

Machine Learning

Supervised Learning. Features. Loss Functions. Cross-validation

Aleksandr Petiushko

ML Research



① Supervised Learning Setting

Content

- ➊ Supervised Learning Setting
- ➋ Objects' features

Content

- ➊ Supervised Learning Setting
- ➋ Objects' features
- ➌ Model outputs

Content

- ➊ Supervised Learning Setting
- ➋ Objects' features
- ➌ Model outputs
- ➍ Loss functions

Content

- ➊ Supervised Learning Setting
- ➋ Objects' features
- ➌ Model outputs
- ➍ Loss functions
- ➎ Cross-validation

Content

- ➊ Supervised Learning Setting
- ➋ Objects' features
- ➌ Model outputs
- ➍ Loss functions
- ➎ Cross-validation
- ➏ Hyperparameters tuning

Machine Learning Paradigms: reminder

Definitions

- X — set of objects
- Y — set of (correct) answers/labels
- $y : X \rightarrow Y$ — the unknown dependency

Machine Learning Paradigms: reminder

Definitions

- X — set of objects
- Y — set of (correct) answers/labels
- $y : X \rightarrow Y$ — the unknown dependency

Machine learning paradigms

- **Supervised** (now)
 - Sufficient amount of training material, i.e. pairs (x_i, y_i)

Machine Learning Paradigms: reminder

Definitions

- X — set of objects
- Y — set of (correct) answers/labels
- $y : X \rightarrow Y$ — the unknown dependency

Machine learning paradigms

- **Supervised** (now)
 - Sufficient amount of training material, i.e. pairs (x_i, y_i)
- Semi-supervised
 - Few labeled data (x_i, y_i) and many unlabeled examples x_j

Machine Learning Paradigms: reminder

Definitions

- X — set of objects
- Y — set of (correct) answers/labels
- $y : X \rightarrow Y$ — the unknown dependency

Machine learning paradigms

- **Supervised** (now)
 - Sufficient amount of training material, i.e. pairs (x_i, y_i)
- Semi-supervised
 - Few labeled data (x_i, y_i) and many unlabeled examples x_j
- *Unsupervised* (in future lectures?)
 - No labeled pairs, only x_i examples

Machine Learning Paradigms: reminder

Definitions

- X — set of objects
- Y — set of (correct) answers/labels
- $y : X \rightarrow Y$ — the unknown dependency

Machine learning paradigms

- **Supervised** (now)
 - Sufficient amount of training material, i.e. pairs (x_i, y_i)
- Semi-supervised
 - Few labeled data (x_i, y_i) and many unlabeled examples x_j
- *Unsupervised* (in future lectures?)
 - No labeled pairs, only x_i examples
- Reinforced
 - Action generation based on interaction with the environment

Problem setting for supervised learning

- Given:
 - $\{(x_1, y_1), \dots, (x_n, y_n)\} \subset X \times Y$ — training set

Problem setting for supervised learning

- Given:
 - $\{(x_1, y_1), \dots, (x_n, y_n)\} \subset X \times Y$ — training set
- Find
 - A decision function $a : X \rightarrow Y$ that approximates the target dependency y .

Problem setting for supervised learning

- Given:
 - $\{(x_1, y_1), \dots, (x_n, y_n)\} \subset X \times Y$ — training set
- Find
 - A decision function $a : X \rightarrow Y$ that approximates the target dependency y .
- Need to clarify:

Problem setting for supervised learning

- Given:
 - $\{(x_1, y_1), \dots, (x_n, y_n)\} \subset X \times Y$ — training set
- Find
 - A decision function $a : X \rightarrow Y$ that approximates the target dependency y .
- Need to clarify:
 - How objects are defined

Problem setting for supervised learning

- Given:
 - $\{(x_1, y_1), \dots, (x_n, y_n)\} \subset X \times Y$ — training set
- Find
 - A decision function $a : X \rightarrow Y$ that approximates the target dependency y .
- Need to clarify:
 - How objects are defined
 - How answers are given

Problem setting for supervised learning

- Given:
 - $\{(x_1, y_1), \dots, (x_n, y_n)\} \subset X \times Y$ — training set
- Find
 - A decision function $a : X \rightarrow Y$ that approximates the target dependency y .
- Need to clarify:
 - How objects are defined
 - How answers are given
 - What does it mean that one dependency approximates another

How objects are defined

Definition

Object = set of features

How objects are defined

Definition

Object = set of features

Feature types

- Categorical feature

How objects are defined

Definition

Object = set of features

Feature types

- Categorical feature
- Binary attribute
 - ▶ A special case of categorical, when category = “does this property exist or not”

How objects are defined

Definition

Object = set of features

Feature types

- Categorical feature
- Binary attribute
 - ▶ A special case of categorical, when category = “does this property exist or not”
- Ordinal attribute
 - ▶ Full (or partial) order within categories

How objects are defined

Definition

Object = set of features

Feature types

- Categorical feature
- Binary attribute
 - ▶ A special case of categorical, when category = “does this property exist or not”
- Ordinal attribute
 - ▶ Full (or partial) order within categories
- Quantitative attribute

How answers are given

Classification tasks

- Binary classification $Y = \{-1, 1\}$ or $Y = \{0, 1\}$

How answers are given

Classification tasks

- Binary classification $Y = \{-1, 1\}$ or $Y = \{0, 1\}$
- Multiclass classification $Y = \{0, 1, \dots, M - 1\}$

How answers are given

Classification tasks

- Binary classification $Y = \{-1, 1\}$ or $Y = \{0, 1\}$
- Multiclass classification $Y = \{0, 1, \dots, M - 1\}$
- Multivalued binary classification $Y = \{0, 1\}^M$

How answers are given

Classification tasks

- Binary classification $Y = \{-1, 1\}$ or $Y = \{0, 1\}$
- Multiclass classification $Y = \{0, 1, \dots, M - 1\}$
- Multivalued binary classification $Y = \{0, 1\}^M$

Regression Tasks

$Y = \mathbb{R}$ or $Y = \mathbb{R}^n$

Loss Function

Definition

Loss function $L(a, x)$ — error value of algorithm a on object x

Loss Function

Definition

Loss function $L(a, x)$ — error value of algorithm a on object x

Loss functions for classification problems

$L(a, x) = [a(x) \neq y]$ — error indicator function (either 0 or 1)

Loss Function

Definition

Loss function $L(a, x)$ — error value of algorithm a on object x

Loss functions for classification problems

$L(a, x) = [a(x) \neq y]$ — error indicator function (either 0 or 1)

Loss functions for regression problems

$L(a, x) = (a(x) - y)^2$ — squared error

Comparison of machine learning models

How do you know that one model is better than another?

To do this, we use a set which is independent of **training** set, which is called **test** set

Comparison of machine learning models

How do you know that one model is better than another?

To do this, we use a set which is independent of **training** set, which is called **test** set

Why even bother with this?

- There are many machine learning algorithms and it is important to understand which one is more applicable to a particular task

Comparison of machine learning models

How do you know that one model is better than another?

To do this, we use a set which is independent of **training** set, which is called **test** set

Why even bother with this?

- There are many machine learning algorithms and it is important to understand which one is more applicable to a particular task
- Even within the same model, there can be many (hyper)parameters to choose from

How to choose the best model

Naive approach

Train models with different parameters and choose the best one on the test

How to choose the best model

Naive approach

Train models with different parameters and choose the best one on the test

Disadvantages of the naive approach

- Since the test usually consists of a random subset of the original sample, the result on the test is also some approximation of a random variable

How to choose the best model

Naive approach

Train models with different parameters and choose the best one on the test

Disadvantages of the naive approach

- Since the test usually consists of a random subset of the original sample, the result on the test is also some approximation of a random variable
- If all models are tested on a test dataset and thus choose the best one, then implicit training will occur on the test, and surprises are possible on another independent test

How to choose the best model

Naive approach

Train models with different parameters and choose the best one on the test

Disadvantages of the naive approach

- Since the test usually consists of a random subset of the original sample, the result on the test is also some approximation of a random variable
- If all models are tested on a test dataset and thus choose the best one, then implicit training will occur on the test, and surprises are possible on another independent test

So... what to do?

In order not to implicitly learn from test data — you need to use **cross-validation**

Cross Validation

General idea

The main idea of cross-validation is to split the training set into two non-overlapping sets (possibly multiple times):

$$X^{learn} = X^{train} \sqcup X^{val}$$

On one of them, training takes place, and on the other, the model is validated.

Cross Validation

General idea

The main idea of cross-validation is to split the training set into two non-overlapping sets (possibly multiple times):

$$X^{learn} = X^{train} \sqcup X^{val}$$

On one of them, training takes place, and on the other, the model is validated.

Why validate?

Usually, any machine learning algorithm contains a whole set of so-called “**hyperparameters**” (i.e. parameters that are not learned, but set initially): dimension, various weighting factors, etc.

And in order to select these parameters “fairly”, without using any test data at all, a validation procedure is carried out.

Cross Validation

Special cases

- ① The simplest cross-validation is **hold-out** control, in which the set is split once:



The diagram illustrates a hold-out cross-validation split. It consists of two adjacent rectangular boxes. The left box is red and contains the word 'Train' in black text. The right box is blue and contains the word 'Validation' in black text. These two boxes are positioned on a light blue background that is part of a larger slide container.

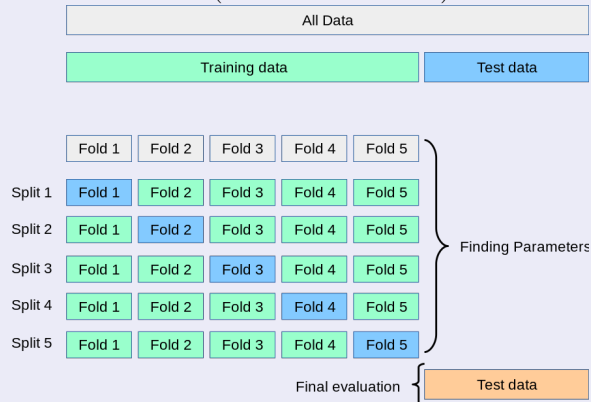
Train

Validation

Cross Validation

Special cases

② k block control (**k-fold** validation)¹:



¹Image source: <https://scikit-learn.org/>

Cross Validation

Special cases

- ③ Control by individual objects (**leave-one-out**, or LOO validation) — a special case of k -fold validation, if k is equal to the cardinality of the training set

Cross Validation

Special cases

- ③ Control by individual objects (**leave-one-out**, or LOO validation) — a special case of k -fold validation, if k is equal to the cardinality of the training set
- ④ **Multiple k -fold** validation — repeat k -fold validation several times with different splits.

Cross-validation: the main stages of correct training

- 1 We come up with a model and hyperparameter space

²Image source: <https://scikit-learn.org/>

Cross-validation: the main stages of correct training

- 1 We come up with a model and hyperparameter space
- 2 We select a test set from the initial data, and divide the remaining set into training and validation

²Image source: <https://scikit-learn.org/>

Cross-validation: the main stages of correct training

- 1 We come up with a model and hyperparameter space
- 2 We select a test set from the initial data, and divide the remaining set into training and validation
- 3 Train the model

²Image source: <https://scikit-learn.org/>

Cross-validation: the main stages of correct training

- 1 We come up with a model and hyperparameter space
- 2 We select a test set from the initial data, and divide the remaining set into training and validation
- 3 Train the model
- 4 We carry out the cross-validation procedure across the space of hyperparameters and find their optimal values

²Image source: <https://scikit-learn.org/>

Cross-validation: the main stages of correct training

- 1 We come up with a model and hyperparameter space
- 2 We select a test set from the initial data, and divide the remaining set into training and validation
- 3 Train the model
- 4 We carry out the cross-validation procedure across the space of hyperparameters and find their optimal values
- 5 We train on the full training set with selected hyperparameters

²Image source: <https://scikit-learn.org/>

Cross-validation: the main stages of correct training

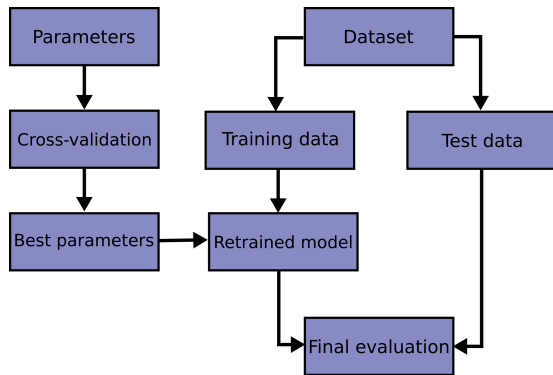
- 1 We come up with a model and hyperparameter space
- 2 We select a test set from the initial data, and divide the remaining set into training and validation
- 3 Train the model
- 4 We carry out the cross-validation procedure across the space of hyperparameters and find their optimal values
- 5 We train on the full training set with selected hyperparameters
- 6 Let's check it on the test!

²Image source: <https://scikit-learn.org/>

Cross-validation: the main stages of correct training

- 1 We come up with a model and hyperparameter space
- 2 We select a test set from the initial data, and divide the remaining set into training and validation
- 3 Train the model
- 4 We carry out the cross-validation procedure across the space of hyperparameters and find their optimal values
- 5 We train on the full training set with selected hyperparameters
- 6 Let's check it on the test!

General scheme²:



²Image source: <https://scikit-learn.org/>

Tuning Hyperparameters

- So far, we can split the examples set into training and testing subsets

Tuning Hyperparameters

- So far, we can split the examples set into training and testing subsets
- We also know how to subdivide the training subset for the validation procedure and even measure the quality during this procedure

Tuning Hyperparameters

- So far, we can split the examples set into training and testing subsets
- We also know how to subdivide the training subset for the validation procedure and even measure the quality during this procedure
- But the main goal is to select the optimal set of hyperparameters!

Tuning Hyperparameters

- So far, we can split the examples set into training and testing subsets
- We also know how to subdivide the training subset for the validation procedure and even measure the quality during this procedure
- But the main goal is to select the optimal set of hyperparameters!
- Usually two approaches are used:
 - ▶ Grid Search: to traverse a predefined range of hyperparameters

Tuning Hyperparameters

- So far, we can split the examples set into training and testing subsets
- We also know how to subdivide the training subset for the validation procedure and even measure the quality during this procedure
- But the main goal is to select the optimal set of hyperparameters!
- Usually two approaches are used:
 - ▶ Grid Search: to traverse a predefined range of hyperparameters
 - ▶ Randomized Search: to generate hyperparameters randomly (according to their given distributions)

Tuning Hyperparameters

- So far, we can split the examples set into training and testing subsets
- We also know how to subdivide the training subset for the validation procedure and even measure the quality during this procedure
- But the main goal is to select the optimal set of hyperparameters!
- Usually two approaches are used:
 - ▶ Grid Search: to traverse a predefined range of hyperparameters
 - ▶ Randomized Search: to generate hyperparameters randomly (according to their given distributions)
 - ▶ Usually there is not much difference — and if you do not need to check **specific** values of hyperparameters in advance, then it is better to limit yourself to a random search

Mandatory external links to read

- ➊ Read the sections 1.2.1, 2.6, (1.3 optionally), 4.1, 5.1, 5.3, 5.7 from "**The Hundred-Page Machine Learning Book**" (see "*References*" course page)
 - ▶ Please don't be afraid: all these sections are 1-2 pages max! After you read it, the items below are to increase your understanding on it
- ➋ Read the reminder about Supervised Learning: [IBM article](#)
 - ▶ This set of readings will remind the main concepts of Supervised Learning
- ➌ Read the section about [How Supervised Learning algorithms work](#)
 - ▶ This will provide the insight on the machinery behind simple ML algorithms
- ➍ Read the page about what is a [feature](#) for an ML algorithm
 - ▶ This will provide the high-level understanding of how feed the inputs into ML algorithms
- ➎ Read the introduction into [loss](#) functions
 - ▶ This will share with you the overview of one of the main concept in ML

Takeaway notes

- ① While doing Supervised Learning, need to define the input and output of the model, and what is the optimization criteria for training

Takeaway notes

- 1 While doing Supervised Learning, need to define the input and output of the model, and what is the optimization criteria for training
- 2 It is necessary to divide the available data into training, validation and test sets — from the very beginning

Takeaway notes

- 1 While doing Supervised Learning, need to define the input and output of the model, and what is the optimization criteria for training
- 2 It is necessary to divide the available data into training, validation and test sets — from the very beginning
- 3 Cross-validation can be of the very different types, but the main goal is the same: to test the generalization ability of the ML model (generalization means performance on an independent data set)

Takeaway notes

- ➊ While doing Supervised Learning, need to define the input and output of the model, and what is the optimization criteria for training
- ➋ It is necessary to divide the available data into training, validation and test sets — from the very beginning
- ➌ Cross-validation can be of the very different types, but the main goal is the same: to test the generalization ability of the ML model (generalization means performance on an independent data set)
- ➍ Hyperparameters tuning is needed for almost every ML model

Thank you!