigma, \mathcal{G}, \mathbb{P})$  probability space. Given (Borel measurable) random vectors  $X : \Sigma \rightarrow \mathbb{R}^d$ ,  $Y : \Sigma \rightarrow \mathbb{R}^k$  with  $\text{im}(X) \subseteq \Omega$  for  $\Omega \subset \mathbb{R}^d$  compact.

## The Mathematical Learning Problem

Let  $(\Sigma, \mathcal{G}, \mathbb{P})$  probability space. Given (Borel measurable) random vectors  $X : \Sigma \rightarrow \mathbb{R}^d$ ,  $Y : \Sigma \rightarrow \mathbb{R}^k$  with  $\text{im}(X) \subseteq \Omega$  for  $\Omega \subset \mathbb{R}^d$  compact.

For any (Borel measurable)  $f : \Omega \rightarrow \mathbb{R}^k$  define the *least squares error*

$$\mathcal{E}(f) := \mathbb{E}[(f(X) - Y)^2].$$

## The Mathematical Learning Problem

Let  $(\Sigma, \mathcal{G}, \mathbb{P})$  probability space. Given (Borel measurable) random vectors  $X : \Sigma \rightarrow \mathbb{R}^d$ ,  $Y : \Sigma \rightarrow \mathbb{R}^k$  with  $\text{im}(X) \subseteq \Omega$  for  $\Omega \subset \mathbb{R}^d$  compact.

For any (Borel measurable)  $f : \Omega \rightarrow \mathbb{R}^k$  define the *least squares error*

$$\mathcal{E}(f) := \mathbb{E}[(f(X) - Y)^2].$$

The learning problem asks for the function  $\hat{f}$  which minimizes  $\mathcal{E}$ .

## The Mathematical Learning Problem

Let  $(\Sigma, \mathcal{G}, \mathbb{P})$  probability space. Given (Borel measurable) random vectors  $X : \Sigma \rightarrow \mathbb{R}^d$ ,  $Y : \Sigma \rightarrow \mathbb{R}^k$  with  $\text{im}(X) \subseteq \Omega$  for  $\Omega \subset \mathbb{R}^d$  compact.

For any (Borel measurable)  $f : \Omega \rightarrow \mathbb{R}^k$  define the *least squares error*

$$\mathcal{E}(f) := \mathbb{E}[(f(X) - Y)^2].$$

The learning problem asks for the function  $\hat{f}$  which minimizes  $\mathcal{E}$ .



## Example 1: Regression

## Example 1: Regression

- Suppose our data is generated from noisy observations

$$y = f(x) + \xi,$$

where  $\xi$  is a r.v. with  $\mathbb{E}(\xi) = 0$ .

## Example 1: Regression

- Suppose our data is generated from noisy observations

$$y = f(x) + \xi,$$

where  $\xi$  is a r.v. with  $\mathbb{E}(\xi) = 0$ .

- Let  $X : \Omega \rightarrow \mathbb{R}$  be a r.v. (independent of  $\xi$ ) and let  $Y := f(X) + \xi$ .



## Example 1: Regression

- Suppose our data is generated from noisy observations

$$y = f(x) + \xi,$$

where  $\xi$  is a r.v. with  $\mathbb{E}(\xi) = 0$ .

- Let  $X : \Omega \rightarrow \mathbb{R}$  be a r.v. (independent of  $\xi$ ) and let  $Y := f(X) + \xi$ .
- We have that

$$\begin{aligned}\mathcal{E}(g) &= \mathbb{E}[(g(X) - Y)^2] = \mathbb{E}[(g(X) - f(X) - \xi)^2] \\ &= \mathbb{E}[(f(X) - g(X))^2] + 2\mathbb{E}[(g(X) - f(X))\xi] + \mathbb{E}\xi^2 \\ &= \mathbb{E}[(f(X) - g(X))^2] + \mathbb{E}\xi^2 \\ &= \|f - g\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2 + \mathbb{V}[\xi].\end{aligned}$$

## Example 1: Regression

- Suppose our data is generated from noisy observations

$$y = f(x) + \xi,$$

where  $\xi$  is a r.v. with  $\mathbb{E}(\xi) = 0$ .

- Let  $X : \Omega \rightarrow \mathbb{R}$  be a r.v. (independent of  $\xi$ ) and let  $Y := f(X) + \xi$ .
- We have that

$$\begin{aligned}\mathcal{E}(g) &= \mathbb{E}[(g(X) - Y)^2] = \mathbb{E}[(g(X) - f(X) - \xi)^2] \\ &= \mathbb{E}[(f(X) - g(X))^2] + 2\mathbb{E}[(g(X) - f(X))\xi] + \mathbb{E}\xi^2 \\ &= \mathbb{E}[(f(X) - g(X))^2] + \mathbb{E}\xi^2 \\ &= \|f - g\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2 + \mathbb{V}[\xi].\end{aligned}$$

- The learning problem finds  $f$ !

## Example 2: Classifications

## Example 2: Classifications

- Suppose that there is a function  $f$  which maps a matrix  $x \in [0, 1]^{256 \times 256}$  to a histogram  $f(x) \in \mathbb{R}_+^{10}$ . We consider the vector  $f(x) / \sum_{i=1}^{10} f(x)_i$  as a histogram describing which digit the image  $x$  represents.

## Example 2: Classifications

- Suppose that there is a function  $f$  which maps a matrix  $x \in [0, 1]^{256 \times 256}$  to a histogram  $f(x) \in \mathbb{R}_+^{10}$ . We consider the vector  $f(x) / \sum_{i=1}^{10} f(x)_i$  as a histogram describing which digit the image  $x$  represents.
- Let  $(X, Y)$  be random vectors on  $\mathbb{R}^{256 \times 256} \times \mathbb{R}_+^{10}$  which generate the measurement data we get to see ( $(X, Y)$  will not be known to us!!!)

## Example 2: Classifications

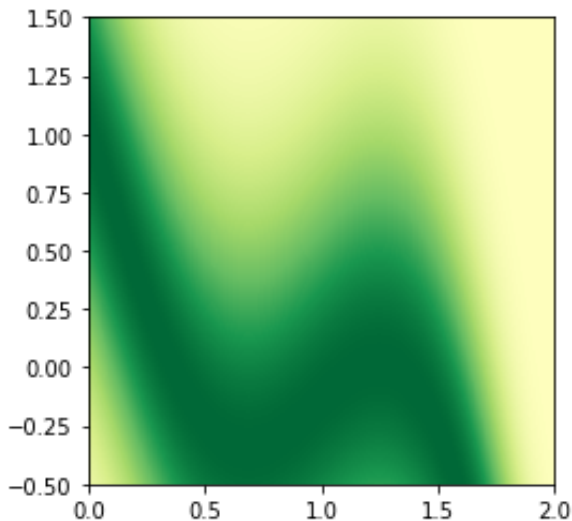
- Suppose that there is a function  $f$  which maps a matrix  $x \in [0, 1]^{256 \times 256}$  to a histogram  $f(x) \in \mathbb{R}_+^{10}$ . We consider the vector  $f(x) / \sum_{i=1}^{10} f(x)_i$  as a histogram describing which digit the image  $x$  represents.
- Let  $(X, Y)$  be random vectors on  $\mathbb{R}^{256 \times 256} \times \mathbb{R}_+^{10}$  which generate the measurement data we get to see ( $(X, Y)$  will not be known to us!!!)
- Now, a function  $f$  as above will in general not exist for our problem. But we can look for the function  $\hat{f}$  which minimizes the least squares error  $\mathcal{E}$  – this will be the optimal explanation of the measurements in terms of a functional relation between  $X$  and  $Y$ !

## A New Look

Suppose that our training data consists of samples according to a given data distribution  $(X, Y)$

## A New Look

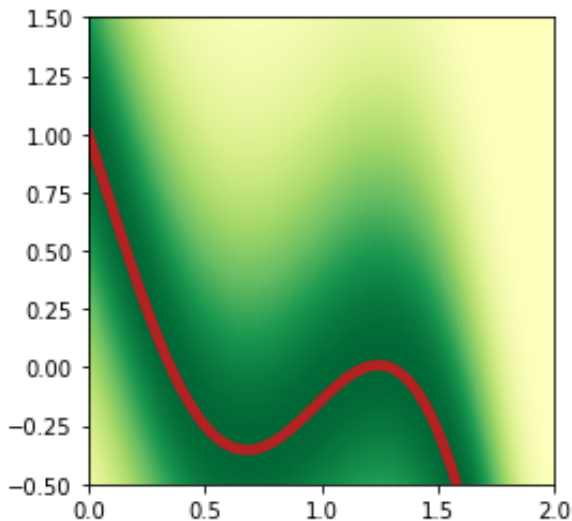
Suppose that our training data consists of samples according to a given data distribution  $(X, Y)$





## A New Look

If we knew the data distribution  $(X, Y)$ , the best functional relation between  $X$  and  $Y$  would simply be  $\mathbb{E}[Y|X = x]$ !

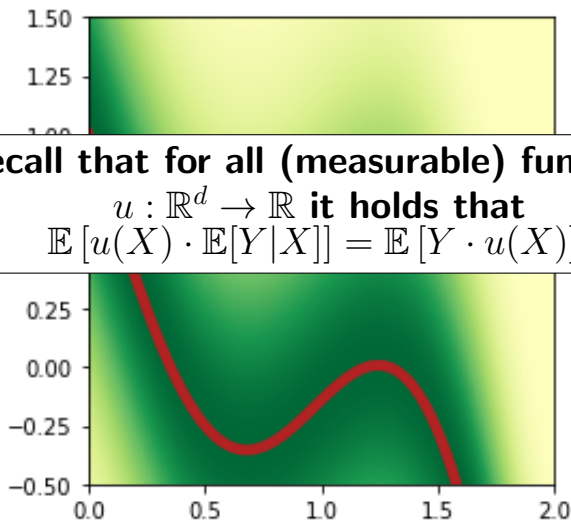


## A New Look

If we knew the data distribution  $(X, Y)$ , the best functional relation between  $X$  and  $Y$  would simply be  $\mathbb{E}[Y|X = x]$ !

**Recall that for all (measurable) functions**

$$u : \mathbb{R}^d \rightarrow \mathbb{R} \text{ it holds that}$$
$$\mathbb{E} [u(X) \cdot \mathbb{E}[Y|X]] = \mathbb{E} [Y \cdot u(X)].$$



# Regression Function

## Theorem (Main Regression Theorem)

Let  $\hat{f} := \mathbb{E}[Y|X]$  be the regression function and  $\sigma^2 := \mathcal{E}(\hat{f})$ . It holds that

$$\mathcal{E}(f) = \|f - \hat{f}\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2 + \sigma^2$$

# Regression Function

## Theorem (Main Regression Theorem)

Let  $\hat{f} := \mathbb{E}[Y|X]$  be the regression function and  $\sigma^2 := \mathcal{E}(\hat{f})$ . It holds that

$$\mathcal{E}(f) = \|f - \hat{f}\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2 + \sigma^2$$

Proof.

$$\begin{aligned}\mathcal{E}(f) &= \mathbb{E}[(f(X) - Y)^2] = \mathbb{E}[(\hat{f}(X) - Y)^2] + \\ &\quad 2 \underbrace{\mathbb{E}[(f(X) - \hat{f}(X)) \cdot (\hat{f}(X) - Y)]}_{=0} + \mathbb{E}[(f(X) - \hat{f}(X))^2].\end{aligned}$$



# Regression Function

## Theorem (Main Regression Theorem)

Let  $\hat{f} := \mathbb{E}[Y|X]$  be the regression function and  $\sigma^2 := \mathcal{E}(\hat{f})$ . It holds that

$$\mathcal{E}(f) = \|f - \hat{f}\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2 + \sigma^2$$

## Proof.

$$\begin{aligned}\mathcal{E}(f) &= \mathbb{E}[(f(X) - Y)^2] = \mathbb{E}[(\hat{f}(X) - Y)^2] + \\ &\quad 2 \underbrace{\mathbb{E}[(f(X) - \hat{f}(X)) \cdot (\hat{f}(X) - Y)]}_{=0} + \mathbb{E}[(f(X) - \hat{f}(X))^2].\end{aligned}$$



## Corollary

The regression function solves the learning problem!

# Regression Function

## Theorem (Main Regression Theorem)

Let  $\hat{f} := \mathbb{E}[Y|X]$  be the regression function and  $\sigma^2 := \mathcal{E}(\hat{f})$ . It holds that

$$\mathcal{E}(f) = \|f - \hat{f}\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2 + \sigma^2$$

Proof.

**Are we done?**

$$\begin{aligned}\mathcal{E}(f) &= \mathbb{E}[(f(X) - Y)^2] = \mathbb{E}[(\hat{f}(X) - Y)^2] + \\ &\quad \underbrace{2\mathbb{E}[(f(X) - \hat{f}(X)) \cdot (\hat{f}(X) - Y)]}_{=0} + \mathbb{E}[(f(X) - \hat{f}(X))^2].\end{aligned}$$



## Corollary

The regression function solves the learning problem!

# Regression Function

## Theorem (Main Regression Theorem)

Let  $\hat{f} := \mathbb{E}[Y|X]$  be the regression function and  $\sigma^2 := \mathcal{E}(\hat{f})$ . It holds that

$$\mathcal{E}(f) = \|f - \hat{f}\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2 + \sigma^2$$

Proof.

Are we done?

$$\mathcal{E}(f) = \mathbb{E}[(\text{😞 We don't know } (X, Y)!!! \\ \underbrace{2\mathbb{E}[(f(X) - \hat{f}(X)) \cdot (\hat{f}(X) - Y)]}_{=0} + \mathbb{E}[(f(X) - \hat{f}(X))^2].$$



## Corollary

The regression function solves the learning problem!

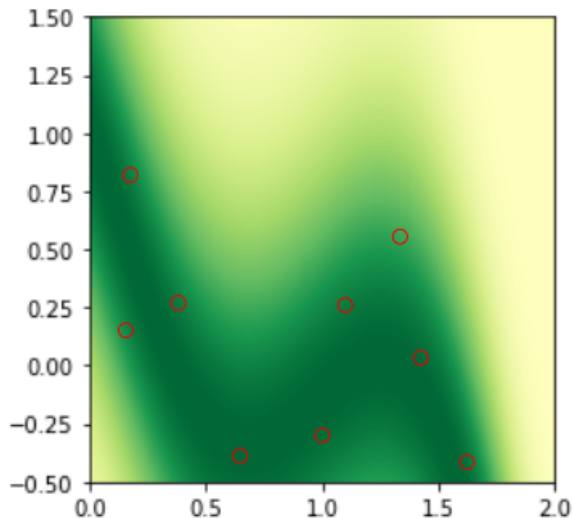
# The Actual Problem

We only have samples.



# The Actual Problem

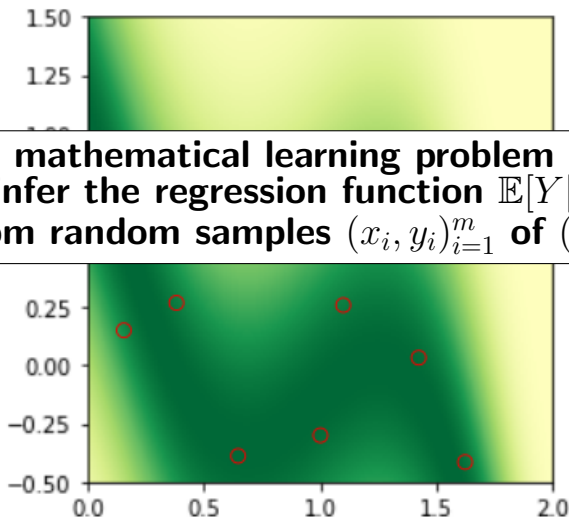
We only have samples.



# The Actual Problem

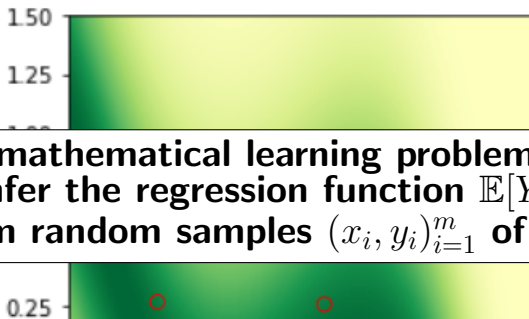
We only have samples.

**A mathematical learning problem seeks to infer the regression function  $\mathbb{E}[Y|X = x]$  from random samples  $(x_i, y_i)_{i=1}^m$  of  $(X, Y)$ .**



# The Actual Problem

We only have samples.



**A mathematical learning problem seeks to infer the regression function  $\mathbb{E}[Y|X = x]$  from random samples  $(x_i, y_i)_{i=1}^m$  of  $(X, Y)$ .**

**More generally we would like to minimize  $\mathbb{E}[\mathcal{L}(f(X), Y)]$  with general loss function.**

$\mathcal{L}(y, y') = (y - y')^2 \rightsquigarrow$  **quadratic loss**

$\mathcal{L}(y, y') = y \log(y') + (1 - y) \log(1 - y') \rightsquigarrow$  **cross-entropy loss.**

## 1.2.2 Empirical Minimization and Hypothesis Space

# Sampling

## Empirical Error

Given  $\mathbf{z} = ((X^{(1)}, Y^{(1)}), \dots, (X^{(m)}, Y^{(m)}))$  be i.i.d. with  $(X^{(1)}, Y^{(1)}) \sim (X, Y)$ . Define the *empirical error*

$$\mathcal{E}_{\mathbf{z}}(f) := \frac{1}{m} \sum_{i=1}^m (f(X^{(i)}) - Y^{(i)})^2.$$

# Sampling

## Empirical Error

Given  $\mathbf{z} = ((X^{(1)}, Y^{(1)}), \dots, (X^{(m)}, Y^{(m)}))$  be i.i.d. with  $(X^{(1)}, Y^{(1)}) \sim (X, Y)$ . Define the *empirical error*

$$\mathcal{E}_{\mathbf{z}}(f) := \frac{1}{m} \sum_{i=1}^m (f(X^{(i)}) - Y^{(i)})^2.$$

Given  $\mathbf{z}$  the empirical error can actually be computed!

# Sampling

## Empirical Error

Given  $\mathbf{z} = ((X^{(1)}, Y^{(1)}), \dots, (X^{(m)}, Y^{(m)}))$  be i.i.d. with  $(X^{(1)}, Y^{(1)}) \sim (X, Y)$ . Define the *empirical error*

$$\mathcal{E}_{\mathbf{z}}(f) := \frac{1}{m} \sum_{i=1}^m (f(X^{(i)}) - Y^{(i)})^2.$$

Given  $\mathbf{z}$  the empirical error can actually be computed!

## Defect

The defect of  $f$  is defined as

$$L_{\mathbf{z}}(f) := \mathcal{E}(f) - \mathcal{E}_{\mathbf{z}}(f).$$

# Sampling

## Empirical Error

Given  $\mathbf{z} = ((X^{(1)}, Y^{(1)}), \dots, (X^{(m)}, Y^{(m)}))$  be i.i.d. with  $(X^{(1)}, Y^{(1)}) \sim (X, Y)$ . Define the *empirical error*

$$\mathcal{E}_{\mathbf{z}}(f) := \frac{1}{m} \sum_{i=1}^m (f(X^{(i)}) - Y^{(i)})^2.$$

Given  $\mathbf{z}$  the empirical error can actually be computed!

## Defect

The defect of  $f$  is defined as

$$L_{\mathbf{z}}(f) := \mathcal{E}(f) - \mathcal{E}_{\mathbf{z}}(f).$$

Can we control the defect? If yes, we actually have some hope of approximating the regression function.



# Data Generating Distribution

We suppose that there exists a probability distribution on  $\mathbb{R}^{784}$  that randomly generates handwritten digits.

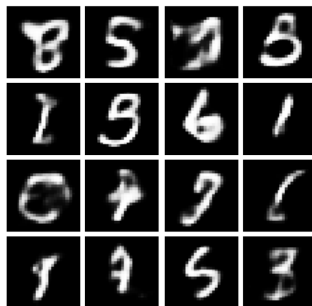
# Data Generating Distribution

We suppose that there exists a probability distribution on  $\mathbb{R}^{784}$  that randomly generates handwritten digits.



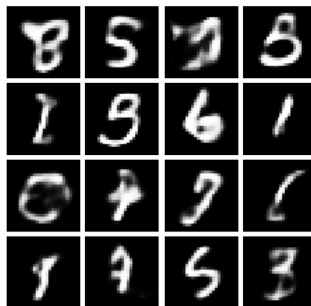
# Data Generating Distribution

We suppose that there exists a probability distribution on  $\mathbb{R}^{784}$  that randomly generates handwritten digits.



# Data Generating Distribution

We suppose that there exists a probability distribution on  $\mathbb{R}^{784}$  that randomly generates handwritten digits.



**Variational Autoencoder Demo**

# Concentration Inequalities

## Bernstein Inequality

Suppose that  $(\xi^{(i)})_{i=1}^m$  i.i.d. with  $\xi^{(1)} \sim \xi$  with mean  $\mathbb{E}(\xi) = \mu$  and  $\mathbb{V}(\xi) = \sigma^2$ . Suppose that  $|\xi - \mu| \leq M$  with probability 1. Then

$$\mathbb{P} \left\{ \left| \frac{1}{m} \sum_{i=1}^m \xi^{(i)} - \mu \right| \geq \varepsilon \right\} \leq 2e^{-\frac{m\varepsilon^2}{2\left(\sigma^2 + \frac{1}{3}M\varepsilon\right)}}.$$

# Bounding the Defect

## Theorem A

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$  and let  $\sigma_f^2 = \mathbb{V}[(f(X) - Y)^2]$ . Suppose that  $|f(X) - Y| \leq M$  almost everywhere. Then, for all  $\varepsilon > 0$ ,

$$\mathbb{P} \{ |L_{\mathbf{z}}(f)| \leq \varepsilon \} \geq 1 - 2e^{-\frac{m\varepsilon^2}{2(\sigma_f^2 + \frac{1}{3}M\varepsilon)}}.$$

# Bounding the Defect

## Theorem A

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$  and let  $\sigma_f^2 = \mathbb{V}[(f(X) - Y)^2]$ . Suppose that  $|f(X) - Y| \leq M$  almost everywhere. Then, for all  $\varepsilon > 0$ ,

$$\mathbb{P} \{ |L_{\mathbf{z}}(f)| \leq \varepsilon \} \geq 1 - 2e^{-\frac{m\varepsilon^2}{2(\sigma_f^2 + \frac{1}{3}M\varepsilon)}}.$$

## Proof.

Apply Bernstein Inequality to  $\xi = (f(X) - Y)^2$ . □

# Bounding the Defect

## Theorem A

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$  and let  $\sigma_f^2 = \mathbb{V}[(f(X) - Y)^2]$ . Suppose that  $|f(X) - Y| \leq M$  almost everywhere. Then, for all  $\varepsilon > 0$ ,

$$\mathbb{P} \{ |L_{\mathbf{z}}(f)| \leq \varepsilon \} \geq 1 - 2e^{-\frac{m\varepsilon^2}{2(\sigma_f^2 + \frac{1}{3}M\varepsilon)}}.$$

## Proof.

Apply Bernstein Inequality to  $\xi = (f(X) - Y)^2$ . □

Are we done?? We could just minimize the empirical error and bound the defect...



# Bounding the Defect

## Theorem A

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$  and let  $\sigma_f^2 = \mathbb{V}[(f(X) - Y)^2]$ . Suppose that  $|f(X) - Y| \leq M$  almost everywhere. Then, for all  $\varepsilon > 0$ ,

$$\mathbb{P} \{ |L_{\mathbf{z}}(f)| \leq \varepsilon \} \geq 1 - 2e^{-\frac{m\varepsilon^2}{2(\sigma_f^2 + \frac{1}{3}M\varepsilon)}}.$$

## Proof.

Apply Bernstein Inequality to  $\xi = (f(X) - Y)^2$ . □

Are we done?? We could just minimize the empirical error and bound the defect...



Any  $f$  vanishing on the sample points makes the empirical error vanish!!!

# Hypothesis Space

## Definition

Let  $\mathcal{H}$  be a compact subset of the Banach space

$\{f : X \rightarrow Y, \text{ continuous}\}$  with norm  $\|f\| := \max_{x \in X} |f(x)|$ . We call  $\mathcal{H}$  *hypothesis space* or *model space*.

# Hypothesis Space

## Definition

Let  $\mathcal{H}$  be a compact subset of the Banach space

$\{f : X \rightarrow Y, \text{ continuous}\}$  with norm  $\|f\| := \max_{x \in X} |f(x)|$ . We call  $\mathcal{H}$  *hypothesis space* or *model space*.

## Best Approximation in $\mathcal{H}$

Define the *best approximation in  $\mathcal{H}$*  via

$$\hat{f}_{\mathcal{H}} := \operatorname{argmin}_{f \in \mathcal{H}} \mathcal{E}(f) = \operatorname{argmin}_{f \in \mathcal{H}} \|\hat{f} - f\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}.$$

# Hypothesis Space

## Definition

Let  $\mathcal{H}$  be a compact subset of the Banach space  $\{f : X \rightarrow Y, \text{ continuous}\}$  with norm  $\|f\| := \max_{x \in X} |f(x)|$ . We call  $\mathcal{H}$  *hypothesis space* or *model space*.

## Best Approximation in $\mathcal{H}$

Define the *best approximation in  $\mathcal{H}$*  via

$$\hat{f}_{\mathcal{H}} := \operatorname{argmin}_{f \in \mathcal{H}} \mathcal{E}(f) = \operatorname{argmin}_{f \in \mathcal{H}} \|\hat{f} - f\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}.$$

## Empirical Regression Function

Given  $\mathbf{z}$  define the *empirical regression function* as

$$\hat{f}_{\mathcal{H}, \mathbf{z}} := \operatorname{argmin}_{f \in \mathcal{H}} \mathcal{E}_{\mathbf{z}}(f).$$

# Hypothesis Space

## Definition

Let  $\mathcal{H}$  be a compact subset of the Banach space

$\{f : X \rightarrow Y, \text{ continuous}\}$  with norm  $\|f\| := \max_{x \in X} |f(x)|$ . We call  $\mathcal{H}$  *hypothesis space* or *model space*.

## Best Approximation in $\mathcal{H}$

Define the *best approximation* in  $\mathcal{H}$  via

$$\hat{f}_{\mathcal{H}} := \operatorname{argmin}_{f \in \mathcal{H}} \mathcal{E}(f) = \operatorname{argmin}_{f \in \mathcal{H}} \|\hat{f} - f\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}.$$

## Empirical Regression Function

Given  $\mathbf{z}$  define the *empirical regression function* as

$$\hat{f}_{\mathcal{H}, \mathbf{z}} := \operatorname{argmin}_{f \in \mathcal{H}} \mathcal{E}_{\mathbf{z}}(f).$$



The empirical regression function can be computed!

### 1.2.3. Bias-Variance Decomposition

# Generalization- and Approximation Error

## Theorem (Bias-Variance Decomposition)

It holds that

$$\|\hat{f}_{\mathcal{H},\mathbf{z}} - \hat{f}\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2 = \left( \mathcal{E}(\hat{f}_{\mathcal{H},\mathbf{z}}) - \mathcal{E}(\hat{f}_{\mathcal{H}}) \right) + \|\hat{f}_{\mathcal{H}} - \hat{f}\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2.$$

The first term is called *generalization error* and the second term is called *approximation error*.

# Generalization- and Approximation Error

## Theorem (Bias-Variance Decomposition)

It holds that

$$\|\hat{f}_{\mathcal{H},\mathbf{z}} - \hat{f}\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2 = \left( \mathcal{E}(\hat{f}_{\mathcal{H},\mathbf{z}}) - \mathcal{E}(\hat{f}_{\mathcal{H}}) \right) + \|\hat{f}_{\mathcal{H}} - \hat{f}\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2.$$

The first term is called *generalization error* and the second term is called *approximation error*.

## Proof.

By the Main Regression Theorem

$$\begin{aligned} \|\hat{f}_{\mathcal{H},\mathbf{z}} - \hat{f}\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2 &= \mathcal{E}(\hat{f}_{\mathcal{H},\mathbf{z}}) - \mathcal{E}(\hat{f}) \\ &= \mathcal{E}(\hat{f}_{\mathcal{H},\mathbf{z}}) - \mathcal{E}(\hat{f}_{\mathcal{H}}) + \mathcal{E}(\hat{f}_{\mathcal{H}}) - \mathcal{E}(\hat{f}) \\ &= \left( \mathcal{E}(\hat{f}_{\mathcal{H},\mathbf{z}}) - \mathcal{E}(\hat{f}_{\mathcal{H}}) \right) + \|\hat{f}_{\mathcal{H}} - \hat{f}\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2. \end{aligned}$$





# Generalization- and Approximation Error

## Theorem (Bias-Variance Decomposition)

It holds that

$$\|\hat{f}_{\mathcal{H},\mathbf{z}} - \hat{f}\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2 = \left( \mathcal{E}(\hat{f}_{\mathcal{H},\mathbf{z}}) - \mathcal{E}(\hat{f}_{\mathcal{H}}) \right) + \|\hat{f}_{\mathcal{H}} - \hat{f}\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2.$$

The first term is called *generalization error* and the second term is called *approximation error*.

## Proof.

By the M

**Our goal is to make the empirical error**

$$\|\hat{f}_{\mathcal{H},\mathbf{z}} - \hat{f}\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2$$

**as small as possible.**

$$= \left( \mathcal{E}(\hat{f}_{\mathcal{H},\mathbf{z}}) - \mathcal{E}(\hat{f}_{\mathcal{H}}) \right) + \|\hat{f}_{\mathcal{H}} - \hat{f}\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2.$$



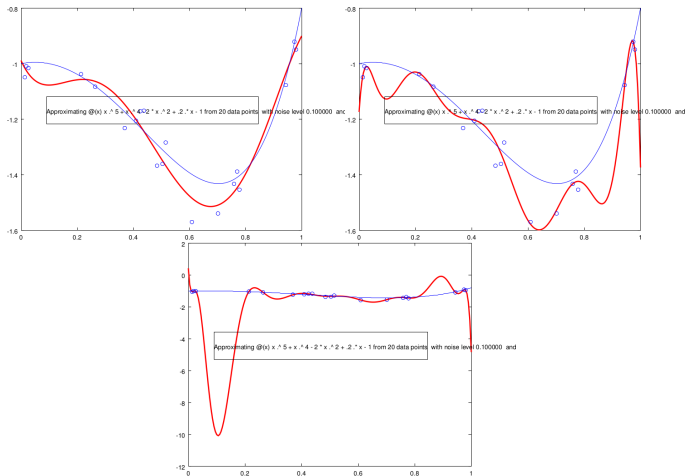


Figure: Blue:  $f_{\mathcal{H}}$ , Red:  $f_{\mathcal{H},z}$ ,  $m = 10$ ,  $\mathcal{H}$  = polynomials of degree 5, 15, 20 (from top left to bottom).

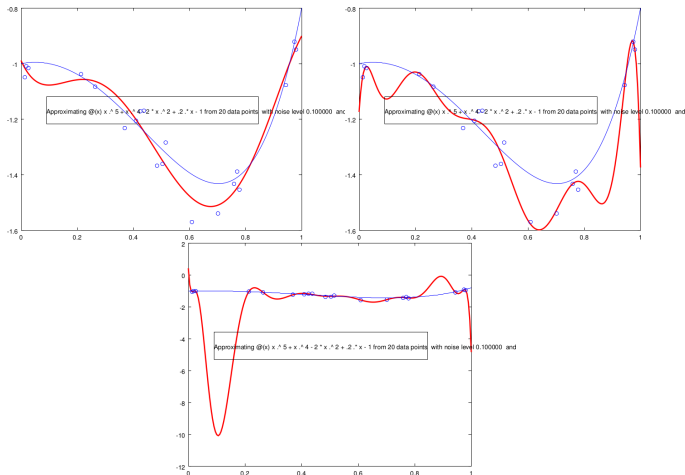


Figure: Blue:  $f_{\mathcal{H}}$ , Red:  $f_{\mathcal{H},z}$ ,  $m = 10$ ,  $\mathcal{H} =$  polynomials of degree 5, 15, 20 (from top left to bottom).

If  $\mathcal{H}$  is too complex, the sampling error increases.

# The Bias-Variance Trade-Off

If we keep the sample size  $m$  fixed and enlarge the hypothesis space  $\mathcal{H}$ , the approximation error will certainly decrease, **BUT** the sample error will increase – this is exactly what we observed experimentally!

# The Bias-Variance Trade-Off

If we keep the sample size  $m$  fixed and enlarge the hypothesis space  $\mathcal{H}$ , the approximation error will certainly decrease, **BUT** the sample error will increase – this is exactly what we observed experimentally!

Bishop [Neural Networks for Pattern Recognition (1995)]

*“A model which is too simple, or too inflexible, will have a large bias, while one which has too much flexibility in relation to the particular data set will have a large variance. Bias and variance are complementary quantities, and the best generalization is obtained when we have the best compromise between the conflicting requirements of small bias and small variance.”*

# The Bias-Variance Trade-Off

If we keep the sample size  $m$  fixed and enlarge the hypothesis space  $\mathcal{H}$ , the approximation error will certainly decrease, **BUT** the sample error will increase – this is exactly what we observed experimentally!

Bishop [Neural Networks for Pattern Recognition (1995)]

*“A model which is too simple, or too inflexible, will have a large bias, while one which has too much flexibility in relation to the particular data set will have a large variance. Bias and variance are complementary quantities, and the best generalization is obtained when we have the best compromise between the conflicting requirements of small bias and small variance.”*

Bias-Variance Problem

What are the precise relations between the number of samples  $m$  and the “capacity” of our hypothesis space  $\mathcal{H}$ ?

### 1.2.4 Bounds on the Generalization Error $\mathcal{E}(\hat{f}_{\mathcal{H},z}) - \mathcal{E}(\hat{f}_{\mathcal{H}})$ .

# Covering Numbers

## Definition

Let  $S$  be a metric space and  $s > 0$ . Define the *covering number*  $\mathcal{N}(S, s)$  to be the minimal  $l \in \mathcal{N}$  such that there exist  $l$  disks in  $S$  with radius  $s$  covering  $S$ .



# Covering Numbers

## Definition

Let  $S$  be a metric space and  $s > 0$ . Define the *covering number*  $\mathcal{N}(S, s)$  to be the minimal  $l \in \mathcal{N}$  such that there exist  $l$  disks in  $S$  with radius  $s$  covering  $S$ .



# Covering Numbers

## Definition

Let  $S$  be a metric space and  $s > 0$ . Define the *covering number*  $\mathcal{N}(S, s)$  to be the minimal  $l \in \mathcal{N}$  such that there exist  $l$  disks in  $S$  with radius  $s$  covering  $S$ .



Scaling of  $\mathcal{N}(S, s)$  with  $s$  is a measure of complexity of  $S$  termed *metric entropy*.

# Abstract Analysis of Generalization Error

## Theorem B

Let  $\mathcal{H} \subset C(X)$  be a hypothesis class. Assume that for all  $f \in \mathcal{H}$  it holds that  $|f(X) - Y| < M$  a.e. Then, for all  $\varepsilon > 0$ ,

$$\mathbb{P} \left( \sup_{f \in \mathcal{H}} |L_{\mathbf{z}}(f)| \leq \varepsilon \right) \geq 1 - \mathcal{N}(\mathcal{H}, \frac{\varepsilon}{8M}) 2e^{-\frac{m\varepsilon^2}{4(2\sigma^2 + \frac{1}{3}M^2\varepsilon)}},$$

where  $\sigma^2 := \sup_{f \in \mathcal{H}} \sigma_f^2$ .

## Proof.

First show that for all  $f, g$  with  $\|f - g\| \leq \tau$  it holds that

$$|\mathcal{E}(f) - \mathcal{E}(g)| \leq 2M\tau \quad \text{and} \quad |\mathcal{E}_{\mathbf{z}}(f) - \mathcal{E}_{\mathbf{z}}(g)| \leq 2M\tau.$$

Cover  $\mathcal{H}$  with balls  $(U_i)_{i=1}^{\mathcal{N}(\mathcal{H}, \epsilon/(8M))}$  with center  $f_i$  of radius  $\frac{\epsilon}{8M}$ . By the estimate above it holds that

$$\left( \sup_{f \in U_i} |\mathcal{E}_{\mathbf{z}}(f) - \mathcal{E}(f)| > \epsilon \right) \Rightarrow (|\mathcal{E}_{\mathbf{z}}(f_i) - \mathcal{E}(f_i)| > \epsilon/2)$$

Then by this fact and Theorem A it holds that

$$\begin{aligned} \mathbb{P} \left( \sup_{f \in \mathcal{H}} |L_{\mathbf{z}}(f)| > \epsilon \right) &\leq \sum_{i=1}^{\mathcal{N}(\mathcal{H}, \epsilon/(8M))} \mathbb{P} \left( \sup_{f \in U_i} |\mathcal{E}_{\mathbf{z}}(f) - \mathcal{E}(f)| > \epsilon \right) \\ &\leq \sum_{i=1}^{\mathcal{N}(\mathcal{H}, \epsilon/(8M))} \mathbb{P} (|\mathcal{E}_{\mathbf{z}}(f_i) - \mathcal{E}(f_i)| > \epsilon/2) \\ &\leq \mathcal{N}(\mathcal{H}, \epsilon/(8M)) 2e^{-\frac{m\epsilon^2}{4(2\sigma^2 + \frac{1}{3}M^2\epsilon)}}. \end{aligned}$$



# Abstract Analysis of Generalization Error

## Lemma

Let  $\varepsilon > 0$  and  $0 < \delta < 1$  such that

$$\mathbb{P}(\sup_{f \in \mathcal{H}} |L_{\mathbf{z}}(f)| \leq \varepsilon) \geq 1 - \delta.$$

Then

$$\mathbb{P}(\mathcal{E}(\hat{f}_{\mathcal{H}, \mathbf{z}}) - \mathcal{E}(\hat{f}_{\mathcal{H}}) \leq 2\varepsilon) \geq 1 - \delta.$$

# Abstract Analysis of Generalization Error

## Lemma

Let  $\varepsilon > 0$  and  $0 < \delta < 1$  such that

$$\mathbb{P}(\sup_{f \in \mathcal{H}} |L_{\mathbf{z}}(f)| \leq \varepsilon) \geq 1 - \delta.$$

Then

$$\mathbb{P}(\mathcal{E}(\hat{f}_{\mathcal{H}, \mathbf{z}}) - \mathcal{E}(\hat{f}_{\mathcal{H}}) \leq 2\varepsilon) \geq 1 - \delta.$$

## Proof.

Suppose that  $\sup_{f \in \mathcal{H}} |L_{\mathbf{z}}(f)| \leq \varepsilon$ . Then  $|\mathcal{E}_{\mathbf{z}}(\hat{f}_{\mathcal{H}, \mathbf{z}}) - \mathcal{E}(\hat{f}_{\mathcal{H}, \mathbf{z}})| \leq \varepsilon$ ,  $|\mathcal{E}_{\mathbf{z}}(\hat{f}_{\mathcal{H}}) - \mathcal{E}(\hat{f}_{\mathcal{H}})| \leq \varepsilon$  and  $\mathcal{E}_{\mathbf{z}}(\hat{f}_{\mathcal{H}, \mathbf{z}}) - \mathcal{E}_{\mathbf{z}}(\hat{f}_{\mathcal{H}}) \leq 0$ . It follows that

$$\begin{aligned} \mathcal{E}(\hat{f}_{\mathcal{H}, \mathbf{z}}) - \mathcal{E}(\hat{f}_{\mathcal{H}}) &= \mathcal{E}(\hat{f}_{\mathcal{H}, \mathbf{z}}) - \mathcal{E}_{\mathbf{z}}(\hat{f}_{\mathcal{H}, \mathbf{z}}) + \mathcal{E}_{\mathbf{z}}(\hat{f}_{\mathcal{H}, \mathbf{z}}) - \mathcal{E}_{\mathbf{z}}(\hat{f}_{\mathcal{H}}) + \\ &\quad \mathcal{E}_{\mathbf{z}}(\hat{f}_{\mathcal{H}}) - \mathcal{E}(\hat{f}_{\mathcal{H}}) \\ &\leq |\mathcal{E}_{\mathbf{z}}(\hat{f}_{\mathcal{H}, \mathbf{z}}) - \mathcal{E}(\hat{f}_{\mathcal{H}, \mathbf{z}})| + |\mathcal{E}_{\mathbf{z}}(\hat{f}_{\mathcal{H}}) - \mathcal{E}(\hat{f}_{\mathcal{H}})| \leq 2\varepsilon. \end{aligned}$$

# Abstract Analysis of Generalization Error

## Theorem C

Let  $\mathcal{H}$  be a hypothesis class. Assume that for all  $f \in \mathcal{H}$  it holds that  $|f(X) - Y| < M$  a.e. Then, for all  $\varepsilon > 0$ ,

$$\mathbb{P} \left( \mathcal{E}(\hat{f}_{\mathcal{H}, \mathbf{z}}) - \mathcal{E}(\hat{f}_{\mathcal{H}}) \leq \varepsilon \right) \geq 1 - \mathcal{N}(\mathcal{H}, \frac{\varepsilon}{16M}) 2e^{-\frac{m\varepsilon^2}{8(2\sigma^2 + \frac{1}{3}M^2\varepsilon)}},$$

where  $\sigma^2 := \sup_{f \in \mathcal{H}} \sigma_f^2$ .

# Abstract Analysis of Generalization Error

## Theorem C

Let  $\mathcal{H}$  be a hypothesis class. Assume that for all  $f \in \mathcal{H}$  it holds that  $|f(X) - Y| < M$  a.e. Then, for all  $\varepsilon > 0$ ,

$$\mathbb{P} \left( \mathcal{E}(\hat{f}_{\mathcal{H}, \mathbf{z}}) - \mathcal{E}(\hat{f}_{\mathcal{H}}) \leq \varepsilon \right) \geq 1 - \mathcal{N}(\mathcal{H}, \frac{\varepsilon}{16M}) 2e^{-\frac{m\varepsilon^2}{8(2\sigma^2 + \frac{1}{3}M^2\varepsilon)}},$$

where  $\sigma^2 := \sup_{f \in \mathcal{H}} \sigma_f^2$ .

## Proof.

Apply Lemma and Theorem B with  $\epsilon \leftrightarrow \epsilon/2$ . □



# Abstract Analysis of Generalization Error

## Question

Given  $\varepsilon, \delta > 0$ , how many samples  $m$  do we need such that the probability that the generalization error is  $\leq \varepsilon$  is at least  $1 - \delta$ ?

# Abstract Analysis of Generalization Error

## Question

Given  $\varepsilon, \delta > 0$ , how many samples  $m$  do we need such that the probability that the generalization error is  $\leq \varepsilon$  is at least  $1 - \delta$ ?

## Answer

By the previous theorem it suffices to choose

$$m \geq \frac{8(4\sigma^2 + \frac{1}{3}M^2\varepsilon)}{\varepsilon^2} \left( \ln(2\mathcal{N}(\mathcal{H}, \frac{\varepsilon}{16M})) + \ln(\frac{1}{\delta}) \right).$$

# Abstract Analysis of Generalization Error

## Question

Given  $\varepsilon, \delta > 0$ , how many samples  $m$  do we need such that the probability that the generalization error is  $\leq \varepsilon$  is at least  $1 - \delta$ ?

## Answer

By the previous theorem it suffices to choose

$$m \geq \frac{8(4\sigma^2 + \frac{1}{3}M^2\varepsilon)}{\varepsilon^2} \left( \ln(2\mathcal{N}(\mathcal{H}, \frac{\varepsilon}{16M})) + \ln(\frac{1}{\delta}) \right).$$

## Question

How to bound the covering number?

## 1.2.5 A Simple Example

# Unit Balls

For  $R \in (0, \infty)$  and  $n \in \mathbb{N}$  let

$$B_{R,n} := \{w \in \mathbb{R}^n : \|w\|_\infty \leq R\}.$$

# Unit Balls

For  $R \in (0, \infty)$  and  $n \in \mathbb{N}$  let

$$B_{R,n} := \{w \in \mathbb{R}^n : \|w\|_\infty \leq R\}.$$

## Lemma

We have

$$\mathcal{N}(B_{R,n}, \tau) \leq \left\lceil \frac{R}{\tau} \right\rceil^n.$$

# Unit Balls

For  $R \in (0, \infty)$  and  $n \in \mathbb{N}$  let

$$B_{R,n} := \{w \in \mathbb{R}^n : \|w\|_\infty \leq R\}.$$

## Lemma

We have

$$\mathcal{N}(B_{R,n}, \tau) \leq \left\lceil \frac{R}{\tau} \right\rceil^n.$$

## Proof.

Just decompose  $B_{R,n}$  into  $\left\lceil \frac{R}{\tau} \right\rceil^n$  cubes...



# Parametrized Hypothesis Classes

Let  $\mathcal{F} : B_{R,n} \rightarrow C(\Omega)$  and  $L \in (0, \infty)$  with

$$\|\mathcal{F}(v) - \mathcal{F}(w)\|_{L^\infty(\Omega)} \leq L\|v - w\|_\infty$$

for all  $v, w \in B_{R,n}$ .



# Parametrized Hypothesis Classes

Let  $\mathcal{F} : B_{R,n} \rightarrow C(\Omega)$  and  $L \in (0, \infty)$  with

$$\|\mathcal{F}(v) - \mathcal{F}(w)\|_{L^\infty(\Omega)} \leq L\|v - w\|_\infty$$

for all  $v, w \in B_{R,n}$ .

## Lemma

We have

$$\mathcal{N}(\mathcal{F}(B_{R,n}), \tau) \leq \left\lceil \frac{LR}{\tau} \right\rceil^n.$$

# Parametrized Hypothesis Classes

Let  $\mathcal{F} : B_{R,n} \rightarrow C(\Omega)$  and  $L \in (0, \infty)$  with

$$\|\mathcal{F}(v) - \mathcal{F}(w)\|_{L^\infty(\Omega)} \leq L\|v - w\|_\infty$$

for all  $v, w \in B_{R,n}$ .

## Lemma

We have

$$\mathcal{N}(\mathcal{F}(B_{R,n}), \tau) \leq \left\lceil \frac{LR}{\tau} \right\rceil^n.$$

## Proof.

Use the Lipschitz property to deduce that a  $\tau/L$  cover of  $B_{R,n}$  induced a  $\tau$  cover of  $\mathcal{F}(B_{R,n})$ . □

## Example: Linear Regression

Let  $\mathcal{F}(w) := \sum_{i=1}^n w_i \varphi_i$ , as in the linear Regression example.

## Example: Linear Regression

Let  $\mathcal{F}(w) := \sum_{i=1}^n w_i \varphi_i$ , as in the linear Regression example.

### Lemma

It holds that

$$\|\mathcal{F}(w) - \mathcal{F}(v)\|_{L^\infty(\Omega)} \leq \left\| \sum_{i=1}^n |\varphi_i| \right\|_{L^\infty(\Omega)} \|v - w\|_\infty.$$

## Example: Linear Regression

Let  $\mathcal{F}(w) := \sum_{i=1}^n w_i \varphi_i$ , as in the linear Regression example.

### Lemma

It holds that

$$\|\mathcal{F}(w) - \mathcal{F}(v)\|_{L^\infty(\Omega)} \leq \left\| \sum_{i=1}^n |\varphi_i| \right\|_{L^\infty(\Omega)} \|v - w\|_\infty.$$

### Corollary

Let  $\mathcal{H} = \{\sum_{i=1}^n w_i \varphi_i : \|w\|_\infty \leq R\}$ . Then

$$\mathcal{N}(\mathcal{H}, \tau) \leq \left\lceil \frac{R \left\| \sum_{i=1}^n |\varphi_i| \right\|_{L^\infty(\Omega)}}{\tau} \right\rceil^n.$$

# Analysis of Linear Regression

## Theorem

Consider linear regression as above and suppose that we have the approximation error estimate

$$\inf_{f \in \mathcal{H}} \|\hat{f} - f\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2 \leq \frac{\epsilon}{2}.$$

Then

$$m \gtrsim \frac{(n \cdot \text{polylog}(\epsilon) + \ln(\frac{1}{\delta}))}{\epsilon^2}$$

independent training samples suffice to get an empirical error  $\epsilon$  with probability  $\geq 1 - \delta$ .

# Analysis of Linear Regression

## Theorem

Consider linear regression as above and suppose that we have the approximation error estimate

$$\inf_{f \in \mathcal{H}} \|\hat{f} - f\|_{L^2(\mathbb{R}^d, d\mathbb{P}_X)}^2 \leq \frac{\epsilon}{2}.$$

Then

$$m \gtrsim \frac{(n \cdot \text{polylog}(\epsilon) + \ln(\frac{1}{\delta}))}{\epsilon^2}$$

independent training samples suffice to get an empirical error  $\epsilon$  with probability  $\geq 1 - \delta$ .

## Rule of Thumb to avoid Overfitting

If we double the dimension (for example polynomial degree) we need to double the number of training samples!

## More Advanced Topics



## More Advanced Topics

- General Loss function (see for example Devroye, Györfi, Lugosi: A Probabilistic Theory of Pattern Recognition)

## More Advanced Topics

- General Loss function (see for example Devroye, Györfi, Lugosi: A Probabilistic Theory of Pattern Recognition)
- Other (sharper) complexity measures such as Rademacher complexity, VC dimension, empirical oscillation, .... (see for example Devroye, Györfi, Lugosi: A Probabilistic Theory of Pattern Recognition)

## More Advanced Topics

- General Loss function (see for example Devroye, Györfi, Lugosi: A Probabilistic Theory of Pattern Recognition)
- Other (sharper) complexity measures such as Rademacher complexity, VC dimension, empirical oscillation, .... (see for example Devroye, Györfi, Lugosi: A Probabilistic Theory of Pattern Recognition)
- Weaker conditions on  $(X, Y)$  (we essentially require bounded noise!!) (see for example Mendelson: Learning without Concentration)

## More Advanced Topics

- General Loss function (see for example Devroye, Györfi, Lugosi: A Probabilistic Theory of Pattern Recognition)
- Other (sharper) complexity measures such as Rademacher complexity, VC dimension, empirical oscillation, .... (see for example Devroye, Györfi, Lugosi: A Probabilistic Theory of Pattern Recognition)
- Weaker conditions on  $(X, Y)$  (we essentially require bounded noise!!) (see for example Mendelson: Learning without Concentration)
- Better learning procedures than ERM (see for example Mendelson: An Optimal Unrestricted Learning Procedure)

## More Advanced Topics

- General Loss function (see for example Devroye, Györfi, Lugosi: A Probabilistic Theory of Pattern Recognition)
- Other (sharper) complexity measures such as Rademacher complexity, VC dimension, empirical oscillation, .... (see for example Devroye, Györfi, Lugosi: A Probabilistic Theory of Pattern Recognition)
- Weaker conditions on  $(X, Y)$  (we essentially require bounded noise!!) (see for example Mendelson: Learning without Concentration)
- Better learning procedures than ERM (see for example Mendelson: An Optimal Unrestricted Learning Procedure)
- Better sampling procedures (see for example Cohen, Migliorati: Optimal Weighted Least Squares Methods)