

Machine Learning I

Supervised learning framework

Souhaib Ben Taieb

March 2, 2022

University of Mons

Table of contents

Components of supervised learning

Examples of learning models

A note on E_{in} and E_{out} of MLE

Training and test errors

Model selection and assessment

Estimation of the out-of-sample error with resampling methods

Table of contents

Components of supervised learning

Examples of learning models

A note on E_{in} and E_{out} of MLE

Training and test errors

Model selection and assessment

Estimation of the out-of-sample error with resampling methods

Input and output variables

The **input variables**¹ are typically denoted using the symbol X . If we observe p different variables, we write $X = (X_1, X_2, \dots, X_p)$. The inputs belong to an *input space* \mathcal{X} .

- Examples: $\mathcal{X} \subseteq \mathbb{R}^p$ or $\mathcal{X} = \{0, 1\}^p$.

The **output variable**² is typically denoted using the symbol Y . The output belongs to an *output space* \mathcal{Y} .

- Regression: $\mathcal{Y} \subseteq \mathbb{R}$
- Classification (with K categories): $\mathcal{Y} = \{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_K\}$
 - Binary classification ($K = 2$): $\mathcal{Y} = \{-1, 1\}$ or $\mathcal{Y} = \{0, 1\}$

¹also called *predictors*, *independent variables*, *features*, *variables* or just *inputs*.

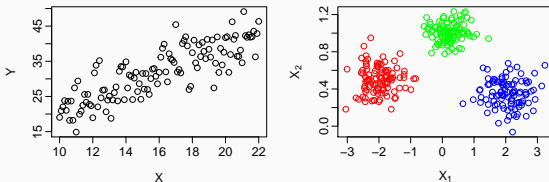
²also called the *response* or *dependent variable*.

The dataset

The **dataset**, also called *training set*, is a set of n input-output pairs , given by

$$\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\} = \{(x_i, y_i)\}_{i=1}^n.$$

Each pair, also called an *example* or a *data point*, belongs to the *data space* $\mathcal{X} \times \mathcal{Y}$.



- Left figure: $\mathcal{X} \subseteq \mathbb{R}$ ($p = 1$) and $\mathcal{Y} \subseteq \mathbb{R}$
- Right figure: $\mathcal{X} = \mathbb{R}^2$ ($p = 2$) and $\mathcal{Y} = \{R, G, B\}$

Data distribution

We assume the data points (x_i, y_i) are i.i.d. realizations from a fixed unknown **data distribution** $p_{X,Y}$, which represents different *sources of uncertainty*. In other words, we have $(X, Y) \sim p_{X,Y}$.

The probability distribution $p_{X,Y}$ can be factorized as

$$p_{X,Y}(x, y) = p_X(x)p_{Y|X}(y|x),$$

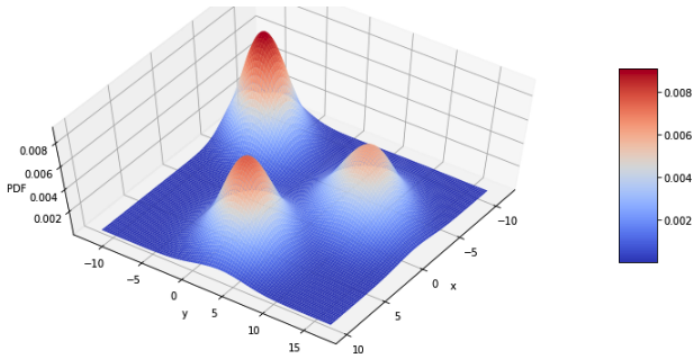
where

- the marginal distribution p_X models uncertainty in the sampling of the inputs.
- the conditional distribution $p_{Y|X}$ describes a stochastic (non-deterministic) relation between inputs and output.

Equivalently, we have

$$X \sim p_X \text{ and } Y|X = x \sim p_{Y|X}(\cdot|x).$$

Data distribution



Source: <https://tinyurl.com/19bdt531>

The loss function

Given a prediction $\hat{y} \in \mathcal{Y}$ and the true (observed) value $y \in \mathcal{Y}$, the **loss function**

$$L : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, \infty),$$

measures how far \hat{y} is from y .

In the continuous case, examples include the squared error loss

$$L(y, \hat{y}) = (y - \hat{y})^2,$$

and the absolute error loss

$$L(y, \hat{y}) = |y - \hat{y}|.$$

In the discrete case, an example is the zero-one loss

$$L(y, \hat{y}) = \mathbb{1}\{y \neq \hat{y}\},$$

where $\mathbb{1}\{\cdot\}$ is the indicator function.

The supervised learning problem

Let us consider a **hypothesis** (“prediction function”)

$$h : \mathcal{X} \rightarrow \mathcal{Y}.$$

The supervised learning problem

Let us consider a **hypothesis** (“prediction function”)

$$h : \mathcal{X} \rightarrow \mathcal{Y}.$$

We define the **out-of-sample error**³ of h as

$$E_{\text{out}}(h) = \mathbb{E}_{x,y}[L(y, h(x))] = \mathbb{E}_x \left[\underbrace{\mathbb{E}_{y|x}[L(y, h(x))|x]}_{E_{\text{out}}(h,x)} \right] \quad (1)$$

The **optimal prediction function** is given by

$$f = \operatorname{argmin}_{h:\mathcal{X} \rightarrow \mathcal{Y}} E_{\text{out}}(h).$$

³also called **expected error** or **expected risk**.

The supervised learning problem

$$f = \operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_{\mathbf{x}, \mathbf{y}}[L(\mathbf{y}, h(\mathbf{x}))] := E_{\text{out}}(h).$$

1. We do not know $p_{\mathcal{X}, \mathcal{Y}}$
2. We search among all functions $h : \mathcal{X} \rightarrow \mathcal{Y}$

The supervised learning problem

Let \mathcal{H} be a **hypothesis set**, i.e. a set of prediction function (hypotheses) under consideration.

Given a dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ where $(x_i, y_i) \stackrel{\text{i.i.d.}}{\sim} p_{X,Y}$, we can solve the following optimization problem:

$$g (= g_{\mathcal{D}}) = \operatorname{argmin}_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n L(y_i, h(x_i)) := E_{\text{in}}(h),$$

where E_{in} is the **in-sample error** or **training error**.

Compare it with

$$f = \operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_{x,y} [L(y, h(x))] := E_{\text{out}}(h).$$

Summary

The **dataset** $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ is composed of n input-output pairs (x_i, y_i) , which are i.i.d. realizations from the **data distribution** $p_{X,Y}$ where $X \in \mathcal{X}$ and $Y \in \mathcal{Y}$.

The **loss function** $L : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, \infty)$ allows us to measure the error we incur in predicting \hat{y} in place of y .

The **hypothesis set** \mathcal{H} is a set of prediction function under consideration. Each hypothesis $h \in \mathcal{H}$ has an **in-sample error** $E_{\text{in}}(h)$, computed using \mathcal{D} , and an **out-of-sample error** $E_{\text{out}}(h)$, which depends on $p_{X,Y}$.

Given \mathcal{H} and using \mathcal{D} , the **learning algorithm** \mathcal{A} picks the best hypothesis g from \mathcal{H} according to the loss function L .

Together, the hypothesis set and the learning algorithm are referred to as the **learning model**.

Table of contents

Components of supervised learning

Examples of learning models

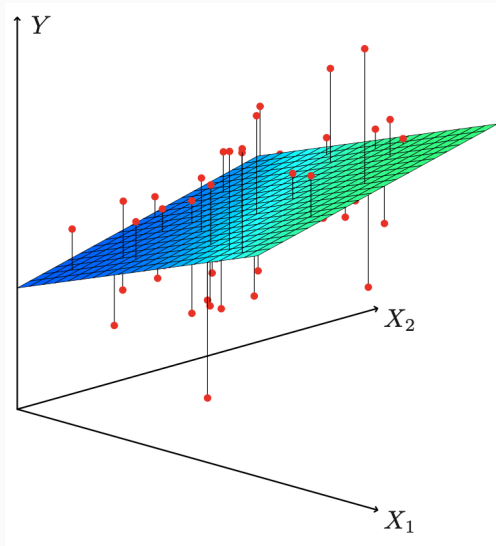
A note on E_{in} and E_{out} of MLE

Training and test errors

Model selection and assessment

Estimation of the out-of-sample error with resampling methods

Linear models



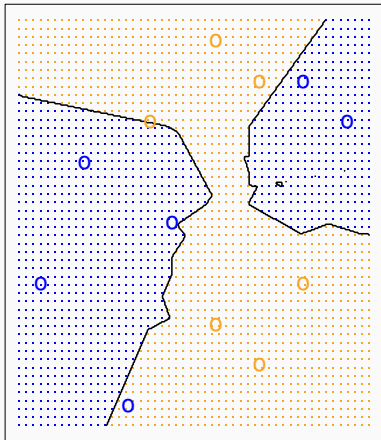
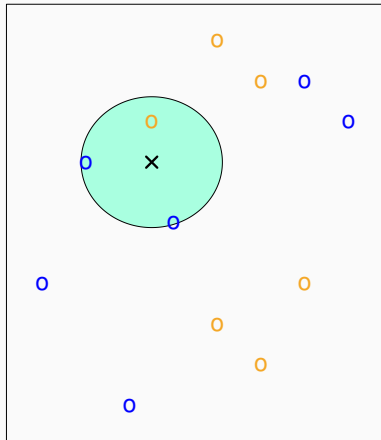
Let us consider a regression problem where $x \in \mathbb{R}^p$.

- The **hypothesis set** for linear models is given by

$$\mathcal{H} = \{h(x) = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p : \beta_0, \beta_1, \dots, \beta_p \in \mathbb{R}\}$$

- One **learning algorithm** is the *(ordinary) least squares* method.

K-Nearest Neighbours model



K-Nearest Neighbours (KNN) is one of the simplest machine learning model for both classification and regression.

K-Nearest Neighbours model

Given a dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$, the number of neighbors $K \leq n$, and a new input x_* :

(1) Find the K nearest points to x_* in \mathcal{D} , denoted \mathcal{N}_* .

- Classification with M categories ($y \in \{C_1, \dots, C_M\}$)

(2)

$$\hat{p}(C_m \mid x = x_*) \approx \frac{1}{K} \sum_{i \in \mathcal{N}_*} \mathbb{1}\{y_i = C_m\} \quad (m = 1, 2, \dots, M)$$

$$h(x_*) = C_{m^*} \text{ where } m^* = \operatorname{argmax}_m \hat{p}(y = C_m \mid x = x_*)$$

- Regression

(2)

$$h(x_*) = \frac{1}{K} \sum_{i \in \mathcal{N}_*} y_i$$

Parametric and non-parametric models

- In a **parametric model**, every hypothesis is uniquely defined by a **fixed number of parameters**.
- In a **non-parametric model**, we can not describe a hypothesis with a fixed number of parameters. Usually the number of “parameters” **grows with the size of the dataset**.
- Both parametric and non-parametric models have **hyper-parameters** (structural parameters), while parametric models also have parameters
 - For KNN, K is a hyper-parameter.
 - For linear models, p is a hyper-parameter, and the coefficients β_j are parameters.
- We also make a distinction between **linear** and **non-linear** models.

Table of contents

Components of supervised learning

Examples of learning models

A note on E_{in} and E_{out} of MLE

Training and test errors

Model selection and assessment

Estimation of the out-of-sample error with resampling methods

Recall that, given a set of possible distributions

$$\mathcal{H} = \{p(y; \theta) : \theta \in \Theta\},$$

the maximum likelihood estimator is given by

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmax}} \sum_{i=1}^n \log p(y_i; \theta) = \underset{\theta \in \Theta}{\operatorname{argmin}} \underbrace{\frac{1}{n} \sum_{i=1}^n -\log p(y_i; \theta)}_{E_{\text{in}}(\theta)}. \quad (2)$$

Since each hypothesis h is completely characterized by θ , this is equivalent to

$$g = \underset{h \in \mathcal{H}}{\operatorname{argmin}} E_{\text{in}}(h).$$

By (strong) law of large numbers,

$$\frac{1}{n} \sum_{i=1}^n -\log p(y_i; \boldsymbol{\theta}) \xrightarrow{n \rightarrow \infty} \mathbb{E}[-\log p(y; \boldsymbol{\theta})]$$

In other words, we can think of maximum likelihood estimation as trying to minimize

$$E_{\text{out}}(\boldsymbol{\theta}) = \mathbb{E}[-\log p(y; \boldsymbol{\theta})]$$

$$\operatorname{argmin}_{\theta \in \Theta} E_{\text{out}}(\theta) = \operatorname{argmin}_{\theta \in \Theta} \mathbb{E}[-\log p(y; \theta)] \quad (3)$$

$$= \operatorname{argmin}_{\theta \in \Theta} \mathbb{E}[\log p(y) - \log p(y; \theta)] \quad (4)$$

$$= \operatorname{argmin}_{\theta \in \Theta} \mathbb{E} \left[\log \frac{p(y)}{p(y; \theta)} \right] \quad (5)$$

$$= \operatorname{argmin}_{\theta \in \Theta} \int \log \frac{p(y)}{p(y; \theta)} p(y) \, dy \quad (6)$$

$$= \operatorname{argmin}_{\theta \in \Theta} \text{KL}(p_{\theta}, p), \quad (7)$$

where $\text{KL}(q, p)$ is the KL-divergence between two distributions q and p , which measures the discrepancy between the two distributions. Note that KL-divergence is not a distance measure (not symmetric).

Table of contents

Components of supervised learning

Examples of learning models

A note on E_{in} and E_{out} of MLE

Training and test errors

Model selection and assessment

Estimation of the out-of-sample error with resampling methods

Training and test errors

Consider

$$f = \operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{Y}} E_{\text{out}}(h),$$

$$g^* = \operatorname{argmin}_{h \in \mathcal{H}} E_{\text{out}}(h),$$

and

$$g = \operatorname{argmin}_{h \in \mathcal{H}} E_{\text{in}}(h).$$

How different are $E_{\text{out}}(g)$ and $E_{\text{out}}(f)$?

The approximation-generalization tradeoff

The difference between the out-of-sample error of g and f can be decomposed as follows

$$E_{\text{out}}(g) - E_{\text{out}}(f) = \underbrace{[E_{\text{out}}(g^*) - E_{\text{out}}(f)]}_{\text{Approximation error}} + \underbrace{[E_{\text{out}}(g) - E_{\text{out}}(g^*)]}_{\text{Estimation error}}$$

- **Approximation error** is how far the entire hypothesis set is from f . Larger hypothesis sets have lower approximation error.
- **Estimation error** is how good g is with respect to the best in the hypothesis set. Larger hypothesis sets have higher estimation error because it is harder to find a good prediction function based on limited data.

This is called the **approximation-generalization** tradeoff.

How does E_{out} relates to E_{in} ?

$$E_{\text{out}}(g) = E_{\text{in}}(g) + [E_{\text{out}}(g) - E_{\text{in}}(g)].$$

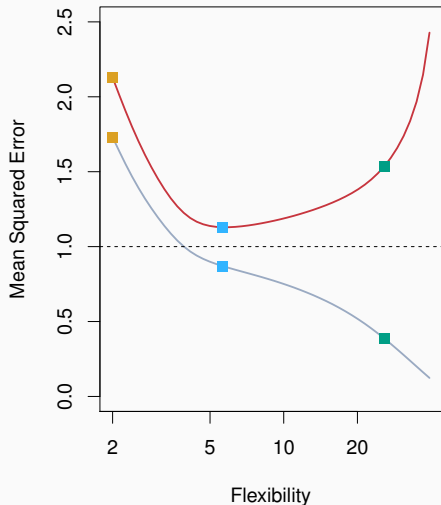
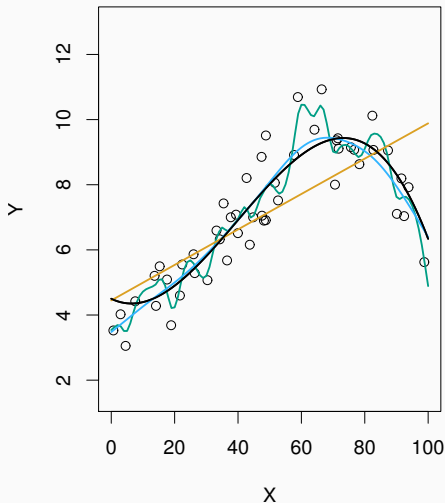
To obtain a small $E_{\text{out}}(g)$, we want

1. small $E_{\text{in}}(g)$
2. small $[E_{\text{out}}(g) - E_{\text{in}}(g)]$

In the following, we will see that the selection of the best hypothesis g by minimizing $E_{\text{in}}(g)$ can be misleading. Let $\mathcal{D}' = \{(x'_i, y'_i)\}_{i=1}^{n'}$ be another sample (**independent** of \mathcal{D}) where $(x'_i, y'_i) \stackrel{\text{i.i.d.}}{\sim} p_{X,Y}$. We define the **testing/test error** as

$$E_{\text{test}}(g) = \frac{1}{n'} \sum_{i=1}^{n'} L(y'_i, g(x'_i)).$$

Training and test errors in regression



Black: true curve

Orange: linear regression

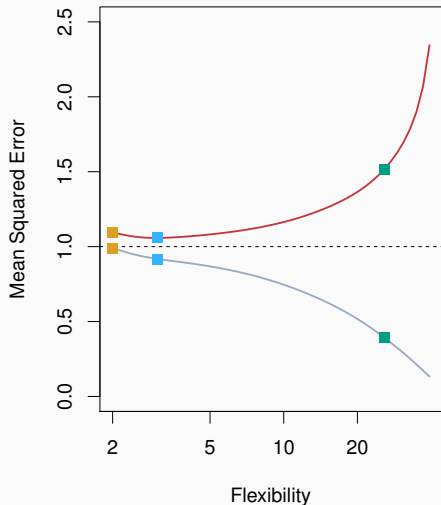
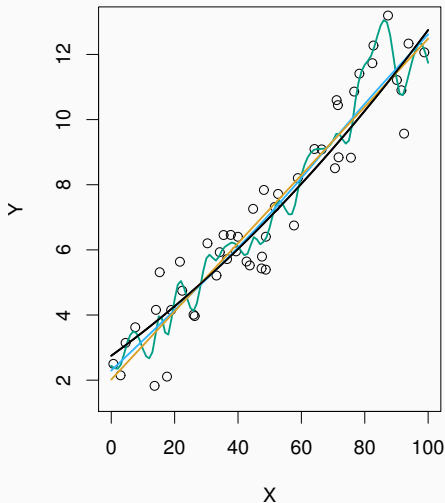
Blue/green: nonlinear regression

Grey: Training MSE

Red: Test MSE

Dashed: Minimum test MSE₇

Training and test errors in regression



Black: true curve

Orange: linear regression

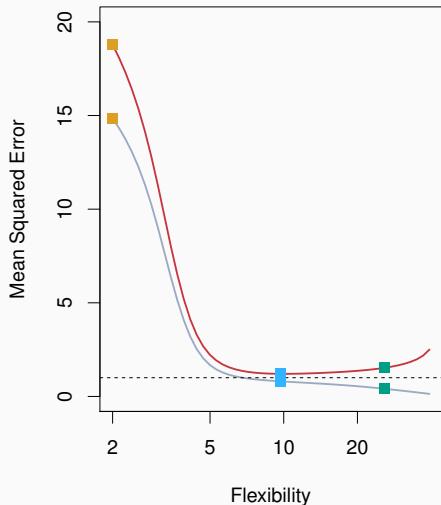
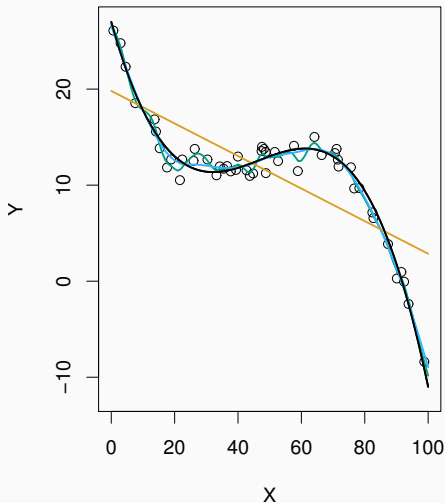
Blue/green: nonlinear regression

Grey: Training MSE

Red: Test MSE

Dashed: Minimum test MSE

Training and test errors in regression



Black: true curve

Orange: linear regression

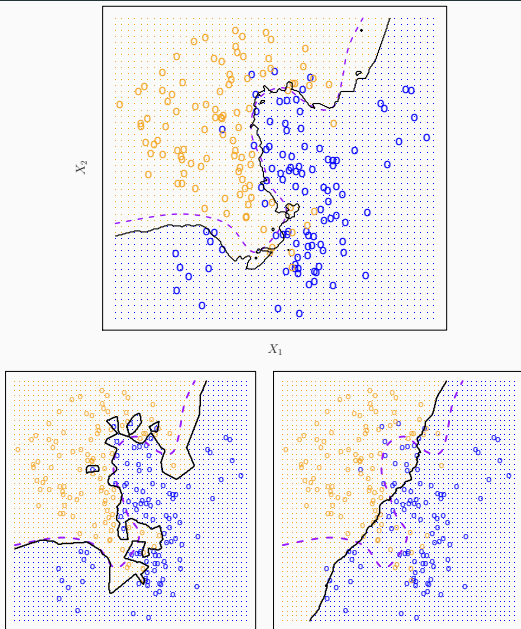
Blue/green: nonlinear regression

Grey: Training MSE

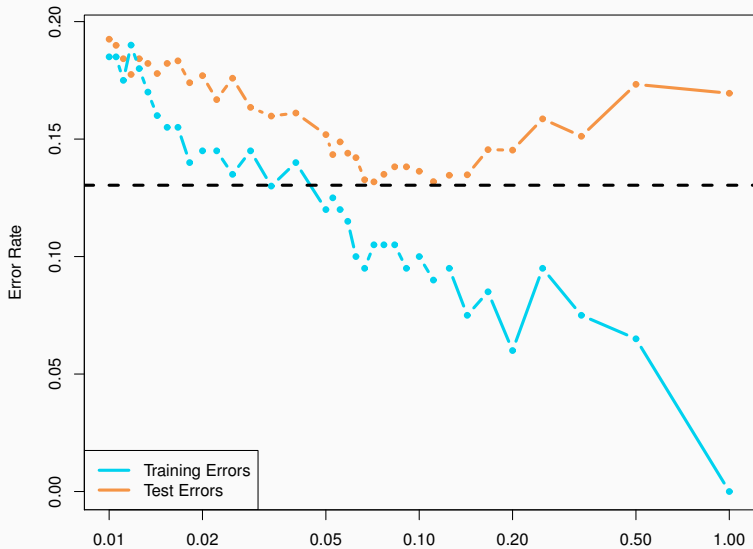
Red: Test MSE

Dashed: Minimum test MSE

Training and test errors in classification



Training and test errors in classification



A fundamental picture

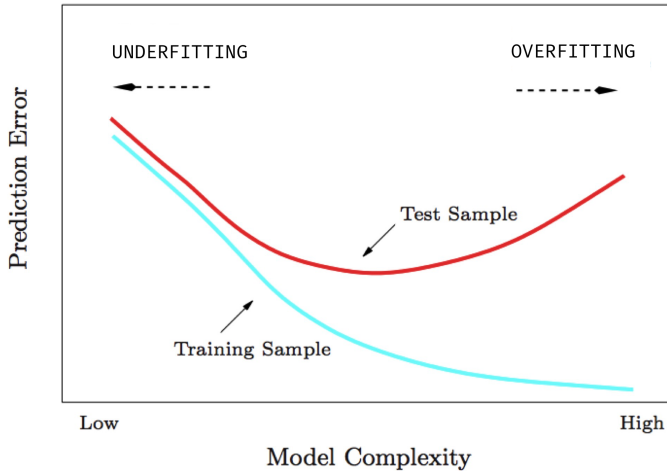


Table of contents

Components of supervised learning

Examples of learning models

A note on E_{in} and E_{out} of MLE

Training and test errors

Model selection and assessment

Estimation of the out-of-sample error with resampling methods

The model selection problem

- How many and which input variables should we choose?
 - How do we choose the number of neighbours K in KNN?
 - ...
 - More generally, how do we choose the **hyper-parameters** of a machine learning model?
- We need a way of **assessing** and **selecting** the best model among multiple competing models.

The model selection/structural identification procedure

1. Model generation

Generate a set of **candidate model structures** among which the best one is to be selected. If applicable, estimate the parameters of each candidate model, i.e. fit the model by minimizing the **in-sample (training) error** (*parametric identification*).

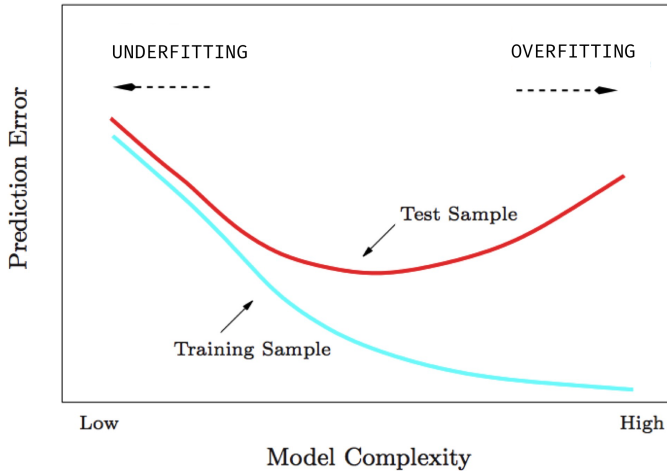
2. Model assessment/validation

Evaluate/validate the model's performance by computing a **validation error**, i.e. an estimate of the out-of-sample error.

3. Model selection

Select the final model structure in the set that has been proposed by **model generation** and assessed by **model validation**. We typically select the model structure that minimizes the **validation error**.

Why not use the in-sample (training) error for model selection?



How to estimate the out-of-sample error?

$$E_{\text{out}}(h) = E_{\text{in}}(h) + \underbrace{[E_{\text{out}}(h) - E_{\text{in}}(h)]}_{\text{overfit penalty}}, \quad h \in \mathcal{H}.$$

1. **Directly estimate it** using resampling methods, i.e. by resampling the data set. Examples include **the validation set**, **cross-validation** and **bootstrap**.
2. **Estimate the overfit penalty** and **add it to the in-sample (training) error**.

Table of contents

Components of supervised learning

Examples of learning models

A note on E_{in} and E_{out} of MLE

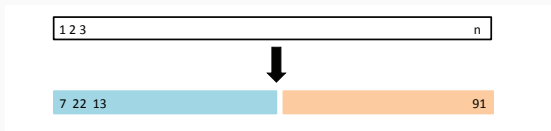
Training and test errors

Model selection and assessment

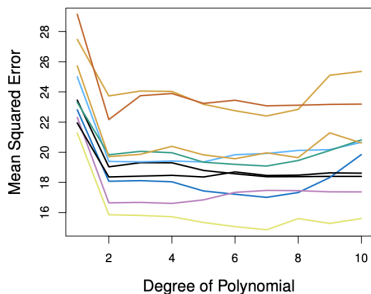
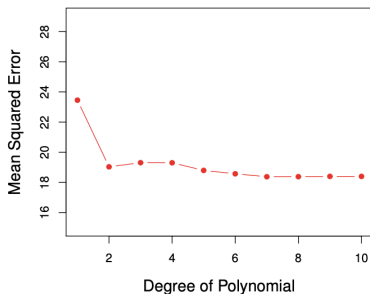
Estimation of the out-of-sample error with resampling methods

Validation-set approach

- We randomly divide the dataset into two parts: a **training set** and a **validation set**.
- The model is fit on the training set, and the fitted model is used to predict the responses for the observations in the validation set
- The resulting **validation-set error** provides an estimate of the **out-of-sample error**.
- This estimate of the out-of-sample error can be highly variable, depending on precisely which observations are included in the training set and the validation set.

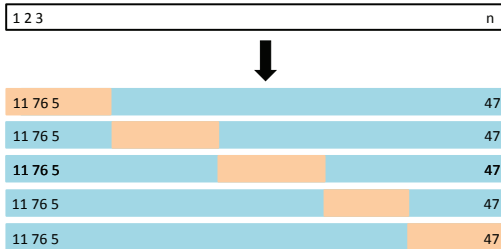


Validation-set approach



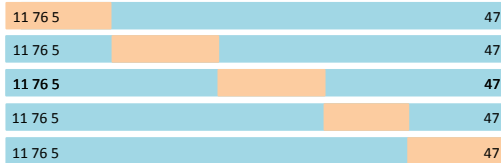
Left panel shows single split; right panel shows multiple splits

K-fold Cross-validation



K-fold Cross-validation

1 2 3 n



1 2 3 n



1 2 3 n



K -fold Cross-validation

- Divide the data set into K different parts.
- Remove one part, fit the model on the remaining $K - 1$ parts, and compute the prediction error on the omitted part.
- Repeat K times taking out a different part each time
- By averaging the K prediction errors we obtain an estimate of the out-of-sample error, i.e. the prediction error for new observations (not used in training).
- Setting $K = n$ yields n -fold or **leave-one out cross-validation** (LOOCV).

K-fold Cross-validation

Let the K parts be A_1, A_2, \dots, A_K , where A_k denotes the indices of the observations in part k . There are n_k observations in part k : if N is a multiple of K , then $n_k = n/K$.

Compute

$$\text{CV} = \sum_{k=1}^K \frac{n_k}{n} \text{Err}_k$$

where, for example,

$$\text{Err}_k = \sum_{i \in A_k} \frac{1}{n_k} (y_i - \hat{y}_i^{(-k)})^2 \quad (\text{squared error loss}),$$

or

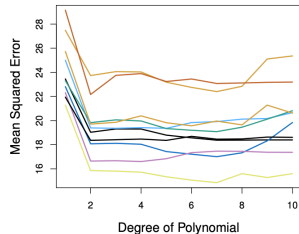
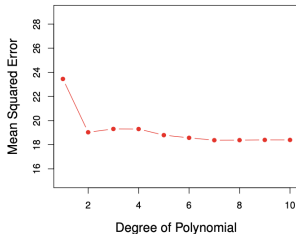
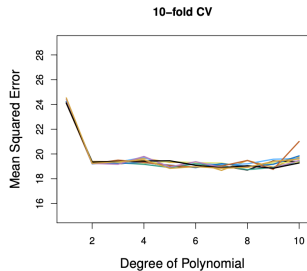
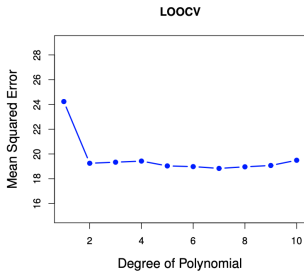
$$\text{Err}_k = \sum_{i \in A_k} \frac{1}{n_k} \mathbb{1}\{y_i \neq \hat{y}_i^{(-k)}\} \quad (\text{zero-one loss}),$$

and $\hat{y}_i^{(-k)}$ is the prediction for observation i , obtained from the model fit with data where part k is removed.

K -fold Cross-validation - Bias and variance tradeoff

- Each training set is only $(K - 1)/K$ as big as the original data set. So the estimates of prediction error will be biased upwards.
- Bias minimized when $K = n$ (LOOCV).
- But variance increases with K since the estimates from each fold are more correlated (as there are more overlapping observations in each part) and hence their average can have higher variance.
- Empirical observation: $K = 5$ or $K = 10$ provide a good compromise for this bias-variance tradeoff.

10-fold Cross-validation



Left panel shows single split; right panel shows multiple splits

K-fold Cross-validation

For K-fold cross-validation, it's very helpful to assign a quantitative notion of variability to the cross-validation error estimate.

Assuming $n_k = n/K$, we can write

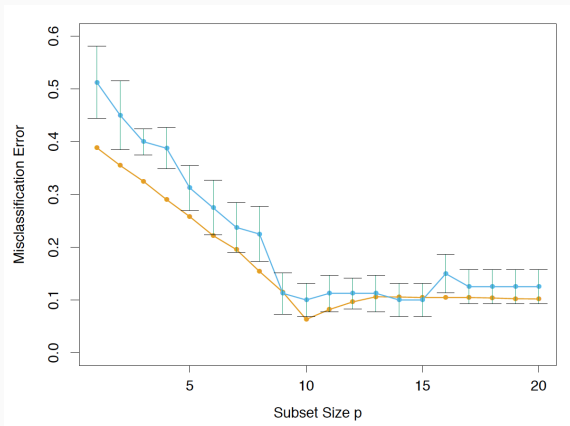
$$\text{Var}(CV) = \text{Var}\left(\sum_{k=1}^K \frac{1}{K} \text{Err}_k\right) \approx \frac{1}{K^2} \sum_{k=1}^K \text{Var}(\text{Err}_1) = \frac{1}{K} \text{Var}(\text{Err}_1)$$

- This is an approximation since $\text{Err}_1, \dots, \text{Err}_K$ are not i.i.d.
- This approximation is valid for small K (e.g., $K = 5$ or 10) but not really for big K (e.g., $K = n$), because then the quantities $\text{Err}_1, \dots, \text{Err}_K$ are highly correlated.
- We can compute the standard deviation or standard error of the cross-validation error estimate as

$$\frac{1}{\sqrt{K}} \text{sd}\{\text{Err}_1, \text{Err}_2, \dots, \text{Err}_K\},$$

where *sd* denotes the *empirical standard deviation*.

The one standard error rule



Choose the simplest model whose CV error is no more than one standard error above the model with the lowest CV error

Cross-validation: right and wrong

Consider a simple regression procedure applied to a dataset with 500 predictors and 50 samples:

1. Find the 5 predictors having the largest correlation with the response
2. Apply linear regression using only these 5 predictors

How to use cross-validation to estimate the out-of-sample error of this procedure?

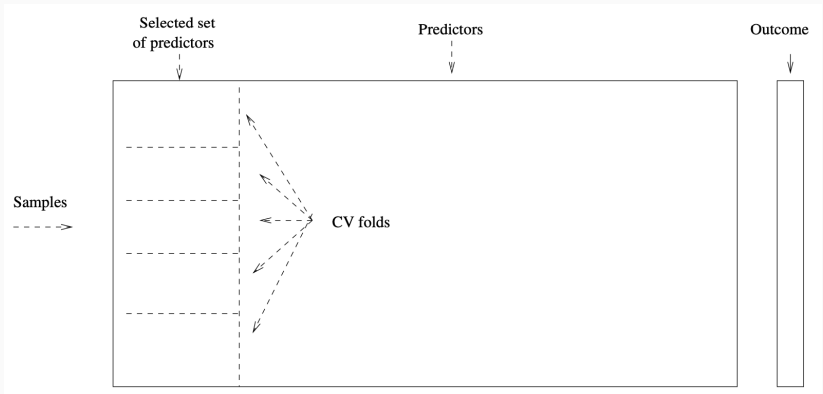
Cross-validation: right and wrong

Consider a simple regression procedure applied to a dataset with 500 predictors and 50 samples:

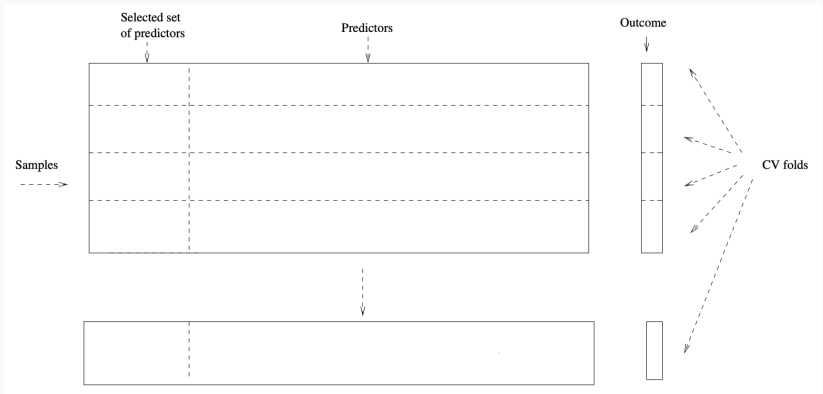
1. Find the 5 predictors having the largest correlation with the response
2. Apply linear regression using only these 5 predictors

How to use cross-validation to estimate the out-of-sample error of this procedure?

Wrong way



Right way



Cross-validation: right and wrong

Apply cross-validation to steps 1 **and** 2 (not just step 2)

→ Every aspect of the procedure that involves using the data — variable selection, scaling, etc — must be cross-validated.