

# Machine Learning for Physical Scientists

## Lecture 2

### The Simplest Supervised Learning: Linear Regression

# Framework of Statistical Learning Theory

(supervised learning)

$X$ : **Instance Space** (e.g.  $\mathbb{R}^{16 \times 16}$  for 16x16 greyscale images)

$Y$ : **Label Space** (e.g.  $\mathbb{R}$  for regression or  $\{1, \dots, k\}$  for multi-class classification)

$\mathcal{D}$ : **Probability Distribution** over  $X \times Y$  (*unknown, but can sample from*)

$\ell : Y \times Y \rightarrow \mathbb{R}_{\geq 0}$  **Loss** or **Cost Function** (e.g.  $\ell(y, \hat{y}) = (y - \hat{y})^2$  for  $Y = \mathbb{R}$ )

## Objective

Given a **training set**  $S = \left\{ (x_i, y_i) \right\}_{i=1}^m$  drawn i.i.d. from  $\mathcal{D}$ , return hypothesis (predictor)

$h : X \rightarrow Y$  that minimizes the **population loss** or **expected risk**:

$$L_{\mathcal{D}}(h) := \mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell(y, h(x))]$$

## Approximate Approach

Predetermine or assume a **hypotheses space**  $\mathcal{H} \subset Y^X$ , and return hypothesis  $h \in \mathcal{H}$  that minimizes **sample loss** or **empirical loss** or **empirical risk**:

$$L_S(h) := \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(x_i))$$

# Framework of Statistical Learning Theory

(supervised learning)

$X$ : **Instance Space** (e.g.  $\mathbb{R}^{16 \times 16}$  for 16x16 greyscale images)

$Y$ : **Label Space** (e.g.  $\mathbb{R}$  for regression or  $\{1, \dots, k\}$  for multi-class classification)

$\mathcal{D}$ : **Probability Distribution** over  $X \times Y$  (*unknown, but can sample from*)

$\ell : Y \times Y \rightarrow \mathbb{R}_{\geq 0}$  **Loss** or **Cost Function** (e.g.  $\ell(y, \hat{y}) = (y - \hat{y})^2$  for  $Y = \mathbb{R}$ )

## Objective

Given a **training set**  $S = \left\{ (x_i, y_i) \right\}_{i=1}^m$  drawn i.i.d. from  $\mathcal{D}$ , return hypothesis (predictor)

$h : X \rightarrow Y$  that minimizes the **population loss** or **expected risk**:

$$L_{\mathcal{D}}(h) := \mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell(y, h(x))]$$

## Approximate Approach

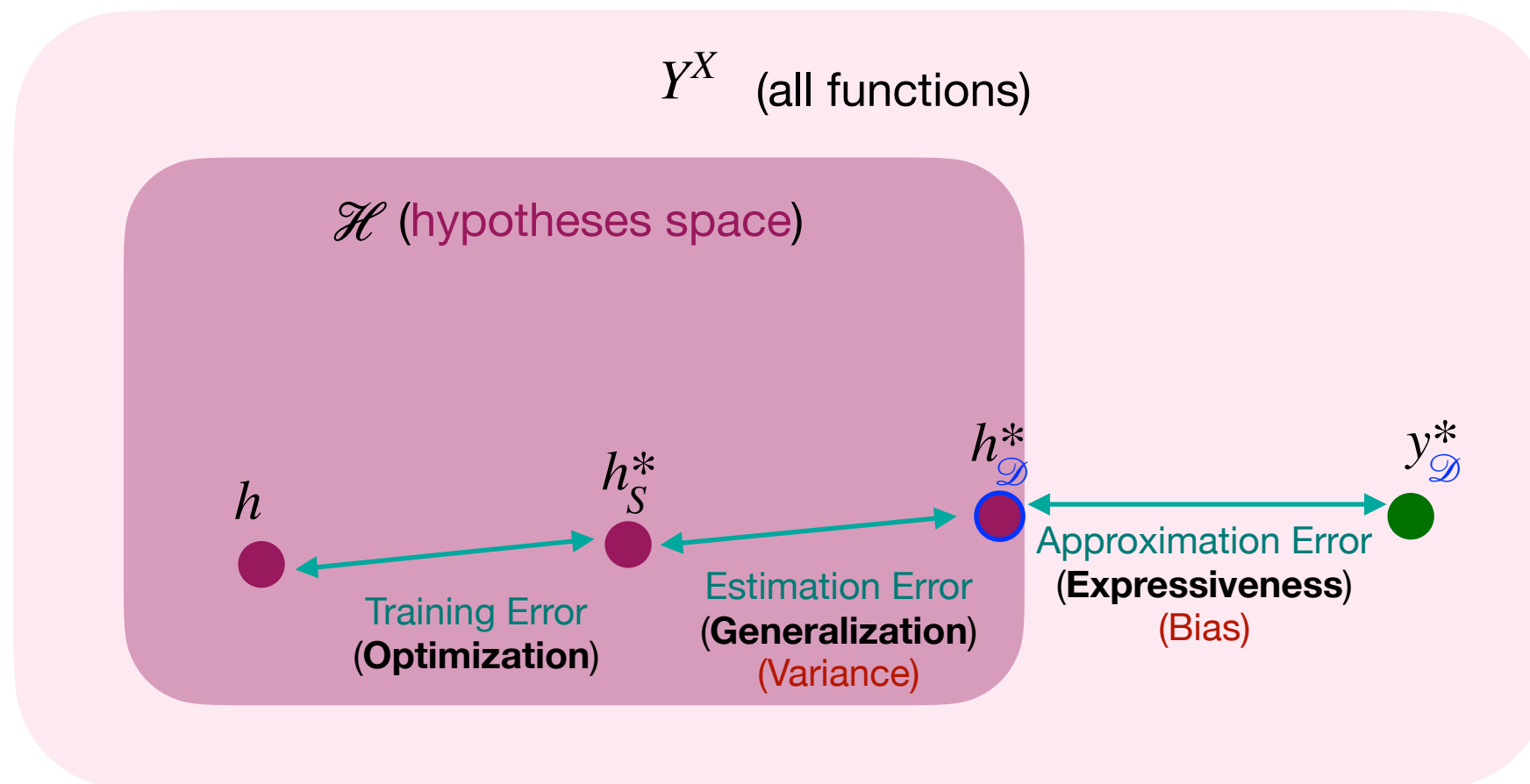
Predetermine or assume a **hypotheses space**  $\mathcal{H} \subset Y^X$ , and return hypothesis  $h \in \mathcal{H}$  that minimizes **sample loss** or **empirical loss** or **empirical risk**:

$$L_S(h) := \frac{1}{m} \sum_{i=1}^m \ell(y_i, h(x_i))$$

Note: An algorithm that searches for empirically optimal  $h_S^*$  is called a “learning algorithm.”

# Jargons in Statistical Learning Theory (SLT)

## (Expressiveness, Generalization, Optimization)



$y_{\mathcal{D}}^*$ : **ground truth** (minimizer of population loss over  $Y^X$ )

$h_{\mathcal{D}}^*$ : **optimal hypothesis** (minimizer of population loss over  $\mathcal{H}$  -infinite data sample)

$h_S^*$ : **empirically optimal hypothesis** (minimizer of sample loss over  $\mathcal{H}$ )

$h$ : **returned hypothesis**

**Note:** For sampling to give a good proxy, we must enforce the *consistency condition* in the infinite sample size limit. Namely,  $\lim_{m \rightarrow \infty} h_S^* = h_{\mathcal{D}}^*$ .

# Linear Regression: the simplest example of supervised learning

Suppose there's a god-given linear relationship between an input  $\mathbf{x} \in \mathbb{R}^d$  and the continuous output (label)  $y \in \mathbb{R}$  according to

$$y = f(\mathbf{x}) + \eta_i = \mathbf{w}_{true}^T \mathbf{x} + \eta_i ,$$

where  $\mathbf{w}_{true} \in \mathbb{R}^d$ , and  $\eta_i$  is an i.i.d. drawn from  $\mathcal{N}(0, \sigma^2)$ .

# Linear Regression: the simplest example of supervised learning

Suppose there's a god-given linear relationship between an input  $\mathbf{x} \in \mathbb{R}^d$  and the continuous output (label)  $y \in \mathbb{R}$  according to

$$y = f(\mathbf{x}) + \eta_i = \mathbf{w}_{true}^T \mathbf{x} + \eta_i ,$$

where  $\mathbf{w}_{true} \in \mathbb{R}^d$ , and  $\eta_i$  is an i.i.d. drawn from  $\mathcal{N}(0, \sigma^2)$ .

Suppose the *loss function* is the square error  $\ell(y, \hat{y}) = (y - \hat{y})^2$ , remember that the objective of SLT is to *search* for the hypothesis function  $h(\mathbf{x})$  that minimizes *the empirical risk*

$$L_S(h) := \frac{1}{m} \sum_{i=1}^m \ell \left( y_i, h(\mathbf{x}^{(i)}) \right),$$

given a *training set*  $S = \left\{ (\mathbf{x}^{(i)}, y_i) \right\}_{i=1}^m$ .

# Linear Regression: the simplest example of supervised learning

Suppose there's a god-given linear relationship between an input  $\mathbf{x} \in \mathbb{R}^d$  and the continuous output (label)  $y \in \mathbb{R}$  according to

$$y = f(\mathbf{x}) + \eta_i = \mathbf{w}_{true}^T \mathbf{x} + \eta_i ,$$

where  $\mathbf{w}_{true} \in \mathbb{R}^d$ , and  $\eta_i$  is an i.i.d. drawn from  $\mathcal{N}(0, \sigma^2)$ .

Suppose the *loss function* is the square error  $\ell(y, \hat{y}) = (y - \hat{y})^2$ , remember that the objective of SLT is to *search* for the hypothesis function  $h(\mathbf{x})$  that minimizes *the empirical risk*

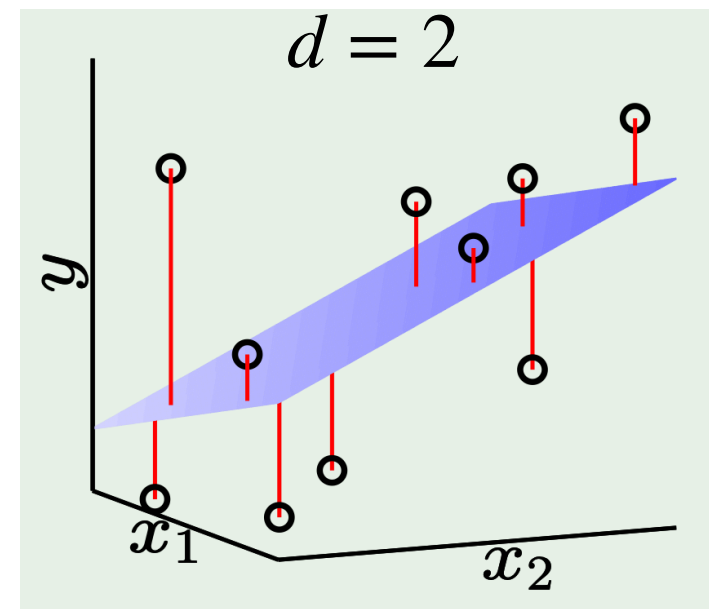
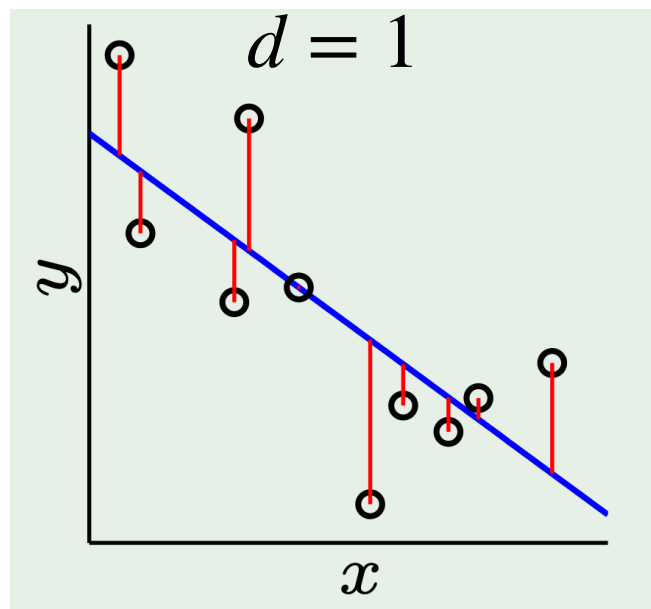
$$L_S(h) := \frac{1}{m} \sum_{i=1}^m \ell \left( y_i, h(\mathbf{x}^{(i)}) \right),$$

given a *training set*  $S = \left\{ (\mathbf{x}^{(i)}, y_i) \right\}_{i=1}^m$ .

For simplicity, suppose we restrict our hypothesis class to be the simplest class possible, i.e. a linear continuous real-valued function (The SAME class as the god-given function, can't be simpler than that!)

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

can we find a learning algorithm that spits out  $h_S^*$ ?



In other words, how do we find a  $d$ -dimensional hyperplane with a minimum square error from training data.

For brevity, let's rewrite the **empirical risk** as an  $L^2$  norm

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^T \mathbf{x}^{(i)} - y_i)^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2,$$

$$\mathbf{X} = \begin{bmatrix} -\mathbf{x}^{(1)\top} & - \\ -\mathbf{x}^{(2)\top} & - \\ \vdots & \\ -\mathbf{x}^{(d)\top} & - \end{bmatrix}$$

where for  $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$  the  $L^p$  norm of  $\mathbf{x}$  is defined as

$$\|\mathbf{x}\|_p = \left( |x_1|^p + \dots + |x_d|^p \right)^{\frac{1}{p}}.$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_d \end{bmatrix}$$



The learning algorithm should be able to spit out the weight  $\mathbf{w}$  that minimizes the **Empirical Risk Minimisation (ERM)** problem, which we'll call an **in-sample error**

$$E_{in}(\mathbf{w}^*) \equiv L_S(h_S^*) = \frac{1}{m} \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

The learning algorithm should be able to spit out the weight  $\mathbf{w}$  that minimizes the **Empirical Risk Minimisation (ERM)** problem, which we'll call an **in-sample error**

$$E_{in}(\mathbf{w}^*) \equiv L_S(h_S^*) = \frac{1}{n} \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

The necessary condition for  $\mathbf{w}$  to be a critical point is

$$\nabla E_{in}(\mathbf{w}^*) = \frac{2}{N} \mathbf{X}^\top (\mathbf{X}\mathbf{w}^* - \mathbf{y}) = \mathbf{0}$$

The learning algorithm should be able to spit out the weight  $\mathbf{w}$  that minimizes the **Empirical Risk Minimisation (ERM)** problem, which we'll call an **in-sample error**

$$E_{in}(\mathbf{w}^*) \equiv L_S(h_S^*) = \frac{1}{n} \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

The necessary condition for  $\mathbf{w}$  to be a critical point is

$$\nabla E_{in}(\mathbf{w}^*) = \frac{2}{N} \mathbf{X}^\top (\mathbf{X}\mathbf{w}^* - \mathbf{y}) = \mathbf{0}$$

In this case, the critical point turns out to be the minimizer (convex optimization)

$$\mathbf{w}^* = \mathbf{X}^\dagger \mathbf{y},$$

where we denote the **Moore-Penrose pseudo-inverse** of  $\mathbf{X}$  as

$$\mathbf{X}^\dagger = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

The learning algorithm should be able to spit out the weight  $\mathbf{w}$  that minimizes the **Empirical Risk Minimisation (ERM)** problem, which we'll call an **in-sample error**

$$E_{in}(\mathbf{w}^*) \equiv L_S(h_S^*) = \frac{1}{m} \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

The necessary condition for  $\mathbf{w}$  to be a critical point is

$$\nabla E_{in}(\mathbf{w}^*) = \frac{2}{N} \mathbf{X}^\top (\mathbf{X}\mathbf{w}^* - \mathbf{y}) = \mathbf{0}$$

In this case, the critical point turns out to be the minimizer (convex optimization)

$$\mathbf{w}^* = \mathbf{X}^\dagger \mathbf{y},$$

where we denote the **Moore-Penrose pseudo-inverse** of  $\mathbf{X}$  as

$$\mathbf{X}^\dagger = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

So we have the first **learning algorithm** to find the unique minimizer of the least-square linear regression problem in **one-shot!**

- Step 1: From training data, construct the *design matrix*  $\mathbf{X}$  and the vector  $\mathbf{y}$ .
- Step 2: Compute the Pseudo-inverse of the design matrix  $\mathbf{X}^\dagger = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$
- Step 3: Perform the matrix vector multiplication  $\mathbf{w}^* = \mathbf{X}^\dagger \mathbf{y}$  to obtain the **empirically optimal hypothesis**  $h_S^*(\mathbf{x}) = \mathbf{w}^{*T} \mathbf{x}$ .

Recall from Linear Algebra that invertibility of the matrix is always tricky, as one needs to check if the matrix has full-rank (typically the case when  $m \gg d$ )...

Recall from Linear Algebra that invertibility of the matrix is always tricky, as one needs to check if the matrix has full-rank (typically the case when  $m \gg d$ )...

However, assuming invertibility, here's an important result you'll prove in homework 1

$$\bar{E}_{in} \equiv \mathbb{E}_D \left[ E_{in}(\mathbf{w}_D^*) \right] = \sigma^2 \left( 1 - \frac{d}{m} \right) \quad \text{(in-sample error)}$$

$$\bar{E}_{out} \equiv \mathbb{E}_D \left[ E_{out}(\mathbf{w}_D^*) \right] = \sigma^2 \left( 1 + \frac{d}{m} \right) \quad \text{(out-of-sample error)}$$

Recall from Linear Algebra that invertibility of the matrix is always tricky, as one needs to check if the matrix has full-rank (typically the case when  $m \gg d$ )...

However, assuming invertibility, here's an important result you'll prove in homework 1

$$\bar{E}_{in} \equiv \mathbb{E}_D \left[ E_{in}(\mathbf{w}_D^*) \right] = \sigma^2 \left( 1 - \frac{d}{m} \right) \quad \text{(in-sample error)}$$

$$\bar{E}_{out} \equiv \mathbb{E}_D \left[ E_{out}(\mathbf{w}_D^*) \right] = \sigma^2 \left( 1 + \frac{d}{m} \right) \quad \text{(out-of-sample error)}$$

This gives us the *generalization error*

$$| \bar{E}_{out} - \bar{E}_{in} | = 2\sigma^2 \left( \frac{d}{m} \right)$$

Recall from Linear Algebra that invertibility of the matrix is always tricky, as one needs to check if the matrix has full-rank (typically the case when  $m \gg d$ )...

However, assuming invertibility, here's an important result you'll prove in homework 1

$$\bar{E}_{in} \equiv \mathbb{E}_D \left[ E_{in}(\mathbf{w}_D^*) \right] = \sigma^2 \left( 1 - \frac{d}{m} \right) \quad (\text{in-sample error})$$
$$\bar{E}_{out} \equiv \mathbb{E}_D \left[ E_{out}(\mathbf{w}_D^*) \right] = \sigma^2 \left( 1 + \frac{d}{m} \right) \quad (\text{out-of-sample error})$$

This gives us the *generalization error*

$$| \bar{E}_{out} - \bar{E}_{in} | = 2\sigma^2 \left( \frac{d}{m} \right)$$

- Large generalization error if  $d \gg m$ . Make sense since you'll likely fit a noisy subspace of the actual high-dimensional hyperplane.
- Even when  $d \approx m$ , noise suppresses learning the actual high-dimensional hyperplane.
- We'll learn how to “*regularise*” learning to not be too sensitive to noise (lower the variance of the learned models).



Schematic of how in-sample error and out-of-sample error are related in linear least square

