

Lecture Notes

CS-E4740 Federated Learning

offered in spring 2023 at Aalto University,

<https://fitech.io/en/>

ROUGH DRAFT - DO NOT DISTRIBUTE FURTHER!

Comments are warmly welcome at alex.jung@aalto.fi and will
be acknowledged properly.

Dipl.-Ing. Dr.techn. Alexander Jung[‡]

March 11, 2023

Abstract

The course **CS-E4740 Federated Learning** discusses theory and algorithms for federated learning (FL) from decentralized data with an intrinsic network structure (= networked data). Such networked data consists of local datasets that are related by different notions of similarities. FL methods collaboratively train local (“personalized”) models for each local dataset. These tailored local models are coupled

*

†

via a network structure that reflects statistical similarities between local datasets. We use the network structure of data to guide this coupling via regularization techniques. In particular, we use different measures for the variation of local models as a regularization term. The resulting generalized total variation minimization provides a unifying design principle for practical FL algorithms. These algorithms solve generalized total variation via distributed optimization methods that are able to cope with limited computational resources and imperfections.

Contents

1	Introduction	18
1.1	Related Courses	20
1.2	Outline of the Course	21
1.3	Acknowledgements	22
1.4	Exercises	23
I	Plain Old Machine Learning	24
2	Three Components of ML	25
2.1	Design Choice: Data [1, Ch. 2.1]	25
2.2	Design Choice: Model [1, Ch. 2.2]	28
2.3	Design Choice: Loss [1, Ch. 2.3]	33
2.3.1	Loss Functions for Numeric Labels	37
2.3.2	Loss Functions for Categorical Labels	39
2.4	Exercises	46
3	A Design Principle for ML	48
3.1	The i.i.d. Assumption	48
3.2	Empirical Risk Minimization	50
3.3	How Much Training is Needed?	53
3.4	Exercises	58
4	Regularization	61
4.1	Overfitting	61
4.2	Regularization via Data, Model and Loss	65

4.3	Federated Learning via Regularization	68
4.4	Exercises	71
5	Gradient Methods	73
5.1	Gradient Descent	74
5.2	Projected Gradient Descent	83
5.3	Perturbed Gradient Descent	84
5.4	Stochastic Gradient Descent	86
5.5	Exercises	90
II	Federated Learning	92
6	Networked Data and Models	93
6.1	The Empirical Graph	94
6.2	Networked Models	96
6.3	Clustering Assumption	98
6.4	Graph Learning Methods	100
6.4.1	Testing Similarity	101
6.4.2	Total Variation Minimization	103
6.5	Exercises	104
7	A Design Principle for FL	105
7.1	Generalized Total Variation	105
7.2	Generalized Total Variation Minimization	106
7.3	Interpretations	107
7.3.1	Regularized Empirical Risk Minimization	108

7.3.2	Multi-task Learning	108
7.3.3	Clustering	108
7.3.4	Locally Weighted Learning	109
7.4	Non-Parametric Models	109
7.5	Exercises	112
8	Main Flavours of Federated Learning	113
8.1	Centralized Federated Learning	114
8.2	Decentralized Federated Learning	114
8.3	Clustered Federated Learning	115
8.4	Personalized Federated Learning	115
8.5	Vertical Federated Learning	116
8.6	Horizontal Federated Learning	117
9	Federated Learning Algorithms	119
9.1	FedRelax	119
9.2	FedSGD	122
9.3	FedAvg	123
9.4	Asynchronous Computation	123
9.5	Exercises	127
III	Trustworthy Federated Learning	128
10	Requirements for Trustworthy AI	129
11	Privacy Protection	131
11.1	Adding Noise	133

11.2 Feature Perturbation	135
11.3 Exercises	140
12 Data Poisoning	141
12.1 Denial-of-Service Attacks	141
12.2 Backdoor Attack	141
12.3 Exercises	143
IV Appendix	144
13 Linear Algebra and Convex Analysis	145
13.1 Convex Sets	145
13.2 Convex Functions	145
13.3 Proximal Operators	146
14 Information Theory	146
15 Fixed Point Theory	146
Glossary	147

Lists of Symbols

Sets and Functions

$a \in \mathcal{A}$	This statement indicates that the object a is an element of the set \mathcal{A} .
$a := b$	This statement defines a to be shorthand for b .
$ \mathcal{A} $	The cardinality (number of elements) of a finite set \mathcal{A} .
$\mathcal{A} \subseteq \mathcal{B}$	\mathcal{A} is a subset of \mathcal{B} .
$\mathcal{A} \subset \mathcal{B}$	\mathcal{A} is a strict subset of \mathcal{B} .
\mathbb{N}	The set of natural numbers $1, 2, \dots$
\mathbb{R}	The set of real numbers x [2].
\mathbb{R}_+	The set of non-negative real numbers $x \geq 0$.
\mathbb{R}_{++}	The set of positive real numbers $x > 0$.
$h(\cdot): \mathcal{A} \rightarrow \mathcal{B} : a \mapsto h(a)$	A function (map) that accepts any element $a \in \mathcal{A}$ from a set \mathcal{A} as input and delivers a well-defined element $h(a) \in \mathcal{B}$ of a set \mathcal{B} . The set \mathcal{A} is the domain of the function h and the set \mathcal{B} is the codomain of h . ML aims at finding (or learning) a function h (“hypothesis”) that reads in the features \mathbf{x} of a data point and delivers a prediction $h(\mathbf{x})$ for its label y .

$\{0, 1\}$	The binary set that consists of the two real numbers 0 and 1.
$[0, 1]$	The closed interval of real numbers x with $0 \leq x \leq 1$.
$\operatorname{argmin} f(\mathbf{w})$	The set of minimizers for a real-valued function $f(\mathbf{w})$.
$\log a$	The logarithm of the positive number $a \in \mathbb{R}_{++}$.

Matrices and Vectors

$\mathbf{I}_{l \times d}$	A generalized identity matrix with l rows and d columns. The entries of $\mathbf{I}_{l \times d} \in \mathbb{R}^{l \times d}$ are equal to 1 along the main diagonal and equal to 0 otherwise.
\mathbf{I}	A square identity matrix whose shape should be clear from the context.
\mathbb{R}^d	The set of vectors $\mathbf{x} = (x_1, \dots, x_d)^T$ consisting of d real-valued entries $x_1, \dots, x_d \in \mathbb{R}$.
$\mathbf{x} = (x_1, \dots, x_d)^T$	A vector of length d . The j th entry of the vector is denoted x_j .
$\ \mathbf{x}\ _2$	The Euclidean (or “ ℓ_2 ”) norm of the vector $\mathbf{x} = (x_1, \dots, x_d)^T \in \mathbb{R}^d$ given as $\ \mathbf{x}\ _2 := \sqrt{\sum_{j=1}^d x_j^2}$.
$\ \mathbf{x}\ $	Some norm of the vector $\mathbf{x} \in \mathbb{R}^d$ [3]. Unless specified otherwise, we mean the Euclidean norm $\ \mathbf{x}\ _2$.
\mathbf{x}^T	The transpose of a vector \mathbf{x} that is considered a single column matrix. The transpose is a single-row matrix (x_1, \dots, x_d) .
\mathbf{X}^T	The transpose of a matrix $\mathbf{X} \in \mathbb{R}^{m \times d}$. A square real-valued matrix $\mathbf{X} \in \mathbb{R}^{m \times m}$ is called symmetric if $\mathbf{X} = \mathbf{X}^T$.
$\mathbf{0} = (0, \dots, 0)^T$	A vector of zero entries.

$(\mathbf{v}^T, \mathbf{w}^T)^T$	The vector of length $d + d'$ obtained by concatenating the entries of vector $\mathbf{v} \in \mathbb{R}^d$ with the entries of $\mathbf{w} \in \mathbb{R}^{d'}$.
$\text{span}\{\mathbf{B}\}$	The span of a matrix $\mathbf{B} \in \mathbb{R}^{a \times b}$, which is the subspace of all linear combinations of columns of \mathbf{B} , $\text{span}\{\mathbf{B}\} = \{\mathbf{B}\mathbf{a} : \mathbf{a} \in \mathbb{R}^b\} \subseteq \mathbb{R}^a$.
\mathbb{S}_+^d	The set of all positive semi-definite (psd) matrices of size $d \times d$.
$\det(\mathbf{C})$	The determinant of the matrix \mathbf{C} .

Probability Theory

$\mathbb{E}_p\{f(\mathbf{z})\}$ The expectation of a function $f(\mathbf{z})$ of a RV \mathbf{z} whose probability distribution is $p(\mathbf{z})$. If the probability distribution is clear from context we just write $\mathbb{E}\{f(\mathbf{z})\}$.

$p(\mathbf{x}, y)$ A (joint) probability distribution of a RV whose realizations are data points with features \mathbf{x} and label y .

$p(\mathbf{x}|y)$ A conditional probability distribution of a RV \mathbf{x} given the value of another RV y [4, Sec. 3.5].

$p(\mathbf{x}; \mathbf{w})$ A parametrized probability distribution of a RV \mathbf{x} . The probability distribution depends on a parameter vector \mathbf{w} . For example, $p(\mathbf{x}; \mathbf{w})$ could be a multivariate normal distribution of a Gaussian RV \mathbf{x} with the parameter vector \mathbf{w} given by the entries of the mean vector $\mathbb{E}\{\mathbf{x}\}$ and the covariance matrix $\mathbb{E}\left\{(\mathbf{x} - \mathbb{E}\{\mathbf{x}\})(\mathbf{x} - \mathbb{E}\{\mathbf{x}\})^T\right\}$.

$\mathcal{N}(\mu, \sigma^2)$ The probability distribution of a scalar normal (“Gaussian”) RV $x \in \mathbb{R}$ with mean (or expectation) $\mu = \mathbb{E}\{x\}$ and variance $\sigma^2 = \mathbb{E}\{(x - \mu)^2\}$.

$\mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$ The multivariate normal distribution of a vector-valued Gaussian RV $\mathbf{x} \in \mathbb{R}^d$ with mean (or expectation) $\boldsymbol{\mu} = \mathbb{E}\{\mathbf{x}\}$ and covariance matrix $\mathbf{C} = \mathbb{E}\{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T\}$.

Machine Learning

r	An index $r = 1, 2, \dots$, that enumerates data points.
m	The number of data points in (the size of) a dataset.
\mathcal{D}	A dataset $\mathcal{D} = \{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ is a list of individual data points $\mathbf{z}^{(r)}$, for $r = 1, \dots, m$.
d	Number of features that characterize a data point.
x_j	The j th feature of a data point. The first feature of a given data point is denoted x_1 , the second feature x_2 and so on.
\mathbf{x}	The feature vector $\mathbf{x} = (x_1, \dots, x_d)^T$ of a data point whose entries are the individual features of a data point. Beside the symbol \mathbf{x} , we sometimes use \mathbf{z} as another symbol to denote a vector whose entries are features of a data point.
\mathbf{z}	We need two different symbols to distinguish between “raw” or “original” and learnt features [1, Ch. 9].
$\mathbf{x}^{(r)}$	The feature vector of the r th data point within a dataset.
$x_j^{(r)}$	The j th feature of the r th data point within a dataset.
\mathcal{B}	A mini-batch (subset) of randomly chosen data points.
B	The size of (the number of data points in) a mini-batch.

y	The label (quantity of interest) of a data point.
$y^{(r)}$	The label of the r th data point.
$(\mathbf{x}^{(r)}, y^{(r)})$	The features and label of the r th data point.
\mathcal{Y}	<p>The label space \mathcal{Y} of a ML method consists of all potential label values that a data point can have. We often use label spaces that are larger than the set of different label values arising in a give dataset (e.g., a training set). We refer to a ML problems (methods) using a numeric label space, such as $\mathcal{Y} = \mathbb{R}$ or $\mathcal{Y} = \mathbb{R}^3$, as regression problems (methods). ML problems (methods) that use a discrete label space, such as $\mathcal{Y} = \{0, 1\}$ or $\mathcal{Y} = \{\text{“cat”}, \text{“dog”}, \text{“mouse”}\}$ are referred to as classification problems (methods).</p>
α	learning rate (step-size) used by gradient-based methods.
$h(\cdot)$	A hypothesis map that reads in features \mathbf{x} of a data point and delivers a prediction $\hat{y} = h(\mathbf{x})$ for its label y .
\mathcal{H}	A hypothesis space or model used by a ML method. The hypothesis space consists of different hypothesis maps $h : \mathcal{X} \rightarrow \mathcal{Y}$ between which the ML method has to choose.

B^2	<p>The squared bias of a learnt hypothesis \hat{h} delivered by a ML algorithm that is fed with data points which are modelled as realizations of RVs. If data is modelled as realizations of RVs, also the delivered hypothesis \hat{h} is the realization of a RV.</p>
V	<p>The variance of the (parameters of the) hypothesis delivered by a ML algorithm. If the input data for this algorithm is interpreted as realizations of RVs, so is the delivered hypothesis a realization of a RV.</p>
$L((\mathbf{x}, y), h)$	<p>The loss incurred by predicting the label y of a data point using the prediction $\hat{y} = h(\mathbf{x})$. The prediction \hat{y} is obtained from evaluating the hypothesis $h \in \mathcal{H}$ for the feature vector \mathbf{x} of the data point.</p>
E_v	<p>The validation error of a hypothesis h, which is its average loss incurred over a validation set.</p>
$\hat{L}(h \mathcal{D})$	<p>The empirical risk or average loss incurred by the predictions of hypothesis h for the data points in the dataset \mathcal{D}.</p>
E_t	<p>The training error of a hypothesis h, which is its average loss incurred over a training set.</p>
t	<p>A discrete-time index $t = 0, 1, \dots$ used to enumerate a sequence to sequential events (“time instants”).</p>
t	<p>An index that enumerates learning tasks within a multi-task learning problem.</p>

λ	A regularization parameter that controls the amount of regularization.
$\lambda_j(\mathbf{Q})$	The j th eigenvalue (sorted either ascending or descending) of a psd matrix \mathbf{Q} . We also use the shorthand λ_j if the corresponding matrix is clear from context.
$\sigma(\cdot)$	The activation function used by an artificial neuron within an artificial neural network (ANN).
$\mathcal{R}_{\hat{y}}$	A decision region within a feature space.
\mathbf{w}	A parameter vector $\mathbf{w} = (w_1, \dots, w_d)^T$ whose entries are parameters of a model. These parameters could be feature weights in linear maps, the weights in ANNs or the thresholds used for splits in decision trees.
$h^{(\mathbf{w})}(\cdot)$	A hypothesis map that involves tunable parameters w_1, \dots, w_d which are stacked into the vector $\mathbf{w} = (w_1, \dots, w_d)^T$.
$\nabla f(\mathbf{w})$	The gradient of a differentiable real-valued function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is the vector $\nabla f(\mathbf{w}) = (\frac{\partial f}{\partial w_1}, \dots, \frac{\partial f}{\partial w_d})^T \in \mathbb{R}^d$ [5, Ch. 9].
$\phi(\cdot)$	A feature map $\phi : \mathcal{X} \rightarrow \mathcal{X}' : \mathbf{x} \mapsto \mathbf{x}' := \phi(\mathbf{x}) \in \mathcal{X}'$.

$\mathcal{G} = (\mathcal{V}, \mathcal{E})$	Empirical graph whose nodes $i \in \mathcal{V}$ carry local datasets and local models.
$i \in \mathcal{V}$	A node in the empirical graph that represents a local dataset and a corresponding local model. It might also be useful to think of node i as a small computer that can collect data and execute computations for model training.
$\mathcal{D}^{(i)}$	The local dataset $\mathcal{D}^{(i)}$ at node $i \in \mathcal{V}$.
m_i	The number of data points (sample size) contained in the local dataset $\mathcal{D}^{(i)}$ at node $i \in \mathcal{V}$.
\mathbf{x}	The feature vector $\mathbf{x} = (x_1, \dots, x_d)^T$ of a generic data point.
\mathcal{X}	The feature space \mathcal{X} is the set of all possible values that the features \mathbf{x} of a data point can take on.
$\mathbf{x}^{(i,r)}$	The feature vector of the r -th data point in the local dataset $\mathcal{D}^{(i)}$.
$\mathbf{v} \in \mathbb{R}^d$	Some numeric vector (one-dimensional array) of length d .

$\ \mathbf{w}\ _2$	The Euclidean norm $\ \mathbf{w}\ _2 := \sqrt{\sum_{j=1}^d w_j^2}$ of the d -dimensional vector \mathbf{w} .
h	A hypothesis map $h : \mathcal{X} \rightarrow \mathcal{Y}$ that reads in the features \mathbf{x} of a data point and delivers a prediction $h(\mathbf{x})$ for its label y .
$\mathcal{Y}^{\mathcal{X}}$	Given two sets \mathcal{X} and \mathcal{Y} , we denote by $\mathcal{Y}^{\mathcal{X}}$ the set of all maps $h : \mathcal{X} \rightarrow \mathcal{Y}$.
\mathcal{H}	A hypothesis space which consists of hypothesis maps h .
$L^{(\text{d})}(\mathbf{x}, h(\mathbf{x}), h'(\mathbf{x}))$	The loss incurred by a “external” hypothesis h' on a data point with features \mathbf{x} and predicted label $h(\mathbf{x})$ that is obtained from some local hypothesis.

1 Introduction

Many important application domains generate decentralized collections of local datasets that are related via an intrinsic network structure [6]. Two timely application domains that generate such networked data are (i) the high-precision management of pandemics and (ii) the Internet of Things (IoT) [7, 8]. Here, local datasets are generated by smartphones, wearables or other IoT devices [9, 10]. These local datasets are related via physical contact networks, social networks [11], co-morbidity networks [12], or communication networks [13].

FL is an umbrella term for distributed optimization techniques to train machine learning (ML) models from decentralized collections of local datasets [14–18]. These methods implement computations, such as gradient computations (see Section 5) for ML model training, at the location of data generation. This design philosophy is different from a naive approach of first collecting all local datasets at a single location where then a single ML model is trained. Distributed training of (several different) ML models at the spatial location of data generation can be beneficial in several aspects [19]:

- **Privacy.** FL methods are appealing for applications involving sensitive data (such as healthcare) as they do not require the exchange of raw data but only model (parameter) updates [16, 17]. By sharing only model updates, FL methods are considered privacy-friendly in the sense of not leaking (too much) sensitive information that is contained in the local datasets (see Section 11).
- **Robustness.** By relying on decentralized data and computation, FL

methods offer robustness against hardware failures (such as “stragglers”) and data poisonings (see Section 12).

- **Parallel Computing.** A main application domain for FL are mobile networks, consisting of humans equipped with smartphones. We can interpret a mobile network as a parallel computer which is constituted by smartphones that can communicate via radio links. This parallel hardware allows to speed up computational tasks such as the computation of gradients required to train ML models (see Section 5).
- **Trading Computation against Communication.** Consider a FL application where local datasets are generated by low-complexity devices at remote locations that cannot be easily accessed. The cost of communicating raw local datasets to some central unit (which then trains a single global ML model) might be much higher than the computational cost incurred by using the low-complexity devices to (partially) train ML models [20].
- **Personalization.** FL can be used to train personalized ML models for collections of local datasets, which might be generated by smartphones (and their users) [21]. A key challenge for ensuring personalization is the heterogeneity of local datasets [22, 23]. Indeed, the statistical properties of different local datasets might vary significantly such they cannot be well modelled as independent and identically distributed (i.i.d.). Each local dataset induces a separate learning task that consists of learning useful parameter values for a local model. This course discusses FL methods to train personalized models via combining the information

carried in decentralized and heterogeneous data (see Section 8.3 and Section 8.4).

1.1 Related Courses

In what follows we briefly explain how this course CS-E4740 relates to selected courses at Aalto University.

- **CS-EJ3211 - Machine Learning with Python.** Teaches the application of basic ML methods using the Python package (library) `scikit-learn` [24]. CS-E4740 couples a network of basic ML methods using regularization techniques to obtain tailored (personalized) ML models for local datasets. This coupling is required to adaptive pool local datasets obtain sufficiently large training sets for the personalized ML models.
- **CS-C3240 - Machine Learning (spring 2022 edition).** Teaches basic theory of ML models and methods [1]. CS-E4740 combines the components of basic ML methods, such as data representation and models, with network models. In particular, instead of a single dataset and a single model (such as a decision tree), we will study networks of local datasets and local models.
- **ABL-E2606 - Data Protection.** This course discusses important legal constraints (“laws”), including the European general data protection regulation (GDPR), for the use of data and, in turn, for the design of trustworthy FL methods (see Part III).

- **MS-C2105 - Introduction to Optimization.** This course teaches basic optimisation theory and how to model applications as (linear, integer, and non-linear) optimization problems. CS-E4740 uses optimization theory and methods to formulate FL problems (see Section 7) and design FL methods (see Section 9).
- **ELEC-E5424 - Convex Optimization.** This course teaches advanced optimisation theory for the important class of convex optimization problems [25]. Convex optimization theory and methods can be used for the study and design of FL algorithms.

1.2 Outline of the Course

This notes are divided into three parts:

- Part I discusses the basic components and principles of (non-federated) ML: Section 2 introduces data, models and loss functions as three main components of ML. Section 3 explains how these components are combined via empirical risk minimization (ERM). The regularization of ERM, via each of its three main components, is then discussed in Section 4. We then explain when and how to solve regularized ERM via simple gradient descent (GD) methods in Section 5.
- Part II shows how to couple (the training of) tailored (personalized) ML models for decentralized data with an intrinsic network structure: Section 6 introduces an empirical graph as a representation for collections of local datasets and corresponding tailored models. The undirected and weighted edges of the empirical graph represent statistical similarities

between local datasets. Section 7 uses the similarity structure to formulate GTV minimization as a main design principle for FL. We then show how main flavours of FL are obtained as special cases of GTVMin in Section 8. Section 9 applies well-known optimization methods to GTVMin, resulting in FL algorithms for distributed training of tailored (“personalized”) models.

- Part III discusses requirements and design aspects for trustworthy FL: Section 10 enumerates seven key requirements for trustworthy artificial intelligence (AI) put forward by the European Union. We then discuss privacy and its protection by FL algorithms in Section 11. Section 12.2 discusses the vulnerability of FL methods against backdoor attacks.

1.3 Acknowledgements

The author is grateful for feedback on the lecture notes received from: Dr. Shamsiat Abdurakhmanova, Dr. Yu Tian, Alexander Pavlyuk, Junyuan Fang, Matti Ahokas, Felix Knop, Iina Hyytiäinen, Maria Montoya Freire and Ariana Marta. Some of the figures have been prepared with the help of Dr. Yu Tian.

1.4 Exercises

Exercise 1.1. Parallel Computing (5 points). Assume that you have a single computer that is able to sum two numbers per second. You can assume that the computer has sufficient storage memory which can be read from and written to at very high speed. How long does the computer take to compute the sum of 100 numbers $\mathcal{D} = \{x^{(1)}, \dots, x^{(100)}\}$.

Now assume that you have two identical such computers which are connected via a high-speed network. How much faster can we compute the sum of the numbers in \mathcal{D} using these two computers? Hint: see [26, Sec. 1.2.3.].

Exercise 1.2. Communication Bottleneck (5 points). Consider two laptops, each having a webcam taking pictures with 1080×1080 pixels. These two laptops are connected via a wireless **Bluetooth** connection with average bitrate of $1 \cdot 10^6$ bits/s [27]. What is the maximum rate at which one laptop can record images and transmit it to the other laptop? You can ignore any compression techniques and assume that images are stored as bitmaps with each pixel encoded by three bytes (red, green and blue intensities on scale $0, \dots, 255$).

Part I

Plain Old Machine Learning

2 Three Components of ML

ML revolves around computing accurate predictions for a quantity of interest. These predictions are computed by evaluating a hypothesis map $h(\mathbf{x})$ which is fed with the features \mathbf{x} of a data point (see Section 2.1). The value $h(\mathbf{x})$ is a prediction for some quantity of interest which we refer to as the label of a data point.

In general, the hypothesis h is learnt or optimized based on the discrepancy between previous predictions and observations. The space of possible hypothesis maps, from which a ML method can choose from, is referred to as hypothesis space or model (see Section 2.2).

To choose or learn a useful hypothesis out of a hypothesis space we need a measure for the quality of the predictions obtained from a hypothesis. To this end, ML methods use loss functions to obtain a quantitative measure for the prediction errors (see Section 2.3).

Different ML methods use different choices for data points (their features and label), the hypothesis space (or model) and loss function. The design choices for these components are guided by computational aspects and statistical aspects. In what follows we discuss each of these design choices for a toy application which is to predict the maximum daytime temperature.

2.1 Design Choice: Data [1, Ch. 2.1]

We consider data as collections of elementary data points that carry relevant information. Each individual data point is characterized by several properties that we group into features and labels.

The features \mathbf{x} are low-level properties that we can measure or compute easily in an automated fashion. Examples for such low-level properties are physical quantities that can be measured with low-cost hardware, such as temperature or humidity. The feature space \mathcal{X} collects all possible values that the features $\mathbf{x} \in \mathcal{X}$ of a data point can take on. This course focuses on ML methods that use an Euclidean space \mathbb{R}^d as the feature space.

Besides their features, data points are characterized by labels y which are high-level properties or quantities of interest. The label space \mathcal{Y} collects the possible values that the labels of a data point can take on, i.e., $y \in \mathcal{Y}$. ML applications involving a “numeric” \mathcal{Y} (e.g., the real numbers \mathbb{R}) are referred to as regression problems. ML applications involving a finite label space \mathcal{Y} are referred to as classification problems.

In contrast to features, determining the label of a data point is more difficult and typically requires to consult a human domain expert. Another source for the label of a data point might be the predictions obtained from a reference (or “teacher”) model that has been trained on a different but related dataset. The FL methods discussed in part II use this approach to leverage similarities between learning tasks arising from datasets with similar statistical properties.

Note that the distinction between features and labels is blurry. Indeed, one and the same property of a data point might be considered as features in one application, while it might be useful to consider it as a label in another application.

Our definition of data points as elementary information-carrying objects is quite abstract and therefore flexible. Indeed, the definition of what we

consider as data points is an important design choice. As an example, we could use data points to represent time periods such as days. One data point would then represent the day 12-Mar-2020. Another data point represents the day 12-Apr-2019. Instead of entire days, we could also define data points as shorter time-periods such as hours or 10-minute intervals.

After we have defined data points we are faced with two further design choices:

- **Which features to use? (feature selection).** Which properties should we use as features of a data point? In the above example, involving a data point 12-Apr-2019, we could use the maximum daytime temperatures during the previous two days as features. Another choice for the features would be to use these temperatures during the previous ten days as features.
- **Which label to use? (learning task).** What property of a data point is a quantity of interest or label y ? For the above weather forecasting example, we could define the label as the average temperature during the next two days. Another choice for the label y would be the average temperature during the next ten days. We could also define several different labels y_1, y_2, \dots , for the same data point. For example, y_j could be the average temperature during the next j days.

The design choices for the data points, their features and label must take into account statistical aspects, computational aspects and trustworthiness of the resulting ML method. For example, using a large number of features will typically result in a higher computational complexity of the resulting ML

method. On the other hand, a large number of features might be necessary to enable learning an accurate hypothesis. However, to ensure explainability of the resulting ML method, we might prefer using few interpretable features [28].

The design choice for the labels of a data point must also take into account the trustworthiness of the resulting ML methods. If a data point represents a user of a social network or online marketplace, then we might avoid labels that allow to predict individual behaviour and, in turn, enable manipulation of the human user.

2.2 Design Choice: Model [1, Ch. 2.2]

Consider a ML application that involves data points, each characterized by features $\mathbf{x} \in \mathcal{X}$ and a label $y \in \mathcal{Y}$. The overall goal of ML is to find or learn a hypothesis map,

$$h : \mathcal{X} \rightarrow \mathcal{Y} : \mathbf{x} \mapsto h(\mathbf{x}). \quad (1)$$

that computes a prediction $\hat{y} = h(\mathbf{x})$ for the label $y \in \mathcal{Y}$ of a data point solely from its features $\mathbf{x} \in \mathcal{X}$.

Practical ML methods have only finite computational resources to evaluate a hypothesis map (1) and evaluate them (see Section 2.3). Any practical ML method must therefore restrict itself to a (tiny) subset of the set $\mathcal{Y}^{\mathcal{X}}$ that contains all possible hypothesis maps. This subset of computationally feasible (“affordable”) hypothesis maps is referred to as the hypothesis space (or model),

$$\mathcal{H} := \{h : h(\mathbf{x}) \text{ can be computed efficiently.}\} \quad (2)$$

As depicted in Figure 1, ML methods typically use a hypothesis space \mathcal{H} that

is a tiny subset of $\mathcal{Y}^{\mathcal{X}}$.

Similar to the features and labels used to characterize data points, also the hypothesis space underlying a ML method is a design choice. The design choice for the hypothesis space involves a trade-off between computational complexity and accuracy of the resulting ML methods [1, Ch. 2]. Note that the informal criterion “can be computed efficiently” depends on the available computational resources. Graphic processing units can efficiently compute matrix-vector operations that constitute a hypothesis map [29]. However, we can also implement ML methods using optical computers that realize ANNs using diffractive layers [30].

The hypothesis space should be large enough to contain at least one hypothesis that allows to accurately predict the label from the features of a data point. However, it should not be too large such that ML methods can efficiently search it for a useful hypothesis. Another limitation for the size or dimension of the hypothesis space is the amount of available data (see Section 3.3). Beside the computational aspects and statistical aspects, a third design aspect for the model is the explainability of the resulting ML method [31, 32].

During the rest of this course, we will focus on FL methods that use linear models [1, Ch. 3]. In its most basic form, linear models can only be used for data points characterized by numeric features $x_1, \dots, x_d \in \mathbb{R}$. It is convenient to collect those features of a data point into a feature vector

$$\mathbf{x} := (x_1, \dots, x_d)^T \in \mathbb{R}^d.$$

Formally, for a given number d of numeric features, a linear model is the

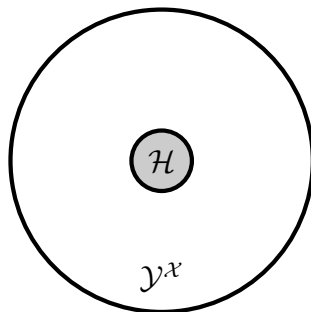


Figure 1: The hypothesis space \mathcal{H} is a (typically very small) subset of the (typically very large) set $\mathcal{Y}^{\mathcal{X}}$ of all possible maps from the feature space \mathcal{X} to the label space \mathcal{Y} .

hypothesis space which consists of all linear maps,

$$\mathcal{H}^{(d)} := \{h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} : \mathbf{w} \in \mathbb{R}^d\}. \quad (3)$$

Note that (103) defines an entire family of hypothesis spaces, which is indexed by the number d of features that are linearly combined to form the prediction $h(\mathbf{x})$. The design choice of d is guided by computational aspects (smaller d means less computation), statistical aspects (larger d might reduce prediction error) and interpretability (favouring a linear model with few features).

Each individual hypothesis in the linear model $\mathcal{H}^{(d)}$ is fully specified by a weight vector $\mathbf{w} \in \mathbb{R}^d$, which we indicate by the notation $h^{(\mathbf{w})}(\mathbf{x}) := \mathbf{w}^T \mathbf{x}$. Thus, the linear model (103) is an example for a parametrized model. Another example for a parametrized model is an ANN [1, Ch. 2, 3].

Upgrading Linear Models via Feature Maps. We can use the linear model $\mathcal{H}^{(d)}$ as a building block for constructing non-linear models. For example, we could use a linear model for the transformed features $\mathbf{x}' = \Phi(\mathbf{x})$



Figure 2: We can extend a linear model (103) by replacing the original feature x of a data point with transformed features $\mathbf{x}' = \Phi(\mathbf{x})$. These transformed features are obtained by applying a feature map Φ to the original features.

(see Figure 2). The resulting model consists of all maps that are concatenations of linear maps with the feature transformation or map Φ . Kernel methods use feature maps obtained from a kernel function. The input layers of an ANN (whose output layer is linear) are another instance of a feature map (Figure 3).

Non-Numeric Features. We can use the linear model (103) only for data points whose features are numeric vectors $\mathbf{x} = (x_1, \dots, x_d)^T \in \mathbb{R}^d$. In some application domains, such as natural language processing [33], there is no obvious natural choice for numeric features. However, since ML methods based on the hypothesis space (103) can be implemented efficiently using widely available software and hardware, it might be useful to construct numerical features even for non-numeric data. For text data, there has been significant progress recently on methods that map a human-generated text into sequences

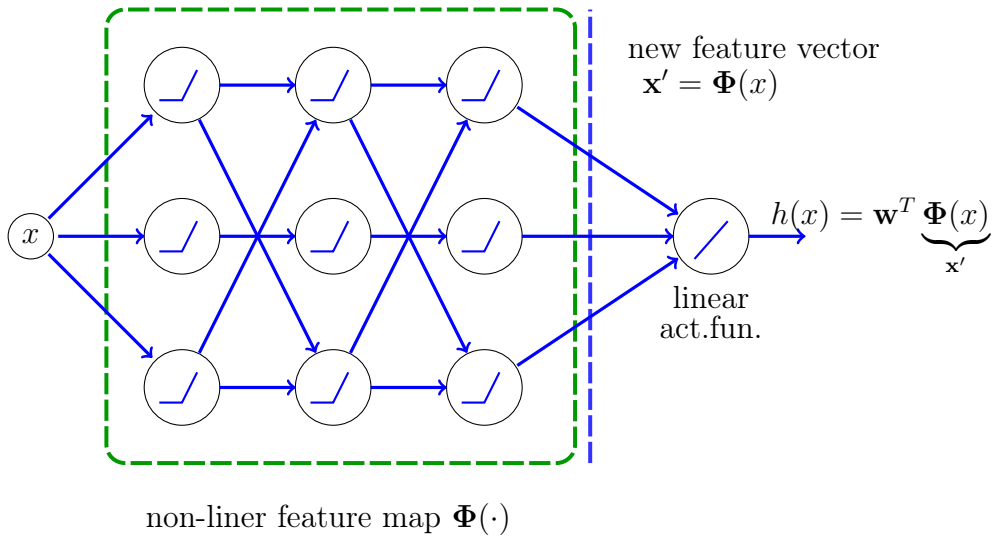


Figure 3: The hypothesis space spanned by an ANN (whose output layer is linear) coincides with the linear model applied to the activations $\Phi(x)$ of the penultimate layer. We can interpret these activations as transformed features that are obtained by applying the input layers to the original features. These input layers implement a non-linear feature map $\Phi(\cdot)$.

of vectors (see [33, Chap. 12] for more details). Moreover, [1, Sec. 9.3] shows how to construct numeric features for data points that are related by some notion of similarity.

2.3 Design Choice: Loss [1, Ch. 2.3]

ML methods find or learn a useful hypothesis out of an entire hypothesis space (or model) \mathcal{H} . To measure the quality of a hypothesis, ML methods use a loss function. Mathematically, a loss function is a map

$$L : \mathcal{X} \times \mathcal{Y} \times \mathcal{H} \rightarrow \mathbb{R}_+ : ((\mathbf{x}, y), h) \mapsto L((\mathbf{x}, y), h)$$

which assigns a pair consisting of a data point (\mathbf{x}, y) , itself characterized by features \mathbf{x} and label y , and a hypothesis $h \in \mathcal{H}$ the non-negative real number $L((\mathbf{x}, y), h)$.

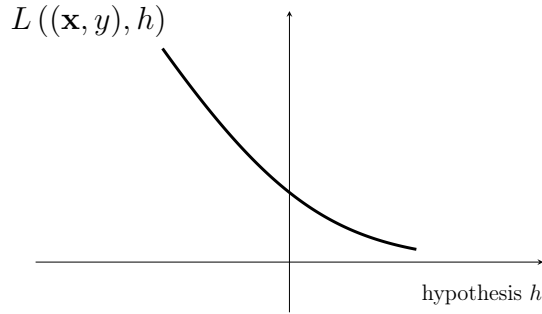


Figure 4: Some loss function $L((\mathbf{x}, y), h)$ for a fixed data point, with features \mathbf{x} and label y , and varying hypothesis h . ML methods try to find (learn) a hypothesis that incurs minimum loss.

The loss value $L((\mathbf{x}, y), h)$ quantifies the discrepancy between the true label y and the prediction $h(\mathbf{x})$. A small (close to zero) value $L((\mathbf{x}, y), h)$

indicates a small discrepancy between the prediction $h(\mathbf{x})$ and the actual (or true) label y of a data point. Figure 4 depicts a loss function for a given data point, with features \mathbf{x} and label y , as a function of the hypothesis $h \in \mathcal{H}$. Note that Figure 4 needs to be understood as a cartoon as we use the horizontal axis to represent a hypothesis space (or model) \mathcal{H} . Note that the elements of \mathcal{H} are hypothesis maps $h : \mathcal{X} \rightarrow \mathcal{Y}$.

We could make Figure 4 more rigorous, by using the horizontal axis to represent a single numeric parameter of a hypothesis map such as its value $h(\mathbf{x})$ for the features of the considered data point (see, e.g., Figure 9). Another option would be to use the horizontal axis to represent the single model parameters of the linear model $\mathcal{H}^{(1)}$.

Much like the choice for the hypothesis space \mathcal{H} used in a ML method, also the loss is a design choice. We will discuss some widely used examples for loss function in Section 2.3.1 and Section 2.3.2. Like data and model, the loss should be chosen in view of three design aspects

- computational aspects,
- statistical aspects,
- trustworthiness of the resulting ML method.

The choice for the loss function impacts the computational complexity of searching the hypothesis space for a hypothesis with minimum loss. Consider a ML method that uses a parametrized hypothesis space and a loss function that is a convex and differentiable (smooth) function of the parameters of a hypothesis. In this case, searching for a hypothesis with small loss can be done efficiently using the gradient-based methods discussed in Section 5.

Some ML methods use a loss function that is neither convex nor differentiable. The minimization of such loss functions tends to be computationally more challenging (requiring more computation). More discussion about the computational complexities of different types of loss functions can be found in [1, Sec. 4.2.].

Beside computational aspects, the choice for the loss function should also take into account statistical aspects. For example, some loss functions result in ML methods that are more robust against outliers [1, Sec. 3.3]. We might also use probabilistic models for the data generated in an ML application to guide the design choice for the loss function. In particular, the maximum likelihood principle suggests to use the logarithm of a likelihood (“log-likelihood”) function as a loss function. This likelihood function is determined by the (parametrized) probability distribution of the data points.

A third aspect for the choice of loss function is its explainability. Section 2.3.2 discusses loss functions to measure the quality of binary classifier. It seems natural to measure the quality of a binary classifier by the average number of wrongly classified data points, which is precisely the average 0/1 loss (7) (see Section 2.3.2).

In contrast to its appealing interpretation as error-rate, the computational aspects of the average 0/1 loss (which is equivalent to accuracy) are less pleasant. Minimizing the average 0/1 loss to learn an accurate hypothesis amounts to a non-convex and non-smooth optimization problem which is typically computationally more challenging compared to minimizing a convex function.

Section 2.3.2 introduces the logistic loss as a computationally attractive

alternative choice for the loss function in binary classification problems. Learning a linear hypothesis with minimum (average) logistic loss amounts to minimizing a differentiable convex function. Section 5 discusses practically useful gradient-based methods to minimize such functions.

To summarize, the computational aspects, the statistical aspects and the explainability of ML methods typically result in conflicting design goals for a loss function. A loss function that has favourable statistical properties might incur a high computational complexity of the resulting ML method. Loss functions that result in computationally efficient ML methods might not allow for an easy interpretation. For example, it is hard to grasp the meaning of a binary classifier achieving an average logistic loss of 10^{-1} . On the other hand, the value of the average 0/1 loss is often considered as an (approximate) error rate. It can be helpful to use different loss functions for the search of a good hypothesis (model training) and for its final evaluation (see Figure 5).

For example, in a binary classification problem, we might use the logistic loss to search for (learn) an accurate hypothesis using the optimization methods in Section 5. After having found (learnt) a hypothesis, we use the average 0/1 loss for the final performance evaluation. The 0/1 loss is appealing for this purpose as it can be interpreted as an error (or misclassification) rate. The loss function used for the final performance evaluation of a learnt hypothesis is sometimes referred to as metric.



Figure 5: Two different loss functions for a given data point and varying hypothesis h . One of these loss functions (solid curve) is used to learn a good hypothesis by minimizing the loss. We use another loss function (dashed curve) to evaluate the performance of the learnt hypothesis \hat{h} .

2.3.1 Loss Functions for Numeric Labels

For ML problems involving data points with numeric labels $y \in \mathbb{R}$, a widely used (first) choice for the loss function is the squared error loss

$$L((\mathbf{x}, y), h) := (y - \underbrace{h(\mathbf{x})}_{=\hat{y}})^2. \quad (4)$$

The squared error loss (4) depends on the features \mathbf{x} only via the predicted label value $\hat{y} = h(\mathbf{x})$. We can evaluate the squared error loss solely using the prediction $h(\mathbf{x})$ and the true label value y . Besides the prediction $h(\mathbf{x})$, no other properties of the features \mathbf{x} influence the value of the squared error loss.

We will slightly abuse notation and use the shorthand $L(y, \hat{y})$ for any loss function that depends on the features \mathbf{x} only via the predicted label $\hat{y} = h(\mathbf{x})$. Figure 6 depicts the squared error loss as a function of the prediction error $y - \hat{y}$. It also depicts the (special case of the) Huber loss

$$L((\mathbf{x}, y), h) := |h(\mathbf{x}) - y|. \quad (5)$$



Figure 6: Two widely used loss functions to measure how well a hypothesis predicts a numeric label are the squared error loss (4) and the Huber loss (5). Note that, for a given hypothesis h , we can evaluate the loss function only if we know the features \mathbf{x} and the label y of the data point.

The squared error loss (4) has appealing computational and statistical properties. For linear predictor maps $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$, the squared error loss is a convex and differentiable function of the parameter vector \mathbf{w} . This allows, in turn, to efficiently search for the optimal linear predictor using efficient iterative optimization methods (see Section 5). The squared error loss also has a useful interpretation in terms of a probabilistic model for the features and labels: Minimizing the squared error loss is equivalent to maximum likelihood estimation within a linear Gaussian model [34, Sec. 2.6.3].

Another loss function used in regression problems is the absolute error loss (5), which is a special case of the Huber loss [1, Sec. 3.3]. It can be shown that ML methods using absolute error loss are more robust against a few outliers in the training set compared to methods using squared error loss (see [1, Sec. 3.3.]). The improved robustness comes at the expense of increased computational complexity of minimizing the non-differentiable loss (5) compared to the convex and differentiable squared error loss (4).

2.3.2 Loss Functions for Categorical Labels

Classification problems involve data points whose labels take on values from a discrete label space \mathcal{Y} . In what follows, unless stated otherwise, we focus on binary classification problems with label space $\mathcal{Y} = \{-1, 1\}$. Classification methods learn a hypothesis or classifier that maps the features \mathbf{x} of a data point to a predicted label $\hat{y} \in \mathcal{Y}$.

It seems that we cannot use linear hypothesis maps $h^{(\mathbf{w})}(\mathbf{x}) := \mathbf{w}^T \mathbf{x}$ (see (103)) for a classification problem. Indeed, the function value $h(\mathbf{x})$ is a real-number and (most likely) does not belong to the label space $\mathcal{Y} = \{-1, 1\}$. However, a simple post-processing step turns a linear hypothesis $h^{(\mathbf{w})}$ into a binary classifier: We classify a data point as

$$\hat{y} = 1 \text{ if } h^{(\mathbf{w})}(\mathbf{x}) > 0, \text{ and } \hat{y} = -1 \text{ otherwise.} \quad (6)$$

Thus, the sign of $h^{(\mathbf{w})}(\mathbf{x})$ determines the predicted label $\hat{y} \in \{-1, 1\}$. We will abuse notation and also refer to the linear map $h^{(\mathbf{w})}(\mathbf{x})$ as a binary classifier.

Besides its sign, we can also use the absolute value $|h^{(\mathbf{w})}(\mathbf{x})|$ to convey useful information. In particular, we might use the absolute value $|h^{(\mathbf{w})}(\mathbf{x})|$ as a measure for the confidence in the prediction \hat{y} (6).

In principle, we can measure the quality of a hypothesis when used to classify data points using the squared error loss (4). However, the squared error loss is a poor quality measure for a binary classifier. Figure 7 illustrates how the squared error loss of a linear hypothesis can be (highly) misleading for evaluating the performance of a binary classifier.

Figure 7 depicts a dataset \mathcal{D} consisting of $m = 4$ data points with binary labels $y^{(r)} \in \{-1, 1\}$, for $r = 1, \dots, m$. The figure also depicts two different



Figure 7: A training set consisting of four data points with binary labels $\hat{y}^{(r)} \in \{-1, 1\}$. Minimizing the squared error loss (4) would prefer the (poor) classifier $h^{(1)}$ over the (reasonable) classifier $h^{(2)}$.

hypotheses $h^{(1)}(x)$ and $h^{(2)}(x)$ for classifying data points (via their signs).

The classifications \hat{y} obtained from the hypothesis $h^{(2)}(x)$ would perfectly match the labels of the four training data points in Figure 7 since $h^{(2)}(x^{(r)}) \geq 0$ if and only if $y^{(r)} = 1$. In contrast, the classifications $\hat{y}^{(r)}$ obtained by thresholding $h^{(1)}(x)$ are incorrect for data points with $y = -1$.

Looking at \mathcal{D} , we might prefer using $h^{(2)}$ over $h^{(1)}$ to classify data points. However, the squared error loss $(y - h^{(2)}(x))^2$ incurred by the (reasonable) classifier $h^{(2)}$ is much larger than the squared error loss incurred by the (poor) classifier $h^{(1)}$. The squared error loss is typically a bad choice for assessing the quality of a hypothesis map that is used for classifying data points into different categories.

Generally speaking, we want the loss function to punish (deliver large

values for) a hypothesis that results in a wrong classification ($\hat{y} \neq y$) with high confidence ($|h(\mathbf{x})|$ is large). On the other hand, a useful loss function should not punish (deliver small values for) a hypothesis that results a correct classification ($\hat{y} = y$) with high confidence ($|h(\mathbf{x})|$ is large). By its very definition, the squared error loss (4) yields always large values if the confidence $|h(\mathbf{x})|$ is large, no matter if the resulting classification (obtained after thresholding) is correct or wrong. This is clearly an undesirable property for a loss function that is used to score a hypothesis that is used to classify data points.

We now discuss some loss functions which have proven useful for assessing the quality of a hypothesis that is used to classify data points. Unless stated otherwise, the formulas for these loss functions are valid only if the label values are the real numbers -1 and 1 , i.e., for the label space $\mathcal{Y} = \{-1, 1\}$. These formulas need to be modified accordingly if different label values are used. For example, instead of the label space $\mathcal{Y} = \{-1, 1\}$, we could equally well use the label space $\mathcal{Y} = \{0, 1\}$, or $\mathcal{Y} = \{\square, \triangle\}$ or $\mathcal{Y} = \{\text{“Class 1”}, \text{“Class 2”}\}$.

A natural choice for the loss function can be based on the requirement that a reasonable hypothesis should deliver correct classifications, $\hat{y} = y$ for any data point. This suggests to learn a hypothesis $h(\mathbf{x})$ with small 0/1 loss

$$L((\mathbf{x}, y), h) := \begin{cases} 1 & \text{if } y \neq \hat{y} \\ 0 & \text{else,} \end{cases} \quad \text{with } \hat{y} = \begin{cases} 1 & \text{if } h(\mathbf{x}) \geq 0 \\ 0 & \text{else.} \end{cases} \quad (7)$$

Figure 9 illustrates the 0/1 loss (7) for a data point with features \mathbf{x} and label $y = 1$ as a function of the value $h(\mathbf{x})$. The 0/1 loss is equal to zero if the hypothesis yields a correct classification $\hat{y} = y$. For a wrong classification $\hat{y} \neq y$, the 0/1 loss is equal to one.

The 0/1 loss (7) is conceptually appealing when data points are interpreted as realizations of i.i.d. random variables (RVs) with common probability distribution $p(\mathbf{x}, y)$ (see Section 3.1). Given m realizations $(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})$ of these i.i.d. RVs, with high probability

$$(1/m) \sum_{r=1}^m L((\mathbf{x}^{(r)}, y^{(r)}), h) \approx p(y \neq h(\mathbf{x})) \quad (8)$$

for sufficiently large sample size m .

A precise formulation of the approximation (8) can be obtained from the law of large numbers [35, Section 1]. We can apply the law of large numbers since the loss values $L((\mathbf{x}^{(r)}, y^{(r)}), h)$ are realizations of i.i.d. RVs. It is customary to indicate the average 0/1 loss of a hypothesis h indirectly via the accuracy $1 - (1/m) \sum_{r=1}^m L((\mathbf{x}^{(r)}, y^{(r)}), h)$.

In view of (8), the 0/1 loss (7) seems a very natural choice for assessing the quality of a classifier if our goal is to enforce correct classifications $\hat{y} = y$. The appealing statistical aspects of the 0/1 loss come at the expense of an increased computational complexity. Indeed, for a given data point (\mathbf{x}, y) , the 0/1 loss (7) is non-convex and non-differentiable when viewed as a function of the hypothesis h . Thus, ML methods that use the 0/1 loss to learn a hypothesis require advanced optimization methods to solve the resulting learning problem (see [1, Sec. 3.8.]).

To avoid the non-convexity of the 0/1 loss (7) we can approximate it by a convex “surrogate” loss function. One such approximation of the 0/1 loss is the hinge loss

$$L((\mathbf{x}, y), h) := \max\{0, 1 - yh(\mathbf{x})\}. \quad (9)$$

Figure 9 depicts the hinge loss (9) as a function of the hypothesis $h(\mathbf{x})$. The



Figure 8: Two examples of loss functions that are convex functions of model parameters. If we can only explore the local shape of a non-differentiable loss function around the current model parameters, it is impossible to decide if we are currently in (a) or (b).

hinge loss (9) becomes minimal (zero) for a correct classification ($\hat{y} = y$) with sufficient confidence $|h(\mathbf{x})| \geq 1$. For a wrong classification ($\hat{y} \neq y$), the hinge loss increases monotonically with increasing confidence $|h(x)|$ in the wrong classification.

While the hinge loss avoids the non-convexity of the 0/1 loss, it still is a non-differentiable function of $h(\mathbf{x})$. Optimizing a non-differentiable loss function is typically harder, requiring more iterations of iterative methods, compared to optimizing a differentiable loss function [36, 37]. Figure 8 depicts an example of a differentiable and a non-differentiable convex loss function, respectively. Iterative optimization methods can typically only explore the loss function locally around a current choice for the model parameters. For the non-differentiable loss function in Figure 8, such algorithms cannot determine how far the current model parameters are from the minimum.

Besides the 0/1 loss and the hinge loss, another popular loss function for binary classification problems is the logistic loss

$$L((\mathbf{x}, y), h) := \log(1 + \exp(-yh(\mathbf{x}))). \quad (10)$$

The logistic loss (10) is used in logistic regression [1, Ch. 3] to measure the usefulness of a linear hypothesis map $h^{(\mathbf{w})}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$.

Figure 9 depicts the logistic loss (10) as a function of the hypothesis $h(\mathbf{x})$. The logistic loss (10) is a convex and differentiable function of $h(\mathbf{x})$. For a correct classification ($\hat{y} = y$), the logistic loss (10) decreases monotonically with increasing confidence $h(\mathbf{x})$. For a wrong classification ($\hat{y} \neq y$), the logistic loss increases monotonically with increasing confidence $|h(\mathbf{x})|$ in the wrong classification.

Both, the hinge loss (9) as well as the logistic loss (10) are convex functions of the weight vector $\mathbf{w} \in \mathbb{R}^d$ for a linear hypothesis map $h^{(\mathbf{w})}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$. In contrast to the hinge loss, the logistic loss (10) is also a differentiable function of the \mathbf{w} .

The convex and differentiable logistic loss function can be minimized using gradient-based methods (see Section 5). In contrast, we cannot use gradient-based methods to minimize the hinge loss since it is not differentiable (it does not have a gradient everywhere). However, we can apply a generalization of GD which is known as subgradient descent [37]. Subgradient descent is obtained from GD by generalizing the concept of a gradient to that of a subgradient [36, 38].

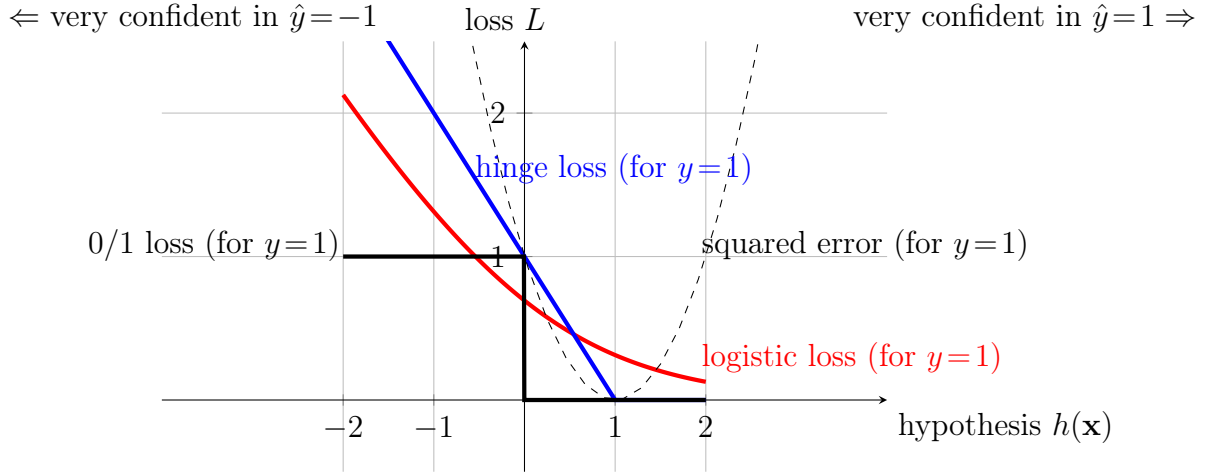


Figure 9: The solid curves depict three widely-used loss functions for measuring the performance of a binary classifier $h^{(\mathbf{w})}(\mathbf{x})$ with label space $\mathcal{Y} = \{-1, 1\}$. A data point with features \mathbf{x} and label $y = 1$ is classified as $\hat{y} = 1$ if $h(\mathbf{x}) \geq 0$ and classified as $\hat{y} = -1$ if $h(\mathbf{x}) < 0$. We can interpret the absolute value $|h(\mathbf{x})|$ as the confidence in the classification \hat{y} . The more confident we are in a correct classification ($\hat{y} = y = 1$), i.e, the more positive $h(\mathbf{x})$, the smaller the loss. Note that each of the three loss functions for binary classification tends monotonically towards 0 for increasing $h(\mathbf{x})$. The dashed curve depicts the squared error loss (4), which increases for increasing $h(\mathbf{x})$.

2.4 Exercises

Exercise 2.1. Data Imputation (10 points). Consider the dataset available under https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_wine.html#sklearn.datasets.load_wine. This dataset consists of data points that represent wine samples. Each data point is characterized by $d = 13$ features x_1, \dots, x_d . Let us assume that the feature value x_1 is missing for several data points. We could then define feature x_1 as the label of a data point and try to predict it from the remaining features. Try to find coefficients w_2, \dots, w_{13} such that $\underbrace{y}_{x_1} \approx \sum_{j=2}^d w_j x_j$.

Exercise 2.2. Design Aspects (5 points). What are three main design aspects for the hypothesis space (model) of a ML method?

Exercise 2.3. 0/1 loss (5 points). Why is the 0/1 loss rarely used as a criterion for learning (optimizing) a classifier?

Exercise 2.4. Perfect Fit ? (20 points). For a prescribed sample size m , can you craft m data points characterized by feature vectors $\mathbf{x}^{(r)} \in \mathbb{R}^d$ and binary label $y^{(r)} \in \{-1, 1\}$ such that it is possible to find a linear hypothesis $h^{(\mathbf{w})}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ that correctly predicts the label of each data point in \mathcal{D} , i.e.,

$$\begin{aligned} h^{(\mathbf{w})}(\mathbf{x}^{(r)}) &\geq 0 \text{ for each data point with } y^{(r)} = 1, \text{ and} \\ h^{(\mathbf{w})}(\mathbf{x}^{(r)}) &< 0 \text{ for each data point with } y^{(r)} = -1 ? \end{aligned} \quad (11)$$

Exercise 2.5. Formulas for Logistic loss (20 points). Modify the formula (10) for the logistic loss, which is valid only for label space $\mathcal{Y} =$

$\{-1, 1\}$, to measure the quality of hypothesis map for a binary classification problem with label space $\mathcal{Y} = \{\square, \triangle\}$.

Exercise 2.6. Decision Tree as Linear Model (20 points). A decision tree is an important example for a non-linear model, i.e., a hypothesis space that contains non-linear hypothesis maps. Consider data points characterized by a single features $x \in \mathbb{R}$ and numeric label y . Discuss how a decision tree, with maximum tree depth 4 and decision nodes implementing threshold tests with thresholds 10, 5, and -4 can be obtained by combining a feature map with a linear model $\mathcal{H}^{(d)}$ (d can be larger than 1).

3 A Design Principle for ML

ML methods learn a hypothesis out of a given hypothesis space (or model) \mathcal{H} that incurs minimum loss when applied to arbitrary data points. To make this informal goal precise we need to specify what we mean by an “arbitrary” data point. One approach to define the notion of “arbitrary” data point is to use a probabilistic model.

Section 3.1 introduces a widely used probabilistic model which is referred to as the i.i.d. assumption. NOTE: the i.i.d. assumption is only a conceptual device for studying ML methods. In practice, we rarely know if the observed data points are really i.i.d. (unless we generate them with an ideal random number generator).

The i.i.d. assumption allows to define the expected loss or risk as a performance measure for a given hypothesis. Section 3.2 introduces ERM as a main design principle for ML. The simple idea of ERM is to approximate the risk of a hypothesis by its average loss on a training set.

Section 3.3 discusses conditions for the training set and hypothesis space such that ERM is likely to succeed.

3.1 The i.i.d. Assumption

Consider a ML application that involves a dataset

$$\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}.$$

The (arguably) most widely-used probabilistic model in ML is to interpret the data points $(\mathbf{x}^{(r)}, y^{(r)})$ as realizations of i.i.d. RVs with a common probability

distribution $p(\mathbf{x}, y)$. This probabilistic model is sometimes referred to as an “i.i.d. assumption” and denoted by

$$(\mathbf{x}^{(r)}, y^{(r)}) \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, y). \quad (12)$$

The probabilistic model (12) enables us to measure the quality of a hypothesis $h \in \mathcal{H}$ by the expected loss or risk [39]

$$\bar{L}(h) := \mathbb{E}\{L((\mathbf{x}, y), h)\} = \int_{\mathbf{x}, y} L((\mathbf{x}, y), h) dp(\mathbf{x}, y). \quad (13)$$

The risk (13) of h is the expected value of the loss $L((\mathbf{x}, y), h)$ incurred when applying the hypothesis h to (the realization of) a random data point with features \mathbf{x} and label y . Note that the computation of the risk (13) requires the joint probability distribution $p(\mathbf{x}, y)$ of the (random) features and label of data points. It seems reasonable to learn a hypothesis h^* that incurs minimum risk,

$$h^* \in \underset{h \in \mathcal{H}}{\operatorname{argmin}} \underbrace{\bar{L}(h)}_{=\mathbb{E}\{L((\mathbf{x}, y), h)\}}. \quad (14)$$

Any hypothesis that solves (14), i.e., that incurs minimal risk, is referred to as a Bayes estimator [39, Chapter 4].

Thus, once we know the underlying probability distribution, the only challenge for learning the optimal hypothesis is the efficient (numerical) solution of the optimization problem (13). Efficient methods to solve the optimization problem (13) include variational methods and random sampling (Monte Carlo) methods [39, 40].

If we do not know the probability distribution underlying the data points, we cannot compute the risk and, in turn, cannot use (14) as a learning

principle. However, we can approximate (or estimate) the risk by an average over the data points in a dataset. We define the empirical risk of a hypothesis $h \in \mathcal{H}$ incurred on a dataset \mathcal{D} as

$$\widehat{L}(h|\mathcal{D}) = (1/m) \sum_{r=1}^m L((\mathbf{x}^{(r)}, y^{(r)}), h). \quad (15)$$

The empirical risk of the hypothesis $h \in \mathcal{H}$ is the average loss on the data points in \mathcal{D} . Note that the empirical risk depends on three components, the loss function $L(\cdot, \cdot)$, the hypothesis h and the (features and labels of the) data points in the dataset \mathcal{D} .

If the data points used to compute the empirical risk (15) are modelled as realizations of i.i.d. RVs whose common probability distribution is $p(\mathbf{x}, y)$, basic results of probability theory tell us that

$$\mathbb{E}\{L((\mathbf{x}, y), h)\} \approx (1/m) \underbrace{\sum_{r=1}^m L((\mathbf{x}^{(r)}, y^{(r)}), h)}_{\stackrel{(15)}{=} \widehat{L}(h|\mathcal{D})} \text{ for sufficiently large } m. \quad (16)$$

The approximation error in (16) can be quantified by different variants of the law of large numbers [4, 35, 41] or large-deviation bounds [42].

3.2 Empirical Risk Minimization

Many ML methods are motivated by (16) which suggests that a hypothesis with small empirical risk (15) will also result in a small expected loss or risk. The minimum risk is achieved by the Bayes estimator of the label y , given the features \mathbf{x} . However, to actually compute the optimal estimator we need the (joint) probability distribution $p(\mathbf{x}, y)$ of features \mathbf{x} and label y . This joint

probability distribution is typically unknown and must be estimated from observed data points.

Many practical ML methods are numerical optimization algorithms to solve ERM

$$\hat{h} := \operatorname{argmin}_{h \in \mathcal{H}} (1/m) \underbrace{\sum_{r=1}^m L((\mathbf{x}^{(r)}, y^{(r)}), h)}_{\hat{L}(h|\mathcal{D})}. \quad (17)$$

The objective function of ERM is the empirical risk (15) which, in turn, approximates the risk $\bar{L}(h)$ via the law of large numbers (16). Note that ERM (17) is parametrized by three components: the dataset \mathcal{D} , the hypothesis space \mathcal{H} and loss function $L(\cdot, \cdot)$. The Python library `scikit-learn` provides implementations of (17) for different choices for the hypothesis space \mathcal{H} and loss function $L(\cdot, \cdot)$ [24].

The objective function of ERM depends on the dataset

$$\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}.$$

In general, $\hat{L}(h|\mathcal{D})$ (viewed as a function of $h \in \mathcal{H}$) depends on each of the feature vectors $\mathbf{x}^{(r)}$ and labels $y^{(r)}$ in \mathcal{D} . By the i.i.d. assumption (12), each data point $(\mathbf{x}^{(r)}, y^{(r)})$ in \mathcal{D} is a realization of a RV with (unknown) probability distribution $p(\mathbf{x}, y)$. This probability distribution crucially affects the computational aspects and statistical aspects of ML methods that use ERM (17).

In general, we do not have (much) control over the validity of the i.i.d. assumption, let alone the underlying probability distribution $p(\mathbf{x}, y)$. However, we can use principled statistical tests to verify if the i.i.d. assumption is likely to be valid [43]. However, we can make different design choices for the

data points (their features \mathbf{x} and label y), hypothesis space \mathcal{H} and the loss function that define ERM. Different ML methods use different choices for these components of ERM which, in turn, result in different computational aspects and statistical aspects of ERM-based ML methods (see [1, Ch. 3]).

By computational aspects of ERM, we mainly refer to the amount of computation required to find (approximate) solutions to (17). In general, this amount depends on the shape of the empirical risk, viewed as a function of the hypothesis (or its parameters) (see Figure 10). This informal claim is made precise in Section 5 where we study the number of iterations required by GD methods for linear regression. For linear regression, the shape of the objective function in (17) is determined by the eigenvalues of a matrix constructed from the feature vectors in (17).

The statistical aspects of (17) include the probability distribution of the solutions to (17) and their deviation from the Bayes estimator (14). The statistical aspects of (17) also include its robustness against violations of the i.i.d. assumption. Different choices for the model and loss function used in (17) yield different levels of robustness for the resulting ML method [44].

For linear regression, which uses a linear model $h^{(\mathbf{w})}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ and squared error loss (4), the generic ERM (17) specializes to

$$\hat{\mathbf{w}} := \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \underbrace{(1/m) \sum_{r=1}^m (\mathbf{w}^T \mathbf{x}^{(r)} - y^{(r)})^2}_{\hat{L}(h^{(\mathbf{w})}|\mathcal{D})}. \quad (18)$$

3.3 How Much Training is Needed?



Figure 10: ERM (17) learns a hypothesis, out of a hypothesis space \mathcal{H} , by minimizing the empirical risk $\hat{L}(h|\mathcal{D})$ over a training set \mathcal{D} . However, the ultimate goal is to learn a hypothesis which predicts well the label of any data point and, in turn, has a small risk $\bar{L}(h)$. Thus, ERM can be expected to deliver a useful hypothesis if the deviation between the empirical risk and the risk is small (uniformly over $h \in \mathcal{H}$).

We have obtained ERM by approximating the risk $\bar{L}(h)$ with the empirical risk $\hat{L}(h|\mathcal{D})$ (see Figure 10). The risk is the ultimate performance criterion but cannot be computed without knowing the underlying probability distribution (see Section 3.1). Instead, we minimize the empirical risk on a given training set constituted by m data points. Thus, we learn a hypothesis \hat{h} by solving ERM (17). Any solution of (17) can be used as the learnt hypothesis \hat{h} .

The solutions of ERM will only deliver a useful hypothesis (with small risk) if two conditions are satisfied:

- the training error $\hat{L}(\hat{h}|\mathcal{D})$ of the learnt hypothesis is small, and,

- the deviation between $\widehat{L}(\hat{h}|\mathcal{D})$ and risk $\bar{L}(\hat{h})$ is small.

We can test the first condition by inspecting the optimal value of ERM. However, we cannot verify the second condition without further assumptions. Indeed, the deviation between empirical risk and risk depends on the statistical properties of the data points and how many of them are used in the training set for ERM. Moreover, the deviation will also depend on the hypothesis space from which ERM learns the hypothesis \hat{h} . We will next analyze this dependency by using a simple probabilistic model that builds on the i.i.d. assumption from Section 3.1.

In what follows, we will interpret any data point, including those in the training set \mathcal{D} of ERM (17), as realizations of i.i.d. RVs with a multivariate normal distribution. In particular, the features $\mathbf{x} \in \mathbb{R}^d$ of data points are realizations of a multivariate normal distribution, $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. The label of any data point is related to the features via a linear Gaussian model,

$$y = \bar{\mathbf{w}}^T \mathbf{x} + \varepsilon, \text{ with } \varepsilon \sim \mathcal{N}(0, \sigma^2), \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (19)$$

The linear Gaussian model (19) postulates that the label of a data point is obtained as a noisy linear combination of its features. We form this linear combination using the entries of a true underlying (but unknown to the ML method) weight vector $\bar{\mathbf{w}} \in \mathbb{R}^d$. The additive noise ε in (19) is a realization of a zero-mean Gaussian RV with variance σ^2 .

We learn a linear hypothesis \hat{h} via the following instance of ERM (17),

$$\hat{h} := \underset{h(\mathbf{w}) \in \mathcal{H}^{(d')}}{\operatorname{argmin}} (1/m) \sum_{r=1}^m (y^{(r)} - h(\mathbf{w})(\mathbf{x}^{(r)}))^2. \quad (20)$$

The linear model $\mathcal{H}^{(d')}$ uses the first $d' \in \{1, \dots, d\}$ features $x_1, \dots, x_{d'}$ of a data point. Thus, we allow the ERM instance (20) to use only a subset of

the d available features $\mathbf{x} = (x_1, \dots, x_d)^T$. The feature vectors $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$ of data points are modelled as i.i.d. realizations $\mathbf{x}^{(r)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$.

Note that the objective function in (20) involves the features and labels of data points in \mathcal{D} . Since we interpret those as realizations of RVs, also the learnt hypothesis \hat{h} and its risk $\bar{L}(\hat{h})$ become realizations of RVs. For a sufficiently large training set, such that $m > d' + 1$, the expectation of the risk $\bar{L}(\hat{h})$ is [1, Ch. 6.3],

$$\mathbb{E}\{\bar{L}(\hat{h})\} = B^2 + \underbrace{(B^2 + \sigma^2)d'/(m - d' - 1) + \sigma^2}_{\text{variance } v}. \quad (21)$$

Here, we used the bias term

$$B^2 = \sum_{j=d'+1}^d \bar{w}_j^2. \quad (22)$$

According to (21), the expected risk $\mathbb{E}\{\bar{L}(\hat{h})\}$ depends on three components:

- the bias term B^2 , which depends on the true underlying weight vector $\bar{\mathbf{w}}$ (not under control of ML engineer) and the dimension d' of the linear model $\mathcal{H}^{(d')}$. The hyper-parameter d' can be chosen freely from $\{1, \dots, d\}$.
- the variance term V which is small if m is significantly larger than d' . As a rule of thumb, the sample size m should be ten times larger than the effective dimension d' of the linear model,

$$m \geq 10 \times d'. \quad (23)$$

- the noise floor σ^2 which is due to the intrinsic noise in the (assumed) relation (19) between features and label of a data point. Under our probabilistic model (19), no ML method can learn a hypothesis with (expected) risk smaller than σ^2 .

We can interpret the component $B^2 + V$ in (21) as a measure for the estimation error incurred by using (20) (via its equivalent form (24)) as an estimator for the true underlying $\bar{\mathbf{w}}$ in (19). If we use the linear model $\mathcal{H}^{(d')}$ with $d' = d$, which implies that the bias vanishes $B^2 = 0$, the estimation error becomes $\sigma^2 d' / (m - d' - 1)$ [45, 46].

The third component σ^2 in (21) arises from the noise contained in the true label y (according to our probabilistic model (19)), which we predict by the learnt hypothesis $h^{(\hat{\mathbf{w}})}(\mathbf{x})$.

Note that (21) is the expected risk of a specific ML method, i.e., using ERM (20) to learn the hypothesis \hat{h} . One might wonder if there is a better ML method that could learn a better linear hypothesis $h'(\mathbf{x}) = (\mathbf{w}')^T \mathbf{x}$, having smaller expected risk: $\mathbb{E}\{\bar{L}(h')\} < \mathbb{E}\{\bar{L}(\hat{h})\}$? It turns out that this is not the case and that the hypothesis learnt via (20) achieves minimum expected risk. Indeed, the expected risk of any linear hypothesis learnt from m data points that follow the probabilistic model (19) is lower bounded by (21) [47].

To summarize, using a specific probabilistic model (19) we analyzed the expected risk $\mathbb{E}\{\bar{L}(\hat{h})\}$ of the linear hypothesis \hat{h} learnt via ERM (20). We decomposed $\mathbb{E}\{\bar{L}(\hat{h})\}$ into different components, one of them being proportional to the ratio $d' / (m - d' - 1)$.

To ensure a small $\mathbb{E}\{\bar{L}(\hat{h})\}$, the size m of the training set used by ERM (20) should be significantly larger (say, ten times) than the dimension d' of the

linear model used in (20). While our analysis only applies to linear models, their results provide useful intuition for non-linear models such as ANN. For example, we can (locally) approximate non-linear models by linear models to determine their effective dimension d' [1, Ch. 2].

3.4 Exercises

Exercise 3.1. Why Probability Theory? (5 points) Try to justify the use of probabilistic models in ML. In particular, why is it useful to interpret data points (such as the rows in a spreadsheet) as realizations of RVs.

Exercise 3.2. Trivial Objective in ERM (5 points) Consider some arbitrary dataset which is used to learn a hypothesis out of \mathcal{H} via ERM (17). For which combinations of loss function and hypothesis space does ERM not depend at all on the dataset.

Exercise 3.3. Data Leakage via ERM (50 points) Consider a dataset \mathcal{D} with m data points, each characterized by a single integer-valued feature $x \in \mathbb{Z}$ and a real-valued label $y \in \mathbb{R}$. We try to learn a hypothesis which maps every integer $x \in \mathbb{Z}$ to a predicted label $h(x) \in \mathbb{R}$. The hypothesis space \mathcal{H} is constituted by all maps $h : \mathbb{Z} \rightarrow \mathbb{R}$, which could be interpreted as discrete-time signals [48]. The learning is based on minimizing the average squared error loss $\widehat{L}(h|\mathcal{D})$ over all hypothesis maps. Is it possible to perfectly recover the dataset just from knowing the values of $\widehat{L}(h|\mathcal{D})$ for all $h \in \mathcal{H}$?

Exercise 3.4. ERM for linear regression (10 points) Show that the solution of (20) is the linear hypothesis $h^{(\widehat{\mathbf{w}})}(\mathbf{x}) = (\widehat{\mathbf{w}})^T \mathbf{x}$ with weight vector $\widehat{\mathbf{w}} \in \mathbb{R}^{d'}$ being a solution to

$$\widehat{\mathbf{w}} \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^{d'}} \sum_{r=1}^m \left(y^{(r)} - \sum_{j=1}^{d'} w_j x_j^{(r)} \right)^2. \quad (24)$$

Exercise 3.5. Uniqueness in ERM for linear regression (20 points) Discuss sufficient conditions on the training set such that (24) has a unique solution.

Exercise 3.6. Estimation Error and Expected Risk (10 points)

Consider a set of m data points whose features and label are realizations of RVs with a probability distribution defined via (19). We use ERM (24), with $d' = d$ to learn the weights of a linear hypothesis map. We can interpret the solution $\hat{\mathbf{w}}$ of (24) as an estimator for the true underlying weights $\bar{\mathbf{w}}$ in (19). The performance of this estimator can be measured via the estimation error $\mathbb{E}\{\|\hat{\mathbf{w}} - \bar{\mathbf{w}}\|_2^2\}$. Discuss the relation between this estimation error and the expected risk $\mathbb{E}\{\bar{L}(h(\hat{\mathbf{w}}))\}$.

Exercise 3.7. Fundamental Limit (10 points) Consider data points characterized by a feature vector $\mathbf{x} = (x_1, \dots, x_5)^T$ and a numeric label $y \in \mathbb{R}$. Assume that the label of a data point is related to its features via a linear Gaussian model (see (19)),

$$y = \sum_{r=1}^5 \bar{w}_j x_j + \varepsilon.$$

Here, \bar{w}_j are fixed weights and $\varepsilon \sim \mathcal{N}(0, 5)$ is a realization of a zero-mean Gaussian RV with variance 5. Is there a ML method that can learn a hypothesis with risk no larger than 5 from an training set with $m = 10$ data points (which are i.i.d. realizations of probability distribution $p(\mathbf{x}, y)$)?

Exercise 3.8. Perfect Fit ? (10 points) Consider the dataset \mathcal{D} that consists of $m = 20$ data points, each characterized by a single numeric feature and label given as

$$x^{(r)} = \sin(4\pi r/m), \text{ and } y^{(r)} = \cos(\pi r/(2m)) \text{ for } r = 1, \dots, m.$$

If we use a sufficiently large hypothesis space \mathcal{H} , is it possible to find a

hypothesis map $\hat{h} \in \mathcal{H}$ that perfectly fits the data points in \mathcal{D} in the sense of

$$\hat{h}(x^{(r)}) = y^{(r)} \text{ for each } r = 1, \dots, m ?$$

4 Regularization

The idea of ERM is to approximate the expected loss or risk of a hypothesis by the empirical risk $\widehat{L}(h|\mathcal{D})$ (17) over a training set \mathcal{D} . Whenever this approximation fails to be accurate, the solution of ERM might be a hypothesis that performs poorly outside the training set used in (17). Figure 10 illustrates the difference between the risk and the empirical risk as its approximation.

Section 4.1 discusses overfitting as an extreme case where the empirical risk (the training error) of a learnt hypothesis is negligible while its risk is unacceptably large. Section 4.2 introduces regularization techniques that reduce the discrepancy between empirical risk and risk by modifying either the data, model or loss function used by a ERM-based ML method. Section 4.3 explains how to use regularization to couple the training of different local models. This coupling of the training of local models is at the heart of the FL methods discussed in part II.

4.1 Overfitting

Consider a ERM-based ML method that learns a hypothesis from a model \mathcal{H} by minimizing the average loss incurred on a training set (see (17)). The behaviour of such a method can often be (partially) predicted by comparing the size of the training set with the size of the model \mathcal{H} . But how to measure the size of the training set and the size of a model ?

It seems quite natural to measure the size of the training set used in ERM (17) is the number m of data points it contains. The number of data points in a dataset is useful measure of its size only if the data points can



Figure 11: Time series of daily vaccinations which we might use to define data points.

be (accurately) modelled as realizations of i.i.d. RVs. However, there might be applications where the number of data points in a dataset is not a very useful measure for its size. Consider data points obtained by the values of a time series, such as the daily vaccinations collected during each day (see Figure 11)

If the time series (or its trend) does not change too rapidly, several consecutive data points will be almost identical (or at least strongly correlated). To measure the size of a dataset obtained by sampling a (correlated) time series, we should also take into account the typical duration during which the time series is approximately constant [49]. However, in what follows we tacitly assume that data points can be reasonably well modelled as realizations of i.i.d. RVs such that their number is a useful measure for the size of dataset.

It seems far less obvious how to measure the size of a hypothesis space

that contains infinitely many hypothesis maps. Probably the most widely used concept to measure the size of infinite hypothesis spaces is the Vapnik–Chervonenkis (VC) dimension. A simplified variant of the VC dimension, referred to as the effective dimension $d_{\text{eff}}(\mathcal{H})$ has been introduced in [1, Ch. 2]. Loosely speaking, $d_{\text{eff}}(\mathcal{H})$ is the maximum number of data points (obtained as i.i.d. realizations of a “well-behaved” probability distribution) such that there is always at least one $h \in \mathcal{H}$ that perfectly fits these data points.

It is instructive to evaluate the effective dimension for the linear model $\mathcal{H}^{(d)}$ that consists of linear hypothesis maps that read in d features of a data point. Here, it can be shown that $d_{\text{eff}}(\mathcal{H}^{(d)}) = d$ (see Exercise 4.5). Another important example is polynomial regression for data points characterized by a single feature x which uses the hypothesis space [1, Sec. 3.2.]

$$\mathcal{H}_d^{(\text{poly})} := \{h(x) = \sum_{j=0}^{d-1} x^j w_j\}.$$

The effective dimension of $\mathcal{H}_d^{(\text{poly})}$ is $d_{\text{eff}}(\mathcal{H}_d^{(\text{poly})}) = d$.

The ratio $d_{\text{eff}}(\mathcal{H})/m$ between effective dimension of the model \mathcal{H} and the size m of the training set underlying ERM is often a useful indicator for overfitting. For example, if $d_{\text{eff}}(\mathcal{H})/m > 1$ then ERM (17) is likely to overfit the training set. Indeed, by the very definition of the effective dimension $d_{\text{eff}}(\mathcal{H})$, for (almost) any given training set

$$\{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\} \text{ of size } m \leq d_{\text{eff}}(\mathcal{H}),$$

there is a hypothesis $\hat{h} \in \mathcal{H}$ that perfectly fits the training set, i.e., $y^{(r)} = \hat{h}(\mathbf{x}^{(r)})$ for $r = 1, \dots, m$. To sum up, ERM with a training set of size $m \leq d_{\text{eff}}(\mathcal{H})$, will (almost always) learn a hypothesis \hat{h} with vanishing training

error.¹

Figure 12 depicts the dependence of the training error $\widehat{L}(\hat{h}|\mathcal{D})$ and the risk $\bar{L}(\hat{h})$ of the hypothesis \hat{h} delivered by ERM using a model with effective dimension $d_{\text{eff}}(\mathcal{H})$ and a training set \mathcal{D} of size m .

Consider a fixed size m of the training set and let us vary the effective dimension $d_{\text{eff}}(\mathcal{H})$ of the model, e.g., by using a varying number of features or using varying number of hidden layers for a ANN. If we increase $d_{\text{eff}}(\mathcal{H})$, the training error typically decreases as we can choose from a larger hypothesis to fit the data points in the training set.

As soon as $d_{\text{eff}}(\mathcal{H}) > m$, the training error vanishes but at the same time, the risk of the learnt hypothesis increases. This regime will typically lead to overfitting by learning a hypothesis that perfectly predicts the labels of data points in the training set but delivers poor predictions for data points outside the training set.

If we use a too small model with $d_{\text{eff}}(\mathcal{H}) \ll m$, we might not be able to find any suitable hypothesis $h \in \mathcal{H}$ that accurately predicts the labels of data points from their features. Here, we ERM will learn a hypothesis \hat{h} with large training error and large risk. We say that the ERM-based method is underfitting the training set.

Note that ensuring $d_{\text{eff}}(\mathcal{H})/m \leq 1$ requires to use at least $m \geq d_{\text{eff}}(\mathcal{H})$ data points in the training set for ERM. However, it might be infeasible to collect such an amount of data points for modern ML methods that use high-dimensional models (such as deep ANNs) with $d_{\text{eff}}(\mathcal{H})$ being in the order

¹If we use a loss function that is zero when the prediction $\hat{y} = h(\mathbf{x})$ coincides with the label y of a data point.



Figure 12: The typical dependence of the training error and the risk of the hypothesis \hat{h} delivered by ERM using a model with effective dimension $d_{\text{eff}}(\mathcal{H})$ and a training set of size m .

of billions [33].

The next section gives a glimpse on regularization techniques for modifying a given model \mathcal{H} and training set $\mathcal{D}^{(\text{train})}$ to decrease the ratio $d_{\text{eff}}(\mathcal{H})/m$. These techniques either reduce $d_{\text{eff}}(\mathcal{H})$ (model pruning) or increase the size m of the training set (data augmentation) without the need to collect new labeled data points.

4.2 Regularization via Data, Model and Loss

Regularization techniques modify ERM (17) to favour a hypothesis that does not perform worse outside the training set. According to Section 4.1, this requires that the ratio $d_{\text{eff}}(\mathcal{H})/m$ between the effective dimension of the model and the size of the training set is sufficiently small (below one).

Obviously, we can decrease this ratio either by using a smaller model (“model pruning”) or by using a larger training set.

In general, we can implement regularization via each of the three main ML components, either as

- data augmentation (see Figure 13): we enlarge the training set \mathcal{D} in (17) by adding new data points obtained by perturbing features or labels.
- model pruning: we modify (17) by reducing the optimization domain, which is the hypothesis space \mathcal{H} , to a (tiny) subset $\mathcal{H}' \subseteq \mathcal{H}$,

$$\hat{h} := \operatorname{argmin}_{h \in \mathcal{H}'} \underbrace{(1/m) \sum_{r=1}^m L((\mathbf{x}^{(r)}, y^{(r)}), h)}_{\hat{L}(h|\mathcal{D})}. \quad (25)$$

- loss penalization: modify the loss function in (17) by adding a scaled regularization term,

$$\hat{h} := \operatorname{argmin}_{h \in \mathcal{H}} \underbrace{(1/m) \sum_{r=1}^m L((\mathbf{x}^{(r)}, y^{(r)}), h)}_{\hat{L}(h|\mathcal{D})} + \lambda \mathcal{R}\{h\} \quad (26)$$

Two well-known instances of (26) are ridge regression and least absolute shrinkage and selection operator (Lasso). These two methods learn a linear hypothesis map $h^{(\mathbf{w})} = \mathbf{w}^T \mathbf{x}$ by minimizing (26) for the squared error loss and regularizer $\mathcal{R}\{h^{(\mathbf{w})}\} = \|\mathbf{w}\|_2^2$ (ridge regression) or $\mathcal{R}\{h^{(\mathbf{w})}\} = \|\mathbf{w}\|_1$ (Lasso).

These three approaches to regularization are essentially equivalent. As a case in point, Figure 14 illustrates the equivalence between data augmentation

and loss penalization. Indeed, the empirical risk incurred over an enlarged dataset obtained by data augmentation coincides with the empirical risk on the original dataset using a penalized loss function. This penalty term corresponds to the average loss on the augmented data points.

Equivalence ridge regression and data augmentation. Consider linear regression methods that learn the weight vector of a linear hypothesis map $h^{(\mathbf{w})}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ by minimizing the average squared error loss on a training set. We augment \mathcal{D} by adding B realizations of i.i.d. Gaussian RVs with zero mean and covariance matrix $\sigma^2 \mathbf{I}$. This results in the augmented dataset \mathcal{D}' which contains each data point of \mathcal{D} along with its B perturbations. As shown in [1, Sec. 7.3.], for $B \rightarrow \infty$, the average squared error loss $\widehat{L}(h^{(\mathbf{w})}|\mathcal{D}')$ on the augmented dataset \mathcal{D}' tends towards $\widehat{L}(h^{(\mathbf{w})}|\mathcal{D}) + \sigma^2 \|\mathbf{w}\|_2^2$. Note that $\widehat{L}(h^{(\mathbf{w})}|\mathcal{D}) + \sigma^2 \|\mathbf{w}\|_2^2$ is an instance of loss penalization (see (26)) using the regularizer $\mathcal{R}\{h^{(\mathbf{w})}\} = \|\mathbf{w}\|_2^2$. Note that $\widehat{L}(h^{(\mathbf{w})}|\mathcal{D}) + \sigma^2 \|\mathbf{w}\|_2^2$ is precisely the average loss used by ridge regression to learn a linear hypothesis map [1, Sec. 3].

Equivalence ridge regression and model pruning. Consider a training set that consists of m data points $(\mathbf{x}^{(r)}, y^{(r)})$, for $r = 1, \dots, m$. Ridge regression learns the weight vector of a linear hypothesis map $h^{(\mathbf{w})}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ by solving the unconstrained optimization problem

$$\min_{\mathbf{w} \in \mathbb{R}^d} (1/m) \sum_{r=1}^m (y^{(r)} - \mathbf{w}^T \mathbf{x}^{(r)})^2 + \lambda \|\mathbf{w}\|_2^2 \text{ with some } \lambda > 0. \quad (27)$$

Let us denote the unique solution of (27) by $\widehat{\mathbf{w}}^{(\lambda)}$ and its squared Euclidean norm by $\rho^{(\lambda)} := \|\widehat{\mathbf{w}}^{(\lambda)}\|_2^2$. Trivially, $\widehat{\mathbf{w}}^{(\lambda)}$ is also a solution to the constrained

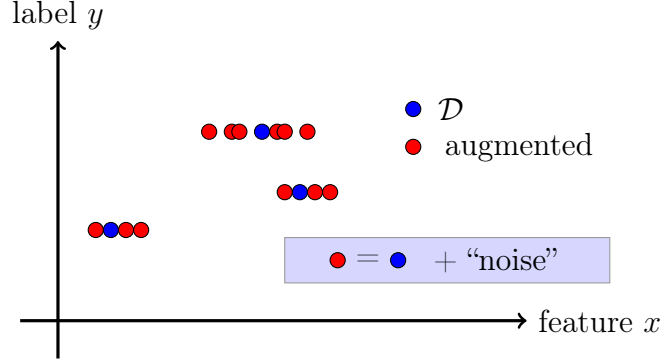


Figure 13: Data augmentation adds new data points to a given dataset \mathcal{D} . These new data points are obtained by perturbing features and labels of the data points in \mathcal{D} .

optimization problem

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d} (1/m) \sum_{r=1}^m (y^{(r)} - \mathbf{w}^T \mathbf{x}^{(r)})^2 \\ \text{s.t. } \|\mathbf{w}\|_2^2 \leq \rho^{(\lambda)}. \end{aligned} \quad (28)$$

Note that (28) is actually an instance of ERM (17) using the pruned linear model $\mathcal{H}^{(\lambda)} := \{h^{(\mathbf{w})} \in \mathcal{H}^{(d)} : \|\mathbf{w}\|_2^2 \leq \rho^{(\lambda)}\} \subseteq \mathcal{H}^{(d)}$.

4.3 Federated Learning via Regularization

Section 9 will present FL algorithms that use regularization to couple the training of several different local models. The main building block of these algorithms is the regularized empirical risk minimization (RERM) instance

$$\hat{\mathbf{w}} := \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \underbrace{\widehat{L}(h^{(\mathbf{w})}|\mathcal{D})}_{\text{training error}} + (\lambda/2) \underbrace{\|\mathbf{w}' - \mathbf{w}\|^2}_{\text{deviation from } \mathbf{w}'} . \quad (29)$$



Figure 14: Equivalence of data augmentation and loss penalization techniques for regularization.

The model parameters $\hat{\mathbf{w}}$ obtained from (29) balances between two (conflicting) goals. On one hand, we want to minimize the empirical risk $\hat{L}(h^{(\mathbf{w})}|\mathcal{D})$. On the other hand, the learnt model parameters $\hat{\mathbf{w}}$ should not deviate too much from given reference model parameters \mathbf{w}' .

The RERM (29) will be used as the elementary computation of FL algorithms in Section 9. In these algorithms, the vector \mathbf{w}' in (29) represents (an aggregation of) model parameters for different (but similar) learning tasks.

Section 6 introduces the empirical graph of a collection of local datasets and corresponding learning tasks that are related by a notion of pairwise similarities. The (weighted) edges of the empirical graph represent (the level of) similarities between local datasets and corresponding learning tasks. Each node of the empirical graph solves RERM (29) to update its current local model parameters. Here, the vector \mathbf{w}' in (29) is an aggregation of the current local model parameters at neighbouring nodes (see Section 9).

Note that solving RERM (29) is to evaluate the proximity operator

$\mathbf{prox}_{f,\lambda}(\mathbf{w}')$ of the empirical risk $f(\mathbf{w}) = \widehat{L}(h^{(\mathbf{w})}|\mathcal{D})$, viewed as a function of the model parameters \mathbf{w} . This interpretation provides a strong conceptual link between the FL methods developed in Section 7.2 and proximal algorithms [50].

4.4 Exercises

Exercise 4.1. Regularization via data augmentation (20 points).

Consider the RERM (29) using some labeled dataset

$$\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\},$$

linear model $\mathcal{H}^{(d)}$ and squared error loss (4). Show that (29) is equivalent to plain ERM (17) for linear regression applied to another dataset \mathcal{D}' which is obtained by adding data points, with carefully chosen features and labels, to the original dataset \mathcal{D} .

Exercise 4.2. Proximity operator for linear regression (10 points).

Consider the update (29) for the specific choice of squared error loss and linear model. Show that this special case of (29) has always a unique solution.

Exercise 4.3. Uniqueness of Ridge regression (10 points). Provide a rigorous proof for the existence and uniqueness of the solution to the unconstrained optimization problem (27).

Exercise 4.4. Ridge regression as model pruning (10 points). Provide a rigorous proof for the equivalence (in both directions!) between the unconstrained optimization problem (27) and the constrained optimization problem (28).

Exercise 4.5. Effective Dimension of Linear model (100 points).

Consider data points with d numeric features $\mathbf{x} = (x_1, \dots, x_d)^T \in \mathbb{R}^d$ and a numeric label $y \in \mathbb{R}$. Show that with probability one there is a linear hypothesis map $h(\mathbf{x}) = \sum_{j=1}^d w_j x_j$ that perfectly fits $m \leq d$ data points,

$$(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)}) \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, y). \quad (30)$$

In other words, these data points are realizations of i.i.d. RVs with a common probability distribution $p(\mathbf{x}, y)$ that is required to be continuous but arbitrary otherwise. The requirement for $p(\mathbf{x}, y)$ to be continuous is important. Indeed, try to find a probability distribution $p(\mathbf{x}, y)$ that is not continuous and for which it is impossible (with probability one) to find a linear map that perfectly fits the data points in (30).

5 Gradient Methods

Section 3.2 formulated ML as ERM, which is a special case of an optimization problem. The properties of this optimization problem crucially depend on the design choices for the hypothesis space and loss function of a ML method. This section explains gradient-based methods for solving the ERM arising from combinations of a parametrized model and a differentiable loss function [1, Ch. 3].

We will illustrate the key principles of gradient-based methods using linear regression, i.e., learning a linear hypothesis map to predict the numeric label of a data point from its numeric features. The restriction to linear regression allows for a rather complete analysis of gradient-based methods which also provides intuition for the behaviour of gradient-based methods applied to other combinations of a parametrized hypothesis space (such as ANNs) and loss function (such as logistic loss).

The core of any gradient-based methods is a gradient step (see Section 5.1) which updates a current choice for model parameters by the scaled negative gradient of the empirical risk (viewed as a function of model parameters). We can interpret the gradient step as the constrained minimization of a local linear approximation of the empirical risk. Another interpretation of the gradient step is that of minimizing a local approximation to the objective function that is a quadratic function (see Exercise 5.10).

Some ML methods require the model parameters to belong to a subset (a constraint set) of the Euclidean space. The resulting ERM is then a constrained optimization problem. Projected GD (see Section 5.2) handles such constraints on the model parameters by projecting the result of a gradient

step back into the constraint set.

Section 5.3 discusses the effects of inexact gradient steps due to different types of errors in the gradient computation. Numerical errors might arise from constraints on computational resources such as a finite communication rate between computational units. Approximation errors might arise when using different objective functions for the implementation and for the analysis of GD methods. This might seem strange at first, but conceptually it might be useful to allow errors in order to extend the analysis of GD beyond the (important but limited) class of smooth and strongly convex objective functions.

Note that the objective function of ERM is an approximation to the risk (viewed as a function of the model parameters). In particular, the gradient of the empirical risk is a stochastic approximation (or “estimate”) for the gradient of the risk. Stochastic gradient descent (SGD) methods (see Section 5.4) are obtained from GD by taking into account the stochastic nature of the gradient errors.

5.1 Gradient Descent

Let us rewrite ERM (18) for linear regression as

$$\hat{\mathbf{w}} \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) := (1/m) \sum_{r=1}^m (y^{(r)} - \mathbf{w}^T \mathbf{x}^{(r)})^2. \quad (31)$$

Our goal is to compute (at least approximately) an optimal parameter vector $\hat{\mathbf{w}}$ that achieves the minimum objective value in (31). Note that the objective function $f(\mathbf{w})$ in (31) is a differentiable function of the model parameters \mathbf{w} . In particular, we can compute a gradient $\nabla f(\mathbf{w})$ for any model parameters



Figure 15: We can approximate a differentiable function $f(\mathbf{w})$ locally around a point $\mathbf{w}^{(r)} \in \mathbb{R}^d$ using a hyperplane. The normal vector $\mathbf{n} = (\nabla f(\mathbf{w}^{(r)}), -1)^T \in \mathbb{R}^{d+1}$ of this approximating hyperplane is determined by the gradient $\nabla f(\mathbf{w}^{(r)})$ [5].

$\mathbf{w} \in \mathbb{R}^d$ [51],

$$\nabla f(\mathbf{w}) = -(2/m) \sum_{r=1}^m \mathbf{x}^{(r)} (y^{(r)} - \mathbf{w}^T \mathbf{x}^{(r)}). \quad (32)$$

We use the gradient (32) to construct a sequence of model parameters $\mathbf{w}^{(0)}, \mathbf{w}^{(1)}, \dots$. To this end, we start from an initial choice for the model parameters $\mathbf{w}^{(0)} \in \mathbb{R}^d$ and repeat the gradient step

$$\begin{aligned} \mathbf{w}^{(r+1)} &:= \mathbf{w}^{(r)} - \alpha_r \nabla f(\mathbf{w}^{(r)}) \\ &\stackrel{(32)}{=} \mathbf{w}^{(r)} + (2\alpha_r/m) \sum_{r=1}^m \mathbf{x}^{(r)} \left(\underbrace{y^{(r)}}_{\text{label}} - \underbrace{(\mathbf{w}^{(r)})^T \mathbf{x}^{(r)}}_{\text{prediction}} \right) \\ &= \mathbf{w}^{(r)} + 2\alpha_r \mathbf{q} - 2\alpha_r \mathbf{Q} \mathbf{w}^{(r)} \\ &\quad \text{with } \mathbf{q} := (1/m) \sum_{r=1}^m \mathbf{x}^{(r)} y^{(r)}, \mathbf{Q} := (1/m) \sum_{r=1}^m \mathbf{x}^{(r)} (\mathbf{x}^{(r)})^T. \end{aligned} \quad (33)$$

The gradient step (33) only involves matrix/vector multiplications and additions which can be implemented rather efficiently via numerical linear algebra

soft- and hardware [51, 52].

To turn the gradient step (33) into a practical ML methods, we need to

- choose a suitable value for the step-size or learning rate α in (33). The learning rate must not be too large to avoid overshooting or moving away from the optimum (see Figure 16-(b)). On the other hand, if the learning rate is chosen too small, the gradient step makes too little progress towards the solutions of (31) (see Figure 16-(a)). Note that in practice we can only afford to repeat the gradient step for a finite number of iterations.
- define a stopping criterion to decide when to stop iterating (33).

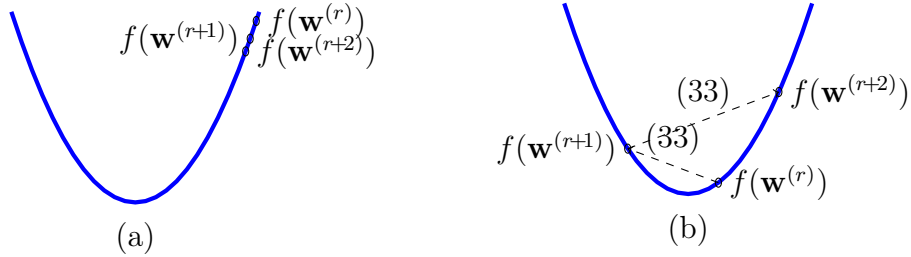


Figure 16: Effect of choosing bad values for the learning rate α in the GD step (33). (a) If the learning rate α in the GD step (33) is chosen too small, the iterations make very little progress towards the optimum or even fail to reach the optimum at all. (b) If the learning rate α is chosen too large the gradient step might “overshoot” such that the iterates $\mathbf{w}^{(r)}$ might diverge from the optimum (it might then happen that $f(\mathbf{w}^{(r+1)}) > f(\mathbf{w}^{(r)})$!).

One approach to choose the learning rate is to start with some initial value (first guess) and monitor the decrease of the objective function. If this decrease

does not agree with the decrease predicted by the gradient, we decrease the learning rate by a constant factor. After we decrease the learning rate we re-consider the decrease of the objective function. We repeat this procedure until a sufficient decrease of the objective function is achieved [36, Sec 6.1].

Another approach to setting the learning rate in (33) is to use a prescribed sequence α_r of learning rates that vary across successive gradient steps. The sequence α_r , for $r = 1, 2, \dots$, is also referred to as a learning rate schedule [53]. In particular, convergence of (33) is guaranteed for any learning rate schedule that satisfies [36, Sec. 6.1]

$$\lim_{r \rightarrow \infty} \alpha_r = 0, \sum_{r=1}^{\infty} \alpha_r = \infty, \text{ and } \sum_{r=1}^{\infty} \alpha_r^2 < \infty. \quad (34)$$

The first condition (34) requires that the learning rate eventually become sufficiently small to avoid overshooting. The third condition (34) requires that the learning rate decays sufficiently fast. Note that the first and third condition in (34) could be satisfied by the trivial learning rate schedule $\alpha_r = 0$ which is clearly not useful as the gradient step has no effect. This trivial schedule is prohibited by the middle condition of (34) which requires that the learning rate α_r is large enough such that the gradient step makes sufficient progress.

For the stopping criterion we might use a fixed number r_{\max} of iterations. This hyper-parameter r_{\max} might be dictated by limited resources (such as computational time) for implementing GD or tuned via validation techniques. Another choice for the stopping criterion could be obtained from monitoring the decrease in the objective value $f(\mathbf{w}^{(r)})$, i.e., we stop repeating the GD step (33) whenever $|f(\mathbf{w}^{(r)}) - f(\mathbf{w}^{(r+1)})| \leq \eta$ for a given tolerance η (which

then becomes a hyper-parameter of the resulting ML method).

The above techniques for choosing the learning rate α and deciding when to stop the GD steps might be useful in practice. However, they might result in sub-optimal use of computational resources for situations where we have more detailed knowledge about the objective function $f(\mathbf{w})$ in (31) and its gradient (see (33)). Indeed, the choice for the learning rate α and the stopping criterion can be guided by the eigenvalues of the matrix (see (33))

$$\mathbf{Q} = (1/m) \sum_{r=1}^m \mathbf{x}^{(r)} (\mathbf{x}^{(r)})^T = (1/m) \mathbf{X} \mathbf{X}^T. \quad (35)$$

Here we used the feature matrix $\mathbf{X} := (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}) \in \mathbb{R}^{d \times m}$. It can be shown that the matrix \mathbf{Q} is psd which implies that all of its eigenvalues are real-valued and non-negative,

$$\lambda_1(\mathbf{Q}) \leq \dots \leq \lambda_d(\mathbf{Q}). \quad (36)$$

It turns out that a GD step (33) reduces the distance $\|\mathbf{w}^{(r)} - \hat{\mathbf{w}}\|_2$ to $\hat{\mathbf{w}}$ (the solution of (31)) by a constant factor [36, Ch. 6],

$$\|\mathbf{w}^{(r+1)} - \hat{\mathbf{w}}\|_2 \leq \kappa^{(\alpha_r)}(\mathbf{Q}) \|\mathbf{w}^{(r)} - \hat{\mathbf{w}}\|_2. \quad (37)$$

Here, we used the contraction factor

$$\kappa^{(\alpha)}(\mathbf{Q}) := \max\{|1 - \alpha 2\lambda_1|, |1 - \alpha 2\lambda_d|\}, \quad (38)$$

which depends on the learning rate and the matrix \mathbf{Q} (35) (via its eigenvalues (36)). The matrix \mathbf{Q} is determined by the feature vectors of the training set in (31), over which we have little (or no) control. However, we can ensure $\kappa^{(\alpha)}(\mathbf{Q}) < 1$ if we choose the learning rate in (33) such that $\alpha_r < 1/\lambda_d$.

For a contraction factor $\kappa^{(\alpha)}(\mathbf{Q}) < 1$, a sufficient condition on number $r^{(\eta)}$ of gradient steps required to ensure an optimization error $\|\mathbf{w}^{(r+1)} - \hat{\mathbf{w}}\|_2 \leq \eta$ is (see (37))

$$r^{(\eta)} \geq \log(\|\mathbf{w}^{(0)} - \hat{\mathbf{w}}\|_2 / \eta) / \log(1/\kappa^{(\alpha)}(\mathbf{Q})). \quad (39)$$

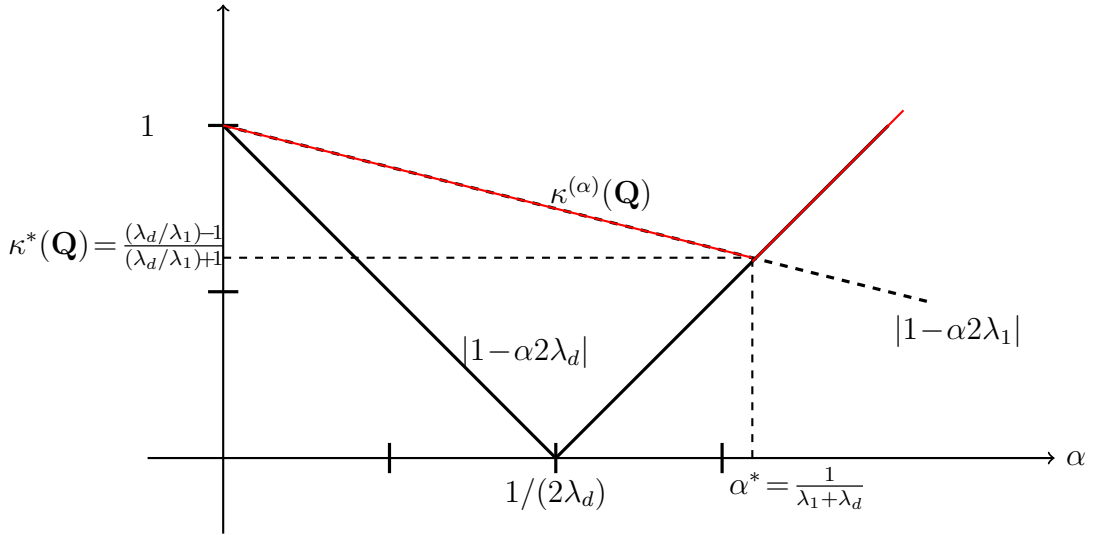


Figure 17: For a given matrix \mathbf{Q} (35), the contraction factor $\kappa^{(\alpha)}(\mathbf{Q})$ varies with the learning rate α . The matrix \mathbf{Q} is determined by the feature vectors of the data points in the ERM (31).

The relation (37) suggests that smaller values of the contraction factor $\kappa^{(\alpha)}(\mathbf{Q})$ result in faster convergence of gradient steps towards a solution of

(31). For a given matrix \mathbf{Q} (35), we can minimize this contraction factor using the learning rate

$$\alpha^{(*)} := \frac{1}{\lambda_1 + \lambda_d}. \quad (40)$$

Inserting the optimal learning rate (40) into (37),

$$\|\mathbf{w}^{(r+1)} - \widehat{\mathbf{w}}\|_2 \leq \underbrace{\frac{(\lambda_d/\lambda_1) - 1}{(\lambda_d/\lambda_1) + 1}}_{:=\kappa^*(\mathbf{Q})} \|\mathbf{w}^{(r)} - \widehat{\mathbf{w}}\|_2. \quad (41)$$

Note that the formula (41) is only valid (meaningful) if the feature vectors $\mathbf{x}^{(r)}$ are such that $\lambda_1 > 0$, i.e., the matrix \mathbf{Q} (35) is invertible. If we cannot guarantee that $\lambda_1 > 0$, we can ensure converge of gradient steps by using a fixed learning rate $\alpha_r = \alpha$ that satisfies [54, Thm. 2.1.14]

$$0 < \alpha < 1/\lambda_d. \quad (42)$$

Note that both, the optimal learning rate (40) and the optimal contraction factor

$$\kappa^*(\mathbf{Q}) := \frac{(\lambda_d/\lambda_1) - 1}{(\lambda_d/\lambda_1) + 1} \quad (43)$$

depend on the eigenvalues of the matrix \mathbf{Q} (35). This matrix, in turn, is determined by the feature vectors $\mathbf{x}^{(r)}$ of the data points used in the ERM (18).

In general we have only little (or no) control over the statistical properties of the features $\mathbf{x}^{(r)}$. However, we might be able to control the eigenvalues (36) of the matrix (35) via using a feature map as a pre-processing step. Two important example for such a feature map are the centering (removing sample mean) and normalization (enforcing unit norm) of features [1, Ch. 5]. Such a pre-processing of features might allow to ensure that the eigenvalues are

contained between lower and upper bounds L and U ,

$$0 < L \leq \lambda_1 \leq \lambda_d \leq U. \quad (44)$$

We can then replace the precise eigenvalues λ_1 and λ_d in (40) and (43) with the bounds L and U , respectively. Moreover, the sufficient condition (42) becomes

$$0 < \alpha < 1/U. \quad (45)$$

According to (41), the ideal case is when all eigenvalues are identical which leads, in turn, to a contraction factor $\kappa^*(\mathbf{Q}) = 0$ (a single gradient step arrives at a solution of (31)!). In this regard, we should use a feature map that results in a matrix \mathbf{Q} (obtained from (35) using the transformed feature vectors) being a scaled identity matrix.

Note that the matrix (35) can be interpreted as an approximation for the covariance matrix of a probability distribution from which the features are drawn i.i.d. (see 3.1). Some authors refer to a RV with covariance matrix being a scaled identity as “white” and, in turn, a feature map that aims at making the covariance matrix a scaled identity as a “whitening” or “decorrelation” transformation [55].

So far we have discussed the properties of GD for solving the ERM (31) arising in linear regression. Let us now study how the behaviour of GD changes by adding a regularization term to the average loss in (31). Thus, we use GD to learn the parameters of a linear hypothesis by solving ridge regression

$$\min_{\mathbf{w} \in \mathbb{R}^d} f^{(\lambda)}(\mathbf{w}) := (1/m) \sum_{r=1}^m \left(y^{(r)} - \mathbf{w}^T \mathbf{x}^{(r)} \right)^2 + \lambda \|\mathbf{w}\|_2^2. \quad (46)$$

The gradient of the objective function $f^{(\lambda)}(\mathbf{w})$ can be computed via

$$\nabla f^{(\lambda)}(\mathbf{w}) = -(2/m) \sum_{r=1}^m \mathbf{x}^{(r)} (y^{(r)} - \mathbf{w}^T \mathbf{x}^{(r)}) + 2\lambda \mathbf{w}. \quad (47)$$

Note that (47) and (46) reduce to (32) and (31), respectively, for the specific choice $\lambda = 0$. The GD step for minimizing $f^{(\lambda)}$ becomes

$$\begin{aligned} \mathbf{w}^{(r+1)} &= \mathbf{w}^{(r)} + 2\alpha_r \mathbf{q} - 2\alpha_r \mathbf{Q}^{(\lambda)} \mathbf{w}^{(r)} \text{ with } \mathbf{q} := (1/m) \sum_{r=1}^m \mathbf{x}^{(r)} y^{(r)}, \\ \mathbf{Q}^{(\lambda)} &:= \lambda \mathbf{I} + (1/m) \mathbf{X} \mathbf{X}^T. \end{aligned} \quad (48)$$

The GD step (48) for ridge regression has exactly the same form as the GD step (33) for linear regression. The convergence analysis (37) also applies to (48). In particular, the convergence speed of (48) is crucially influenced by the eigenvalues

$$0 \leq \tilde{\lambda}_1 \leq \dots \leq \tilde{\lambda}_d \quad (49)$$

of the matrix $\mathbf{Q}^{(\lambda)}$. Basic results from linear algebra tell us that the eigenvalues (49) of the matrix $\mathbf{Q}^{(\lambda)} = \mathbf{Q} + \lambda \mathbf{I}$ and the eigenvalues (36) of the matrix $\mathbf{Q} = (1/m) \mathbf{X} \mathbf{X}^T$ are related via $\tilde{\lambda}_j = \lambda_j + \lambda$. This implies, in turn, that $\kappa^{(*)}(\mathbf{Q}^{(\lambda)}) \leq \kappa^{(*)}(\mathbf{Q})$ and $\lim_{\lambda \rightarrow \infty} \kappa^{(*)}(\mathbf{Q}^{(\lambda)}) = 0$.

Thus, the stronger the regularization in (46), the quicker the gradient steps converge to a solution (46). However, this beneficial computational aspects of using a large λ might result in a poor hypothesis $\hat{h}(\mathbf{x}) = (\hat{\mathbf{w}}^{(\lambda)})^T \mathbf{x}$, with $\hat{\mathbf{w}}^{(\lambda)}$ denoting the solution of (46). Indeed, using excessive regularization (too large λ) might incur a large bias (see Section 3.3 and [1, Ch. 6])

5.2 Projected Gradient Descent

Section 5.1 discussed methods for learning the parameter vector \mathbf{w} of a linear hypothesis by ERM (31). For some applications it might be useful to constrain the parameter vector \mathbf{w} to belong to a subset $\mathcal{C} \subset \mathbb{R}^d$. Indeed, restricting the parameter vectors to \mathcal{C} is a form of model pruning which can help to avoid overfitting (see Section 4).

Another use case for constraining the model parameter are distributed FL methods. In particular, Section 9.3 discusses a FL method for coupling the training of local model parameters $\mathbf{w}^{(i)}$, for $i \in \{1, \dots, n\}$ such that they eventually coincide. This coupling can be implemented by requiring the stacked parameter vector $\mathbf{w} = (\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(n)})^T$ to belong to the subset

$$\mathcal{C} = \left\{ (\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(n)})^T : \mathbf{w}^{(1)} = \dots = \mathbf{w}^{(n)} \right\}.$$

Let us now show how to adapt the basic GD step (33) to solve the ERM

$$f^* = \min_{\mathbf{w} \in \mathcal{C}} f(\mathbf{w}) := (1/m) \sum_{r=1}^m (y^{(r)} - \mathbf{w}^T \mathbf{x}^{(r)})^2. \quad (50)$$

We assume that the constraint set $\mathcal{C} \subseteq \mathbb{R}^d$ is such that we can efficiently compute the projection

$$P_{\mathcal{C}}(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}' \in \mathcal{C}} \|\mathbf{w} - \mathbf{w}'\|_2 \text{ for any } \mathbf{w} \in \mathbb{R}^d. \quad (51)$$

A suitable modification of the gradient step (33) to solve the constrained

ERM (50) is [36]

$$\begin{aligned}
\mathbf{w}^{(r+1)} &:= P_{\mathcal{C}}(\mathbf{w}^{(r)} - \alpha \nabla f(\mathbf{w}^{(r)})) \\
&\stackrel{(32)}{=} P_{\mathcal{C}}(\mathbf{w}^{(r)} + (2\alpha/m) \sum_{r=1}^m \mathbf{x}^{(r)} (y^{(r)} - (\mathbf{w}^{(r)})^T \mathbf{x}^{(r)})) \\
&= P_{\mathcal{C}}(\mathbf{w}^{(r)} + \mathbf{q} - \mathbf{Q} \mathbf{w}^{(r)}) \text{ with } \mathbf{q} := (2\alpha/m) \sum_{r=1}^m \mathbf{x}^{(r)} y^{(r)}, \\
&\quad \mathbf{Q} := (2\alpha/m) \sum_{r=1}^m \mathbf{x}^{(r)} (\mathbf{x}^{(r)})^T. \tag{52}
\end{aligned}$$

The projected GD step (52) amounts to first computing an ordinary GD step $\mathbf{w}^{(r)} \mapsto \mathbf{w}^{(r)} - \alpha \nabla f(\mathbf{w}^{(r)})$ and then projecting back to the constraint set \mathcal{C} . Note that we re-obtain the basic GD step (33) from the projected GD step (52) for the extreme case where $\mathcal{C} = \mathbb{R}^d$.

The approaches for choosing the learning rate α in - and defining a stopping criterion for - the basic GD step (33) explained in Section 5.1 can also be used for the projected GD step (52). Moreover, somewhat surprisingly, the convergence speed of the projected GD is characterized by the very same bound (37) as we derived for the unconstrained GD in Section 5.1 [36, Ch. 6]. However, the bound (37) is only telling about the number of GD steps required to achieve a guaranteed level of sub-optimality $|f(\mathbf{w}^{(r)}) - f^*|$. Each projected GD step (52) might require significantly more computation than the basic GD step, as it requires to compute the projection (51).

5.3 Perturbed Gradient Descent

Consider the gradient step (33) for solving the ERM (31) arising in linear regression. In what follows we assume that the feature vectors $\mathbf{x}^{(r)}$ used in

the ERM (31) are such that the matrix \mathbf{Q} (see (35)) is invertible. Thus, the smallest eigenvalue of the matrix \mathbf{Q} is positive: $\lambda_1 > 0$. Moreover, we use a fixed learning rate $\alpha_r = \alpha$ in (33) which is sufficiently small such that the contraction factor (38) satisfies $\kappa^{(\alpha)}(\mathbf{Q}) < 1$.

In practice, we might not be able to compute all the quantities in (33) exactly. For example, the sum $\sum_{r=1}^m \mathbf{x}^{(r)}(y^{(r)} - (\mathbf{w}^{(r)})^T \mathbf{x}^{(r)})$ might require to gather data points from distributed storage locations. These storage locations might become unavailable during the computation (33) due to hardware or software failures or limited connectivity.

To take into account different sources of imperfections during the computation of (33), we introduce an error term,

$$\begin{aligned}
\mathbf{w}^{(r+1)} &:= \mathbf{w}^{(r)} - \alpha \nabla f(\mathbf{w}^{(r)}) + \boldsymbol{\varepsilon}^{(r)} \\
&\stackrel{(32)}{=} \mathbf{w}^{(r)} + (2\alpha/m) \sum_{r=1}^m \mathbf{x}^{(r)}(y^{(r)} - (\mathbf{w}^{(r)})^T \mathbf{x}^{(r)}) + \boldsymbol{\varepsilon}^{(r)} \\
&= \mathbf{w}^{(r)} + 2\alpha \mathbf{q} - 2\alpha \mathbf{Q} \mathbf{w}^{(r)} + \boldsymbol{\varepsilon}^{(r)} \text{ with } \mathbf{q} := (1/m) \sum_{r=1}^m \mathbf{x}^{(r)} y^{(r)}, \\
&\quad \mathbf{Q} := (1/m) \sum_{r=1}^m \mathbf{x}^{(r)} (\mathbf{x}^{(r)})^T. \tag{53}
\end{aligned}$$

A straight-forward analysis of the perturbed gradient step (53) yields the following bound on the deviation between $\mathbf{w}^{(r)}$ and the solution $\widehat{\mathbf{w}}$ of (31) (see, e.g., [56])

$$\|\mathbf{w}^{(r)} - \widehat{\mathbf{w}}\|_2 \leq (\kappa^{(\alpha)}(\mathbf{Q}))^r \|\mathbf{w}^{(0)} - \widehat{\mathbf{w}}\|_2 + \sum_{r'=1}^r (\kappa^{(\alpha)}(\mathbf{Q}))^{r'} \|\boldsymbol{\varepsilon}^{(r-r')}\|_2. \tag{54}$$

This bound applies for any number of iterations $r = 1, 2, \dots$ of the perturbed gradient step (53).

5.4 Stochastic Gradient Descent

Section 5.1 explained how to use gradient steps to construct a sequence of model parameters that are (hopefully) increasingly accurate approximations to a solution of ERM (31) for linear regression. Each gradient step (33) for linear regression requires to compute the sum

$$- (2/m) \sum_{r=1}^m \mathbf{x}^{(r)} \left(y^{(r)} - (\mathbf{w}^{(r)})^T \mathbf{x}^{(r)} \right) \stackrel{(32)}{=} \nabla f(\mathbf{w}^{(r)}). \quad (55)$$

The sum (55) depends on the data points used in the ERM (31) via their feature vectors $\mathbf{x}^{(r)}$ and labels $y^{(r)}$, for $r = 1, \dots, m$. Moreover, the sum (55) also depends on the current choice $\mathbf{w}^{(r)}$ for the model parameters (which we hope to improve by the gradient step). Note that we need to compute a sum of the form (55) for each single gradient step (33).

Computing the sum in (55) might be computationally challenging (or even impossible) for several reasons:

- **Big Data.** Computing the sum (55) might become infeasible for very large training sets, with m in the order of billions. We might then simply not have enough time to compute the sum (55) within the available time budget for completing a gradient step. For example, if we use GD methods to learn optimal steering of a self-driving car, we might need to execute the gradient steps in real-time.
- **Infrastructure Cost.** Second, for decentralized data which is stored in different data centres located all over the planet, the sum (55) would require an excessive amount of network resources.
- **Finite Bandwidths.** Moreover, the finite transmission rate of com-

munication networks limits the speed with which the gradient steps (33) can be executed in real-time (e.g., while steering an autonomous vehicle).

The idea of SGD (and its variants) is to replace the exact gradient $\nabla f(\mathbf{w})$ (55) in (33) by an approximation $g(\mathbf{w}) \approx \nabla f(\mathbf{w})$,

$$\mathbf{w}^{(r+1)} = \mathbf{w}^{(r)} - \alpha_r g(\mathbf{w}^{(r)}). \quad (56)$$

We hope that the approximation $g(\mathbf{w}^{(r)})$ can be computed quicker than the exact gradient (55) while being sufficiently accurate.

The word “stochastic” in the name SGD hints already at the use of probabilistic models in the construction of the gradient approximation $g(\mathbf{w}) \approx \nabla f(\mathbf{w})$. The most basic form of SGD constructs the approximation $g(\mathbf{w})$ by replacing the sum in (55) with a randomly selected component,

$$g(\mathbf{w}) := -2\mathbf{x}^{(\hat{l})} \left(y^{(\hat{l})} - \mathbf{w}^T \mathbf{x}^{(\hat{l})} \right). \quad (57)$$

The index \hat{l} is chosen randomly from the set $\{1, \dots, m\}$. The resulting SGD step (or iteration) updates given model parameters $\mathbf{w}^{(r)}$ to obtain a new (hopefully improved) model parameters

$$\mathbf{w}^{(r+1)} = \mathbf{w}^{(r)} - \alpha 2\mathbf{x}^{(\hat{l}_r)} \left(y^{(\hat{l}_r)} - (\mathbf{w}^{(r)})^T \mathbf{x}^{(\hat{l}_r)} \right). \quad (58)$$

Every update (58) uses a “fresh” (randomly chosen or drawn) index \hat{l}_r . The indices \hat{l}_r , for $r = 1, 2, \dots$ are realizations of i.i.d. RVs whose common probability distribution is the uniform distribution over the index set $\{1, \dots, m\}$,

$$p(\hat{l}_r = r) = 1/m, \text{ for each } r \in \{1, \dots, m\}.$$

The update (58) avoids the computation of the sum (55) but instead only requires the computation of a random index \hat{l}_r . The resulting computational savings might be significant when the training set consists of a large number of data points.

The saving in computational complexity of SGD comes at the cost of introducing the approximation error (or “gradient noise”)

$$\Delta g^{(r)} := \nabla f(\mathbf{w}^{(r)}) - g(\mathbf{w}^{(r)}). \quad (59)$$

Indeed, we can interpret the basic SGD update (56) as a perturbed gradient step (53) with perturbation given by the scaled gradient noise (59),

$$\mathbf{w}^{(r+1)} := \mathbf{w}^{(r)} - \alpha_r \nabla f(\mathbf{w}^{(r)}) + \boldsymbol{\varepsilon}^{(r)} \text{ with } \boldsymbol{\varepsilon}^{(r)} = \alpha_r \cdot \Delta g^{(r)}. \quad (60)$$

Note that the perturbation in the SGD update (60) is the product of learning rate α_r and the gradient noise (59). Thus, it seems that we should use a very small learning rate α_r to attenuate the gradient noise. However, we cannot make the learning rate arbitrarily small (the extreme case would be $\alpha_r = 0$) since the SGD updates (56) would then make too little progress towards the solution of (31) (see Figure 16-(a)). One useful choice for the learning rate of SGD is $\alpha_r = 1/r$ [57].

Mini-Batch SGD. Different variants of SGD have been developed with the aim to reduce the gradient noise $\Delta g^{(r)}$ in the SGD step (60). One such variant is referred to as mini-batch SGD. The idea is to use more than one (randomly chosen) data point for the construction of a gradient approximation,

$$g(\mathbf{w}) = -(2/B) \sum_{r' \in \mathcal{B}} \mathbf{x}^{(r')} (y^{(r')} - (\mathbf{w})^T \mathbf{x}^{(r')}) \quad (61)$$

For each new iteration of SGD (56), mini-batch SGD generates a new batch

$$\mathcal{B} = \left\{ \left(\mathbf{x}^{(\hat{l}_1)}, y^{(\hat{l}_1)} \right), \dots, \left(\mathbf{x}^{(\hat{l}_B)}, y^{(\hat{l}_B)} \right) \right\}$$

of B randomly selected data points is generated.

Online Learning. We can use SGD (56) also in ML applications where data points $(\mathbf{x}^{(t)}, y^{(t)})$ are gathered sequentially over time $t = 1, 2, \dots$. For example, data points could be the snapshots of the on-board camera in a self-driving car. During each time instant t , the camera delivers a new data point $(\mathbf{x}^{(t)}, y^{(t)})$ that results in a new component $f^{(t)}(\mathbf{w}) := (y^{(t)} - \mathbf{w}^T \mathbf{x}^{(t)})^2$ for the objective function in (31).

We can use the SGD step (56) to obtain an online learning method (see [1, Sec. 4.7]). The idea is to use the gradient of the new component $f^{(t)}(\mathbf{w}) = (y^{(t)} - \mathbf{w}^T \mathbf{x}^{(t)})^2$ as an estimate for the gradient of the objective function (31). This online version of SGD computes, at each time instant t ,

$$\mathbf{w}^{(t+1)} := \mathbf{w}^{(t)} + 2\alpha_t \mathbf{x}^{(t)} \left(y^{(t)} - (\mathbf{w}^{(t)})^T \mathbf{x}^{(t)} \right). \quad (62)$$

We can use the online gradient step (62) also to learn model parameters from a given training set \mathcal{D} under an i.i.d. assumption. Indeed, we can identify the time index t with the sample index r that enumerates the data points in \mathcal{D} . Since the data point $(\mathbf{x}^{(t)}, y^{(t)}) \in \mathcal{D}$ is a realization of a RV with probability distribution $p(\mathbf{x}, y)$, the gradient $\nabla f^{(t)}(\mathbf{w}) = -\mathbf{x}^{(t)}(y^{(t)} - \mathbf{w}^T \mathbf{x}^{(t)})$ is an unbiased estimator for the gradient $\nabla \bar{L}(\mathbf{w})$ of the risk $\bar{L}(\mathbf{w}) = \mathbb{E} \left\{ (y - \mathbf{w}^T \mathbf{x})^2 \right\}$.

5.5 Exercises

Exercise 5.1. Positive Semidefinite (10 points). Develop a rigorous proof for the fact that the matrix \mathbf{Q} defined in (33) is psd.

Exercise 5.2. Optimality Condition (20 points). Show that any solution $\hat{\mathbf{w}}$ of ERM (31) must satisfy

$$\mathbf{Q}\hat{\mathbf{w}} = \mathbf{q}, \text{ with vector } \mathbf{q} \text{ and matrix } \mathbf{Q} \text{ defined in (33).}$$

Exercise 5.3. GD on a Parabola (5 points). Consider GD $w^{(r+1)} = w^{(r)} - \alpha \nabla f(w^{(r)})$ with a fixed learning rate α applied to the objective function $f(w) = (3/2)w^2 + 4w - 10^{10}$. Characterize the set of values for the learning rate α such that $w^{(r)}$ converges to an optimum $\bar{w} \in \operatorname{argmin}_{w \in \mathbb{R}} f(w)$.

Exercise 5.4. Matrix Calculus (5 points). Provide a rigorous proof for the second equality in (35).

Exercise 5.5. Learning rate Schedule (10 points). Find an example for a learning rate schedule α_r that satisfies all three conditions in (34).

Exercise 5.6. Projecting onto Ray (5 points). Try to develop a closed-form expression for the projection (51) of a vector $\mathbf{w} \in \mathbb{R}^3$ onto a constraint set being a ray,

$$\mathcal{C} = \{c(1, 1, 1)^T : c \in \mathbb{R}_+\}.$$

Exercise 5.7. Projecting onto Interval (5 points). Try to develop a closed-form expression for the projection (51) of a real number $w \in \mathbb{R}$ (which can be identified with a vector \mathbf{w} of length one) onto the constraint set $\mathcal{C} = [-4, 1]$ which is a closed interval.

Exercise 5.8. Projecting onto Ball (20 points). Try to develop a closed-form expression for the projection (51) of a vector $\mathbf{w} \in \mathbb{R}^d$ onto the constraint set $\mathcal{C} = \{\mathbf{w}' \in \mathbb{R}^d : \|\mathbf{w}'\|_2 \leq 3\}$.

Exercise 5.9. Unbiased (20 points). Due to the randomly chosen index \hat{l} , the gradient approximation (57) in SGD becomes (the realization of) a RV. This implies, in turn, that the gradient noise (59) is also (the realization of) a RV. What is the mean vector of this RV ?

Exercise 5.10. Local Quadratic Approximation (10 points). Consider a differentiable function $f(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ and some vector $\mathbf{w}^{(r)} \in \mathbb{R}^d$. The gradient step

$$\mathbf{w}^{(r+1)} := \mathbf{w}^{(r)} - \alpha_r \nabla f(\mathbf{w}^{(r)}). \quad (63)$$

computes the updated vector $\mathbf{w}^{(r+1)}$ that hopefully satisfies $f(\mathbf{w}^{(r+1)}) < f(\mathbf{w}^{(r)})$. Show that the updated vector delivered by the gradient step (63) minimizes the quadratic function

$$\tilde{f}^{(r)}(\mathbf{w}) := f(\mathbf{w}^{(r)}) + \mathbf{w}^T \nabla f(\mathbf{w}^{(r)}) + (1/(2\alpha_r)) \|\mathbf{w} - \mathbf{w}^{(r)}\|_2^2.$$

Part II

Federated Learning

6 Networked Data and Models

Many important application domains such as healthcare generate collections of local datasets. Indeed, healthcare providers (hospitals, blood labs, etc.) maintain their own local patient databases. These local datasets are related via statistical similarities, e.g., the frequency of flu patients at nearby hospitals might be correlated. Section 6.1 introduces the empirical graph as a useful representation of collections of local datasets along with their similarities.

Section 6.2 augments the empirical graph by assigning a separate local hypothesis space (or local model) to each node. We could train the local models for each node separately using the corresponding local datasets as training set in ERM (see Section 3.2). However, the local dataset might be too small to train a (high-dimensional) model such as a large ANN or a linear model with many features (see Section 3.3). We use the network structure of the empirical graph to pool local datasets to obtain a sufficiently large training set to train each local model via ERM.

This pooling can be motivated by a clustering assumption as discussed in Section 6.3. As an alternative to directly pool (or collect) local datasets at well-connected nodes, Section 7 puts forward regularization techniques to enforce similar local model parameters at well-connected nodes.

The empirical graph might be designed manually by using domain expertise that provides measures for the similarity between local datasets. However, for some applications it might be useful to learn the empirical graph in a data-driven fashion (see Section 6.4).

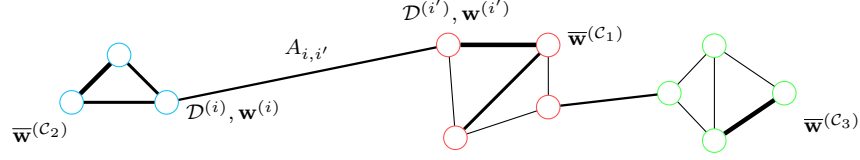


Figure 18: Example of an empirical graph whose nodes $i \in \mathcal{V}$ carry local datasets $\mathcal{D}^{(i)}$ and local models (and their model parameters) $\mathbf{w}^{(i)}$. In this example, the empirical graph is partitioned into three disjoint clusters $\mathcal{C}_1, \mathcal{C}_2, \mathcal{C}_3$, with cluster-wise optimal model parameters $\bar{\mathbf{w}}^{(\mathcal{C})}$ (see Section 6.3). FL methods typically do not know the cluster structure but only pair-wise similarities as edge weights $A_{i,i'}$ (indicated by thickness).

6.1 The Empirical Graph

An empirical graph is an undirected weighted graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ whose nodes $\mathcal{V} := \{1, \dots, n\}$ represent local datasets $\mathcal{D}^{(i)}$, for $i \in \mathcal{V}$. Each node $i \in \mathcal{V}$ of the empirical graph \mathcal{G} carries a separate local dataset $\mathcal{D}^{(i)}$. To build intuition, think of a local dataset $\mathcal{D}^{(i)}$ as a labelled dataset

$$\mathcal{D}^{(i)} := \left\{ (\mathbf{x}^{(i,1)}, y^{(i,1)}), \dots, (\mathbf{x}^{(i,m_i)}, y^{(i,m_i)}) \right\}. \quad (64)$$

Here, $\mathbf{x}^{(i,r)}$ and $y^{(i,r)}$ denote, respectively, the features and the label of the r th data point in the local dataset $\mathcal{D}^{(i)}$. Note that the size m_i of the local dataset might vary between different nodes $i \in \mathcal{V}$. Figure 18 depicts an example for an empirical graph.

An undirected edge $\{i, i'\} \in \mathcal{E}$ in the empirical graph indicates that the local datasets $\mathcal{D}^{(i)}$ and $\mathcal{D}^{(i')}$ have similar statistical properties. We quantify the level of similarity by a positive edge weight $A_{i,i'} > 0$. The neighbourhood

of a node $i \in \mathcal{V}$ is $\mathcal{N}^{(i)} := \{i' \in \mathcal{V} : \{i, i'\} \in \mathcal{E}\}$.

Our (rather) informal notion of similarity could be made precise via a probabilistic model that interprets the local dataset $\mathcal{D}^{(i)}$ as realizations of RVs with some parametrized probability distribution $p^{(i)}(\mathcal{D}^{(i)}; \mathbf{w}^{(i)})$. The similarity between local datasets $\mathcal{D}^{(i)}$ and $\mathcal{D}^{(i')}$ could then be defined in terms of distance $\|\mathbf{w}^{(i)} - \mathbf{w}^{(i')}\|_2$ between the parameters of the probability distributions.

Note that the undirected edges $\{i, i'\}$ of an empirical graph encode a symmetric notion of similarity between local datasets. If the local dataset $\mathcal{D}^{(i)}$ at node i is similar to the local dataset $\mathcal{D}^{(i')}$ at node i' , then also the local dataset $\mathcal{D}^{(i')}$ is similar to the local dataset $\mathcal{D}^{(i)}$.

Despite the symmetric notion of similarity represented by the edges \mathcal{E} of an empirical graph, it will be convenient for the formulation and analysis of FL algorithms to orient the edges in \mathcal{E} . In particular, we define the head and tail of an undirected edge $e = \{i, i'\}$ as $e_+ := \min\{i, i'\}$ and $e_- := \max\{i, i'\}$, respectively. The directed edge of an empirical graph are

$$\{(i, i') : i, i' \in \mathcal{V}, i < i' \text{ and } \{i, i'\} \in \mathcal{E}\}. \quad (65)$$

We abuse notation and use \mathcal{E} to denote both, the set of undirected edges and its oriented version (65).

The empirical graph of networked data is a design choice which is guided by computational aspects and statistical aspects of the resulting ML method. For example, using an empirical graph with a small number of edges (“sparse graphs”) typically results in a smaller computational complexity. Indeed, the amount of computation required by the FL methods developed in Section 9 is proportional to the number of edges in the empirical graph. On the other

hand, the empirical graph should contain sufficient number of edges between nodes that carry statistically similar local datasets. This allows GTVMin techniques to adaptively pool local datasets into clusters of (approximately) homogeneous data (see Section 6.3).

6.2 Networked Models

Consider data with empirical graph \mathcal{G} whose nodes $i \in \mathcal{V}$ carry local datasets $\mathcal{D}^{(i)}$. For each node $i \in \mathcal{V}$, we wish to learn a useful hypothesis $\hat{h}^{(i)}$ from a local hypothesis space $\mathcal{H}^{(i)}$. The learnt hypothesis should incur a small average loss over a local dataset $\mathcal{D}^{(i)}$.

The collection of local models $\mathcal{H}^{(i)}$, for each node $i \in \mathcal{V}$, constitutes a networked model $\mathcal{H}^{(\mathcal{G})}$ over an empirical graph \mathcal{G} . We also allow for heterogeneous networked models where different nodes carry different local models $\mathcal{H}^{(i)}$. For example, $\mathcal{H}^{(i)}$ might be a linear model $\mathcal{H}^{(d)}$, while $\mathcal{H}^{(i')}$ might be a decision tree for some other node $i' \neq i$.

We measure the usefulness of a particular hypothesis $h \in \mathcal{H}^{(i)}$ using a local loss function

$$L_i(\cdot) : \mathcal{H}^{(i)} \rightarrow \mathbb{R}_+ : h \mapsto L_i(h). \quad (66)$$

If node $i \in \mathcal{V}$ carries a local dataset of the form (64), we might define a local loss function via the average loss

$$L_i(h) := (1/m_i) \sum_{r=1}^{m_i} L((\mathbf{x}^{(i,r)}, y^{(i,r)}), h). \quad (67)$$

However, the FL methods developed in Section 7 and Section 9 can be applied also to loss functions that are not of the form (67).

Our focus will be on parametrized networked models where each $\mathcal{H}^{(i)}$ is parametrized by a common finite-dimensional Euclidean space \mathbb{R}^d . This setting covers some widely-used ML models such as (regularized) generalized linear models or linear time series models [58–62]. For a parametrized networked model, the local hypothesis $h^{(i)}$ at any node $i \in \mathcal{V}$ is fully determined by a local parameter vector $\mathbf{w}^{(i)} \in \mathbb{R}^d$. We find it convenient to collect local parameters $\mathbf{w}^{(i)}$ as the values of a map $\mathbf{w} : i \mapsto \mathbf{w}^{(i)}$. These maps constitute the networked model parameter space

$$\mathcal{W} := \left\{ \mathbf{w} : i \mapsto \mathbf{w}^{(i)} \text{ such that } h^{(\mathbf{w}^{(i)})} \in \mathcal{H}^{(i)} \right\}. \quad (68)$$

The usefulness of a specific choice for the local model parameter $\mathbf{w}^{(i)}$ is measured by the local loss function $L_i(\mathbf{w}^{(i)}) := L_i(h^{(\mathbf{w}^{(i)})})$. The local loss functions $L_i(\mathbf{w}^{(i)})$ are an important design choice within FL systems. This design choice is related to the design choice for the local datasets and (the parametrization of) local models $\mathcal{H}^{(i)}$ (see (67)).

As for other ML components (see Section 2), there are statistical aspects and computational aspects guiding the choice for the local loss functions. Statistically, we want the loss function to favour local model parameters that result in a robust and accurate hypothesis for each node. Computationally, we want a local loss function that can be minimized efficiently using the available computational resources (see Section 5 and Section 9).

The FL methods in Section 9.3 and 9.2 require local loss functions that are differentiable and allow for efficient computation of their gradient at arbitrary points. Indeed, these FL methods use GD steps as their elementary computation.

The FL method in Section 9.1 requires that the local loss function $L_i(\mathbf{w}^{(i)})$ allows for efficient solution of the regularized problem (see Section 4.3)

$$\min_{\mathbf{w}^{(i)} \in \mathbb{R}^d} L_i(\mathbf{w}^{(i)}) + \lambda \|\mathbf{w}^{(i)} - \mathbf{w}'\|_2^2 \text{ with some } \mathbf{w}' \in \mathbb{R}^d. \quad (69)$$

Note that the basic computation (69) used by the FL method in Section 9.1 is itself an optimization problem which might be solved by the gradient-based methods in Section 5.

6.3 Clustering Assumption

Many applications generate local datasets which do not carry sufficient statistical power to guide learning of model parameters $\mathbf{w}^{(i)}$ (see Section 3.3). As a case in point, consider a local dataset $\mathcal{D}^{(i)}$ of the form (64), with feature vectors $\mathbf{x}^{(r)} \in \mathbb{R}^d$ with $m_i \ll d$. We would like to learn the parameter vector $\mathbf{w}^{(i)}$ of a linear hypothesis $h(\mathbf{x}) = \mathbf{x}^T \mathbf{w}^{(i)}$.

Linear regression methods [1, Sec. 3.1] learn the parameter vector by minimizing the average squared error loss

$$\min_{\mathbf{w}^{(i)}} L_i(\mathbf{w}^{(i)}) = (1/m_i) \sum_{r=1}^{m_i} (y^{(r)} - (\mathbf{w}^{(i)})^T \mathbf{x}^{(r)})^2.$$

However, for $m_i \ll d$ (the high-dimensional regime) the above minimum is not unique and might also provide a poor hypothesis. In particular, a parameter vector with minimum average loss over $\mathcal{D}^{(i)}$ might incur unacceptably large prediction errors for data points outside $\mathcal{D}^{(i)}$ [1, Ch. 6].

Training linear models in the high-dimensional regime requires regularization techniques such as those used by ridge regression or Lasso [34]. These methods implement regularization by adding a penalty term to the average

loss of a hypothesis incurred over a training set. This penalty term is an approximation or estimate for the increase in average loss when the hypothesis is applied to data points outside the training set.

Section 7 proposes different measures for the variation of local hypotheses $h^{(i)}$, for $i \in \mathcal{V}$, over edges in the empirical graph. The resulting measures for the total variation of a networked hypothesis are then used as a penalty term to regularize the training of local models. We will see that adding this penalty term results in an adaptive pooling of local datasets into clusters for which a common cluster-wise hypothesis is learnt.

Pooling local datasets to form clusters only makes sense if they have similar statistical properties that can be captured by a common hypothesis. In particular, we require the local loss functions at nodes $i \in \mathcal{C}$ in the same cluster \mathcal{C} to have approximately identical minimizers.

Assumption 1 (Clustering (informal)). *Consider local datasets $\mathcal{D}^{(i)}$ represented by the nodes of an empirical graph \mathcal{G} . There is a partition*

$$\mathcal{P} = \{\mathcal{C}_1, \dots, \mathcal{C}_{|\mathcal{P}|}\} \text{ with } \mathcal{C}_c \cap \mathcal{C}_{c'} = \emptyset \text{ and } \mathcal{V} = \mathcal{C}_1 \cup \dots \cup \mathcal{C}_{|\mathcal{P}|}, \quad (70)$$

of the nodes \mathcal{V} into disjoint clusters $\mathcal{C}_c \in \mathcal{P}$. The local datasets at the nodes in the same cluster conform to an i.i.d. assumption (see Section 3.1) with a cluster-specific probability distribution $p^{(c)}(\mathbf{x}, y)$.

Assume we would know which nodes belong to the same cluster \mathcal{C} . It would then be sensible to pool the local datasets, or equivalently add the corresponding local loss functions (67), and learn model parameters by minimizing

the resulting cluster-wise loss function,

$$\bar{\mathbf{w}}^{(\mathcal{C})} = \underset{\mathbf{v} \in \mathbb{R}^d}{\operatorname{argmin}} f^{(\mathcal{C})}(\mathbf{v}) \quad \text{with} \quad f^{(\mathcal{C})}(\mathbf{v}) := \sum_{i \in \mathcal{C}} \underbrace{L_i(\mathbf{v})}_{:= (1/m_i) \sum_{(\mathbf{x}, y) \in \mathcal{D}^{(i)}} L((\mathbf{x}, y), h(\mathbf{v}))}. \quad (71)$$

Note that (71) cannot be implemented in practice since we typically do not know the cluster \mathcal{C} for the i.i.d. assumption is valid (approximately). Instead, we will use regularization to couple the minimizing of individual local loss functions. This regularization will be implemented via measures for the variation of local model parameters across the weighted edges of the empirical graph.

Section 7 defines generalized total variation (GTV) as a measure for the variation of local model parameters. This results then in GTVMin as a flexible design principle for FL methods. These methods use the edges of the empirical graph to couple the training of local models by enforcing small variations of local model parameters across edges. It turns out that if nodes in the same cluster are connected by sufficiently many edges and only few edges connect nodes in different clusters, then GTVMin captures the underlying cluster structure in Assumption 1 and results in local model parameters that approximately solve (71).

6.4 Graph Learning Methods

We have introduced the empirical graph as a representation of collections of local datasets and their similarity structure. It is important to note that the empirical graph is essentially a design choice that is guided by computational aspects and statistical aspects. Computationally, we prefer a sparse empirical

graph with few edges. Indeed, the FL methods presented in Section 9 can be implemented as message passing over the empirical graph which means that the amount of computation is proportional to the number of edges. Statistically, we want the empirical graph to reflect an underlying cluster structure (see Assumption 1). This might require to have a sufficiently large number of edges connecting nodes in the same cluster.

6.4.1 Testing Similarity

An edge between nodes $i, i' \in \mathcal{V}$ should represent a statistical similarity between local datasets $\mathcal{D}^{(i)}$ and $\mathcal{D}^{(i')}$. To make this notion precise we can use a probabilistic model (see Section 3.1) for (the generation of) local datasets. In particular, we interpret the data points in $\mathcal{D}^{(i)}$ as realizations of i.i.d. RVs with common probability distribution $p^{(i)}(\mathbf{x}, y)$. Similarly, we interpret the data points in $\mathcal{D}^{(i')}$ as realizations of i.i.d. RVs with common probability distribution $p^{(i')}(\mathbf{x}, y)$.

We can then define the presence of an edge (and its weight) between nodes $i, i' \in \mathcal{V}$ using various measures for the distance $D(p^{(i)}, p^{(i')})$ between $p^{(i)}(\mathbf{x}, y)$ and $p^{(i')}(\mathbf{x}, y)$. One example for such a distance measure is the Kullback-Leibler divergence which can be defined for any pair of probability distribution under mild technical conditions.

In particular, the KL divergence is between two multivariate normal distributions $p^{(i)}, p^{(i')}$, each having mean $\boldsymbol{\mu}^{(i)}, \boldsymbol{\mu}^{(i')} \in \mathbb{R}^d$ and non-singular

covariance matrix $\mathbf{C}^{(i)}, \mathbf{C}^{(i')}$, respectively, is

$$D(p^{(i)}, p^{(i')}) = (1/2) \left(\text{tr} \{ (\mathbf{C}^{(i')})^{-1} \mathbf{C}^{(i)} \} - d \right. \\ \left. + (\boldsymbol{\mu}^{(i')} - \boldsymbol{\mu}^{(i)})^T (\mathbf{C}^{(i')})^{-1} (\boldsymbol{\mu}^{(i')} - \boldsymbol{\mu}^{(i)}) \right. \\ \left. + \log [\det (\mathbf{C}^{(i')}) / \det (\mathbf{C}^{(i)})] \right). \quad (72)$$

In practice, we typically do not know the underlying probability distribution of data points, including their mean vector and covariance matrix, of a local dataset $\mathcal{D}^{(i)}$. However, we can estimate these parameters using the sample mean and the sample covariance matrix,

$$\hat{\boldsymbol{\mu}}^{(i)} = (1/m) \sum_{r=1}^m \mathbf{x}^{(i,r)} \text{ (sample mean)} \\ \hat{\mathbf{C}}^{(i)} = (1/m) \sum_{r=1}^m (\mathbf{x}^{(i,r)} - \hat{\boldsymbol{\mu}}^{(i)}) (\mathbf{x}^{(i,r)} - \hat{\boldsymbol{\mu}}^{(i)})^T \text{ (sample covariance)}. \quad (73)$$

We can plug in the estimates (73) into (72) to obtain a distance measure $\hat{D}(p^{(i)}, p^{(i')})$ that can be computed directly from local datasets.

We can use a measure $\hat{D}(p^{(i)}, p^{(i')})$ for the dis-similarity between any two local datasets $\mathcal{D}^{(i)}, \mathcal{D}^{(i')}$ to construct an empirical graph (with nodes $\mathcal{V} = \{1, \dots, n\}$) in several ways.

- **thresholding:** Add an edge between nodes $i, i' \in \mathcal{V}$ if $\hat{D}(p^{(i)}, p^{(i')})$ is below a prescribed threshold η ,

$$\{i, i'\} \in \mathcal{E} \text{ if and only if } \hat{D}(p^{(i)}, p^{(i')}) < \eta. \quad (74)$$

- **nearest neighbours:** Fix a number $k \in \{0, \dots, n-1\}$ which roughly determines the maximum number of neighbours for any node $i \in \mathcal{V}$ in

the empirical graph. Then iterate over all nodes $i \in \mathcal{V}$ (in some order), adding k edges between i and those other k nodes $i' \in \mathcal{V} \setminus \{i\}$ with smallest distance $\widehat{D}(p^{(i)}, p^{(i')})$.

- **fully connected weighted empirical graph** Add an edge $\{i, i'\}$ between every pair $i, i' \in \mathcal{V}$ of two distinct nodes. Determine the edge weight $A_{i,i'}$ based on $\widehat{D}(p^{(i)}, p^{(i')})$,

$$A_{i,i'} = g(\widehat{D}(p^{(i)}, p^{(i')})) \text{ with some profile } g(\cdot) : \mathbb{R}_+ \rightarrow \mathbb{R}_+. \quad (75)$$

One popular choice for the profile in (75) is $g(a) := \exp(-a^2)$ [63, Sec. 8], [64]. The threshold rule (74) is obtained from (75) using the profile function $g(a) = 1$ if $a < \eta$ and $g(a) = 0$ if $a \geq \eta$.

6.4.2 Total Variation Minimization

Section 7 defines GTV as a (family of) measure(s) for the variation of local model parameter across edges of the empirical graph. The FL methods in Section 9 are all based on optimizing the GTV, viewed as a function of the local model parameters. However, if we have sufficiently large local datasets $\mathcal{D}^{(i)}$, we might be able to learn useful model parameters separately using ordinary ML methods [24]. We can then use these model parameters to learn the edges of the empirical graph by minimizing the GTV, viewed as a function of the empirical graph [65–67].

6.5 Exercises

Exercise 6.1. Get Weather Data (20 points). Construct an empirical graph \mathcal{G} for the data collected at Finnish Meteorological Institute (FMI) weather stations. The empirical graph should contain at least $n = 5$ nodes which represent five different FMI weather stations. The local dataset $\mathcal{D}^{(i)}$ at node i is constituted by $m = 5$ data points that represent the daily minimum and maximum temperatures during five consecutive days starting with 31 January 2014. The empirical graph should be stored as a `networkx.Graph` containing the local datasets as node attributes (in the form of `numpy` arrays).

Exercise 6.2. Nearest Neighbour Graph (10 points). One approach to learn an empirical graph $\mathcal{G}^{(\text{nn})}$ for a collection of n local datasets is to first compute a measure for the pairwise distances between local datasets and then iterating over all nodes i and connected it with the k nearest nodes (see Section 6.4). Here, $k \in \{0, \dots, n - 1\}$ is a hyper-parameter that controls the number of edges in the resulting empirical graph. Can we be sure that each node i in $\mathcal{G}^{(\text{nn})}$ has at most k neighbours ?

7 A Design Principle for FL

Section 6.2 introduced the empirical graph whose nodes carry local datasets and corresponding local models. This section develops a flexible design principle for FL methods to train local models at nodes of an empirical graph. These methods (see Section 9) combine the information carried by local datasets with the information in their network structure, which is encoded in the weighted edges of the empirical graph.

We couple the training of local models by enforcing approximately identical local model parameters at nodes that carry similar local datasets. Section 7.1 defines the GTV as a measure for the discrepancy between local model parameters at nodes that carry (statistically) similar local datasets.

Section 7.2 then introduces GTVMin to learn local model parameters that optimally balance between their empirical risk incurred on local datasets and their GTV. Section 7.3 discusses several useful interpretations of GTVMin. Section 7.4 extends GTV and GTVMin from parametric to heterogeneous non-parametric local models (such as decision trees).

7.1 Generalized Total Variation

Consider a parametrized networked model with model parameters $\mathbf{w} \in \mathcal{W}$ (see (68)). We can measure the variation of local model parameters $\mathbf{w} \in \mathcal{W}$ across the edges in \mathcal{G} via their variation $\mathbf{u} : e \in \mathcal{E} \mapsto \mathbf{u}^{(e)} \in \mathbb{R}^d$. Formally, the variation assigns each directed edge $e \in \mathcal{E}$ the difference $\mathbf{u}^{(e)} := \mathbf{w}^{(e_+)} - \mathbf{w}^{(e_-)}$. A quantitative measure for the overall variation of the networked model

parameters \mathbf{w} is the GTV

$$\|\mathbf{w}\|_{\text{GTV}} := \sum_{(i,i') \in \mathcal{E}} A_{i,i'} \phi(\mathbf{w}^{(i')} - \mathbf{w}^{(i)}) \text{ with penalty function } \phi(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}. \quad (76)$$

Note that GTV (76) represents an entire ensemble of variation measures. This ensemble is parametrized by the penalty function $\phi(\cdot)$ which we only require to be convex but arbitrary otherwise. Two popular choices are $\phi(\mathbf{v}) := \|\mathbf{v}\|_2$, which is used by network Lasso [68], and $\phi(\mathbf{v}) := (1/2) \|\mathbf{v}\|_2^2$ which is used by “MOCHA” [18]. Another recent FL method based on GTV regularization uses the choice $\phi(\mathbf{v}) := \|\mathbf{v}\|_1$ [69].

Different choices for the penalty function offer different trade-offs between computational aspects and statistical aspects of the resulting FL algorithms (see Section 9). As a case in point, the network Lasso (which uses $\phi(\mathbf{u}) = \|\mathbf{u}\|_2$) is computationally more challenging than MOCHA (which used penalty $\phi(\mathbf{u}) = (1/2)\|\mathbf{u}\|_2^2$). On the other hand, network Lasso is more accurate in learning models for data with specific network structures (such as chains) that are challenging for method that use the smooth penalty $\phi(\mathbf{v}) := (1/2)\|\mathbf{v}\|_2^2$ [70, 71].

7.2 Generalized Total Variation Minimization

We can enforce similar local model parameters $\mathbf{w}^{(i)} \approx \mathbf{w}^{(i')}$, at nodes i, i' connected by edge $\{i, i'\}$ with large weight $A_{i,i'}$, by favouring a small GTV $\|\mathbf{w}\|_{\text{GTV}}$. The extreme case of vanishing GTV is achieved by using identical local model parameters $\mathbf{w}^{(i)} = \mathbf{a}$ for all nodes $i \in \mathcal{V}$. However, we also need to take into account the local loss functions $L_i(\mathbf{w}^{(i)})$ whose minimizers differ

for nodes in different clusters (see Assumption 1). It seems reasonable to learn local model parameters $\mathbf{w}^{(i)}$ by balancing between (the sum of) local loss functions and GTV (76),

$$\hat{\mathbf{w}} \in \arg \min_{\mathbf{w} \in \mathcal{W}} \sum_{i \in \mathcal{V}} L_i(\mathbf{w}^{(i)}) + \lambda \|\mathbf{w}\|_{\text{GTV}}. \quad (77)$$

The regularization parameter $\lambda > 0$ in (77) steers the preference for local model parameters $\mathbf{w}^{(i)}$ with small GTV versus incurring small total loss $\sum_{i \in \mathcal{V}} L_i(\mathbf{w}^{(i)})$. The choice of λ can be guided by cross validation [1, Ch. 8], [34] or by the analysis of solutions to (77) under a clustering assumption (see Section 6.3 and Section 6.3).

By increasing the value of λ in (77) we can enforce the resulting local model parameters $\hat{\mathbf{w}}^{(i)}$ to be constant over increasingly larger subsets of nodes. Choosing λ larger than some critical value, that depends on the shape of the local loss functions and the network structure of \mathcal{G} , results in $\hat{\mathbf{w}}^{(i)}$ being constant over all nodes $i \in \mathcal{V}$. Thus, we effectively train a single global model for all nodes (see Sections 8.1 and 8.2).

GTVMin (77) unifies and considerably extends some well-known methods for distributed optimization and learning. In particular, the network Lasso [68] is obtained from (77) for the choice $\phi(\mathbf{v}) := \|\mathbf{v}\|_2$. The MOCHA method [18] is obtained from (77) for the choice $\phi(\mathbf{v}) := \|\mathbf{v}\|_2^2$. Another special case of (77), obtained for the choice $\phi(\mathbf{v}) := \|\mathbf{v}\|_1$, has been studied recently [69].

7.3 Interpretations

We next discuss some useful relations between GTVMin (77) and basic ML techniques.

7.3.1 Regularized Empirical Risk Minimization

Note that GTVMin (77) is an instance of RERM that uses GTV (76) as regularizer. We measure the empirical risk incurred by the networked model parameters $\mathbf{w} \in \mathcal{W}$ by the sum of the local loss functions $\sum_{i \in \mathcal{V}} L_i(\mathbf{w}^{(i)})$. The interpretation of GTVMin as a form of RERM suggests to apply optimization methods for structured optimization problems to FL applications that can be formulated as GTVMin [72].

7.3.2 Multi-task Learning

Consider networked data with empirical graph \mathcal{G} whose nodes $i \in \mathcal{V}$ carry local datasets $\mathcal{D}^{(i)}$. Each node also carries a local model $\mathcal{H}^{(i)}$ which is trained from $\mathcal{D}^{(i)}$. The training of the local model $\mathcal{H}^{(i)}$ is a separate learning task, one for each node $i \in \mathcal{V}$. GTVMin couples these learning tasks via requiring a small GTV of the resulting local model parameters $\mathbf{w}^{(i)}$. Thus, we can interpret GTVMin as a special case of multi-task learning [73, 74]

7.3.3 Clustering

It can be shown that the solutions of GTVMin are approximately piecewise constant over well-connected subsets of nodes in the empirical graph [58, 61, 70, 75]. Thus, we can interpret GTVMin as a generalization of graph clustering methods such as spectral clustering or flow-based clustering [63, 76].

In contrast to basic graph clustering methods, the clustering delivered by GTVMin is determined not only by the network structure of the empirical graph. Indeed, the resulting cluster structure also depends on the local loss functions $L_i(\cdot)$ which are, in turn, determined by the local datasets $\mathcal{D}^{(i)}$.

We obtain convex clustering [77, 78] as special case of GTVMin (77) with the local loss functions

$$L_i(\mathbf{w}^{(i)}) = \|\mathbf{w}^{(i)} - \mathbf{x}^{(i)}\|_2^2, \text{ for all nodes } i \in \mathcal{V} \quad (78)$$

and GTV penalty $\phi(\mathbf{u}) = \|\mathbf{u}\|_p$ being a p -norm $\|\mathbf{u}\|_p := (\sum_{j=1}^d |u_j|^p)^{1/p}$ with some $p \geq 1$. The feature vectors $\mathbf{x}^{(i)}$ in (78) characterize data points that we wish to cluster. Each cluster is formed by all data points $\mathbf{x}^{(i)}$ that have identical local model parameters $\mathbf{w}^{(i)}$.

7.3.4 Locally Weighted Learning

Another interpretation of GTVMin is that of locally weighted learning [79]. Indeed, the local model parameters delivered by GTVMin are approximately constant over well-connected subsets of nodes in the empirical graph [58, 61, 70, 75]. Thus, for each node i in a given cluster \mathcal{C} , the solution $\widehat{\mathbf{w}}^{(i)}$ of GTV approximates the solution $\overline{\mathbf{w}}^{(\mathcal{C})}$ of (71), $\widehat{\mathbf{w}}^{(i)} \approx \overline{\mathbf{w}}^{(\mathcal{C})}$. Note that the cluster-wide optimization (71) can be written as a locally weighted learning problem [79, Sec. 3.1.2]

$$\overline{\mathbf{w}}^{(\mathcal{C})} = \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{i' \in \mathcal{V}} \rho_{i'} L_{i'}(\mathbf{w}). \quad (79)$$

The locally-weighted learning problem (79) involves the weights $\rho_{i'}$ which are equal to 1 if $i' \in \mathcal{C}$ and equal to 0 otherwise.

7.4 Non-Parametric Models

So far, we focused on networked federated learning (NFL) methods for local models that are parametrized by local model parameters $\mathbf{w}^{(i)}$, for each node

$i \in \mathcal{V}$. However, many important ML methods such as decision trees are non-parametric. We will now modify GTVMin for parametric local models to non-parametric local models (hypothesis spaces) $\mathcal{H}^{(i)}$ for nodes $i \in \mathcal{V}$. To this end, we first extend the concept of GTV to obtain a measure for the variation of local hypothesis maps $h^{(i)} \in \mathcal{H}^{(i)}$ from non-parametric models $\mathcal{H}^{(i)}$.

The idea is to measure the variation between $h^{(i)}$ and $h^{(i')}$, across an edge $\{i, i'\}$, via the discrepancy between their predictions on a common test-set

$$\mathcal{D}^{(\text{test})} = \left\{ \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m')} \right\}. \quad (80)$$

The test-set $\mathcal{D}^{(\text{test})}$ must be shared with each node i, i' of an edge $\{i, i'\}$. Each of these two nodes can then compute the discrepancy

$$\begin{aligned} d(h^{(i)}, h^{(i')}) := \\ (1/m') \sum_{r=1}^{m'} L^{(\text{d})}(\mathbf{x}^{(r)}, h^{(i')}(\mathbf{x}^{(r)}), h^{(i)}(\mathbf{x}^{(r)})). \end{aligned} \quad (81)$$

In principle, we can use different test-sets to compute the discrepancy over different edges in the empirical graph.

We are now in the position to generalize GTV (76) for parametric local models to non-parametric local models. This generalization is obtained by summing the discrepancies (81) over all edges \mathcal{E} ,

$$\text{GTV}\{h\} := \sum_{\{i, i'\} \in \mathcal{E}} A_{i, i'} d(h^{(i)}, h^{(i')}). \quad (82)$$

Note that $\text{GTV}\{h\}$ is parametrized by the choice for the loss function $L^{(\text{d})}(\cdot, \cdot, \cdot)$ used to compute the discrepancy $d(h^{(i)}, h^{(i')})$ (81). Note that the loss function $L^{(\text{d})}(\mathbf{x}, h(\mathbf{x}), h'(\mathbf{x}))$ depends on the hypothesis maps h, h' only via their predictions $h(\mathbf{x}), h'(\mathbf{x})$.

It is important to note that the loss function $L^{(d)}(\cdot, \cdot, \cdot)$ in (81) can be different from the loss function $L_i(\cdot)$ (67) used to train a local model. For example, the local models might be trained using a loss function for multi-class classification into 10 different classes while the loss function $L^{(d)}(\cdot, \cdot, \cdot)$ might be the squared error loss.

We obtain GTVMin for non-parametric local models by using $\text{GTV}\{h\}$ as regularization term in RERM,

$$\hat{h} \in \operatorname{argmin}_{h \in \mathcal{H}^{(\mathcal{G})}} \sum_{i \in \mathcal{V}} L_i(h^{(i)}) + \lambda \text{GTV}\{h\}. \quad (83)$$

Here, we use a networked hypothesis space $\mathcal{H}^{(\mathcal{G})} = \mathcal{H}^{(1)} \times \dots \times \mathcal{H}^{(n)}$ whose elements are maps that assign a hypothesis $h^{(i)} \in \mathcal{H}^{(i)}$ to each node $i \in \mathcal{V}$. Note that (83) is parametrized by the choice for the loss function used to compute the discrepancy (81). Different choices for the loss function in (81) result in different computational aspects and statistical aspects of (83). If we use the squared error loss to measure the discrepancy (81), (83) specializes to

$$\hat{h} \in \operatorname{argmin}_{h \in \mathcal{H}^{(\mathcal{G})}} \sum_{i \in \mathcal{V}} L_i(h^{(i)}) + (\lambda/m') \sum_{(i,i') \in \mathcal{E}} A_{i,i'} \sum_{r=1}^{m'} \left(h^{(i)}(\mathbf{x}^{(r)}) - h^{(i')}(\mathbf{x}^{(r)}) \right)^2. \quad (84)$$

7.5 Exercises

Exercise 7.1. Generalizing TV to Decision Trees Discuss possible modifications of GTVMin (77) to applications that use decision trees as local models.

8 Main Flavours of Federated Learning

Let us now discuss how different choices for the empirical graph and local loss functions used in GTVMin result in main flavours of FL [80]. The most basic FL setup is obtained when the empirical graph is a star with the centre node acting as a server or coordinator (see Section 8.1). Section 8.2 explains decentralized FL where an arbitrary empirical graph is used for message passing implementations of FL algorithms to collaboratively train a single global model.

Note that GTVMin uses a regularization parameter λ to control the (strength of the) coupling between the local models at nodes that are connected by an edge in the empirical graph. Section 8.3 discusses clustered FL which uses this parameter to steer the clustering structure of the learnt local models.

One extreme case of clustered FL is when the empirical graph becomes one single cluster. This corresponds to decentralized FL of a global model that is shared among all nodes. Another extreme case is personalized FL, which is discussed in Section 8.4, where each node in the empirical graph becomes a separate cluster.

Section 8.5 discusses vertical FL from different local datasets that correspond to the same data points but use (have access to) different features. Section 8.6 discusses horizontal FL from local datasets that contain different data points but use the same features to characterize them [81, Ch. 4].

8.1 Centralized Federated Learning

Probably the most basic FL setup uses a server which is connected by separate links to each client. The clients generate local datasets which are used to train a single global model \mathcal{H} . FL methods use the server to maintain the current model parameters [14]. These global model parameters are broadcasted to the clients which compute model parameters updates based on their local datasets (e.g., using gradient steps). These local updates are sent back to the server which aggregates them to obtain new global model parameters.

The above server-based (centralized) FL setup correspond to GTVMin for an empirical graph being a star graph (see Figure 19).

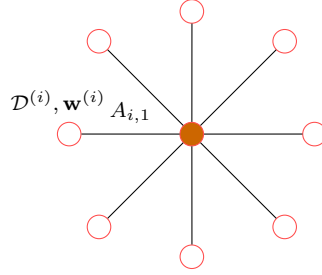


Figure 19: Star graph $\mathcal{G}^{(\text{star})}$ with centre node representing a server and peripheral nodes representing clients that generate local datasets.

8.2 Decentralized Federated Learning

[82] empirical graph is arbitrary; GTVMin parameter chosen so large that all local model parameters equal;

8.3 Clustered Federated Learning

Clustered FL uses different forms of a clustering assumption that requires local datasets, and their associated learning tasks, to form clusters (see Assumption 1) [22, 23, 73, 83, 84]. Local datasets belonging to the same cluster have similar statistical properties and, in turn, similar optimal parameter values for the corresponding local models. Our recent work offers a precise characterization of the cluster structure and local loss functions that allow FL methods from Section 9 to pool local datasets that form a cluster of statistically homogeneous local datasets [75]. This characterization also provides some guidance for choosing the regularization parameter λ in GTVMin (77).

8.4 Personalized Federated Learning

Using a sufficiently small regularization parameter in GTVMin (77) allows the local model parameters to vary over the nodes $i \in \mathcal{V}$. However, these trained personalized local models might still be coupled partially. For example, if local models are deep ANNs we might enforce the parameters of input layers to be identical while the parameters for the deeper layers might be different for each local dataset.

The partial parameter sharing for local models can be implemented in many different ways [85, Sec. 4.3.]. We can use GTVMin (77) for partial parameter sharing by using a specific choice for the GTV penalty function. In particular, we could construct the penalty function as a combination of two terms,

$$\phi(\mathbf{w}^{(i)} - \mathbf{w}^{(i')}) := \lambda^{(1)}\phi^{(1)}(\mathbf{w}^{(i)} - \mathbf{w}^{(i')}) + \lambda^{(2)}\phi^{(2)}(\mathbf{w}^{(i)} - \mathbf{w}^{(i')}). \quad (85)$$

Each component $\phi^{(1)}, \phi^{(2)}$ measures different parts (subsets) of the variation $\mathbf{w}^{(i)} - \mathbf{w}^{(i')}$ of local model parameters at connected nodes $\{i, i'\} \in \mathcal{E}$.

Moreover, we might use different regularization strengths $\lambda^{(1)}$ and $\lambda^{(2)}$ for different penalty components in (85) to enforce different subsets of the model parameters to be clustered with different granularity (cluster size). For local models being deep ANNs, we typically enforcing the low-level layers (closer to the input) to have same model parameters (weights and bias terms), while deeper (closer to the output) layers can have different model parameters. Yet another technique for partial sharing of model parameters is to train a hyper-model which, in turn, is used to initialize the training of local models [86].

For heterogeneous collections of local models, including non-parametric local models such as decision trees, we can implement a partial coupling by carefully choosing the loss function in the discrepancy measure (81). In particular, we might use a loss function $L^{(d)}(\cdot, \cdot, \cdot)$ that is different from the loss used to measure the quality of a hypothesis $h^{(i)}$. Loosely speaking, the loss function $L^{(d)}(\cdot, \cdot, \cdot)$ might not be directly related to the ultimate learning task of the local models (which could be binary classification), but instead correspond to an auxiliary task (such as dimensionality reduction) which partially couples local models.

8.5 Vertical Federated Learning

Vertical FL refers to applications with local datasets that are constituted by the same data points but characterizing them with different features [87]. In particular, vertical FL applications revolve around an underlying global

dataset

$$\mathcal{D}^{(\text{global})} := \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}.$$

Each data point in the global dataset is characterized by d' features $\mathbf{x}^{(r)} = (x_1^{(r)}, \dots, x_{d'}^{(r)})^T$. The global dataset can only be accessed indirectly via local datasets that use different subsets of the feature vectors $\mathbf{x}^{(r)}$.

Formally, the local dataset $\mathcal{D}^{(i)}$ contains the pairs $(\mathbf{x}^{(i,r)}, y^{(i,r)})$, for $r = 1, \dots, m$. The labels $y^{(i,r)}$ are identical to the labels in the global dataset, $y^{(i,r)} = y^{(r)}$. The feature vectors $\mathbf{x}^{(i,r)}$ are obtained by a subset $\mathcal{F}^{(i)} := \{j_1, \dots, j_d\}$ of the original d' features in $\mathbf{x}^{(r)}$,

$$\mathbf{x}^{(i,r)} = (x_{j_1}^{(r)}, \dots, x_{j_d}^{(r)})^T.$$

8.6 Horizontal Federated Learning

Horizontal FL refers to applications with local datasets $\mathcal{D}^{(i)}$, for $i \in \mathcal{V}$ that contain data points characterized by the same features [88]. The local datasets $\mathcal{D}^{(i)}$ are obtained by different subsets of an underlying global dataset

$$\mathcal{D}^{(\text{global})} := \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}.$$

The local dataset $\mathcal{D}^{(i)}$ is constituted by the data points of $\mathcal{D}^{(\text{global})}$ with indices in $\{r_1, \dots, r_{m_i}\}$,

$$\mathcal{D}^{(i)} := \{(\mathbf{x}^{(r_1)}, y^{(r_1)}), \dots, (\mathbf{x}^{(r_{m_i})}, y^{(r_{m_i})})\}.$$

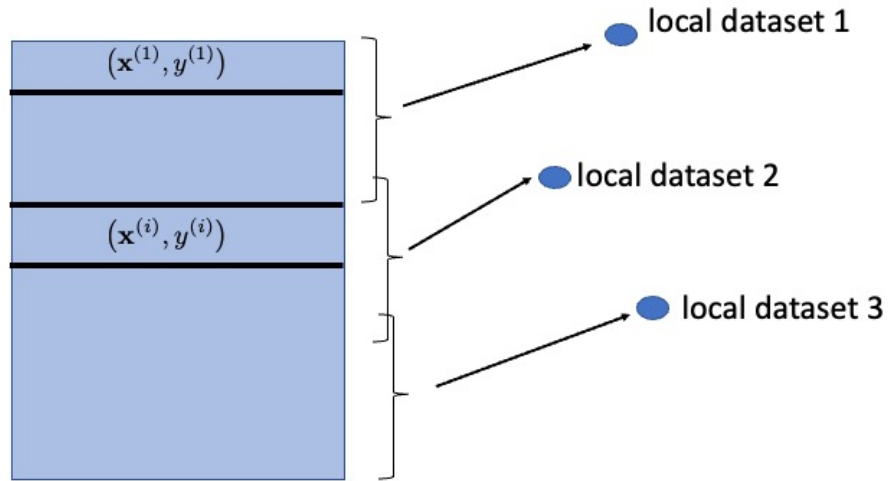


Figure 20: Horizontal FL uses the same features to characterize data points in different local datasets. Different local datasets are constituted by different subsets of an underlying global dataset. As an example consider local datasets that are constituted by patient records using the same type of personal information (social security number, home address). The underlying global dataset consists of the union of all patient records.

9 Federated Learning Algorithms

Section 7 introduced GTVMin as a design principle for FL methods that couple the training of local models on their corresponding local datasets. We obtain FL algorithms by the application of optimization methods to solve GTVMin.

Throughout this chapter we will focus on the special case of GTVMin (77) for the penalty function $\phi(\mathbf{u}) = \|\mathbf{u}\|_2^2$,

$$\hat{\mathbf{w}} \in \arg \min_{\mathbf{w} \in \mathcal{W}} \sum_{i \in \mathcal{V}} L_i(\mathbf{w}^{(i)}) + \lambda \sum_{\{i, i'\} \in \mathcal{E}} A_{i, i'} \left\| \mathbf{w}^{(i)} - \mathbf{w}^{(i')} \right\|_2^2. \quad (86)$$

We obtain different FL methods by applying different iterative optimization methods to solve GTVMin (86). These optimization methods construct a sequence

$$\hat{\mathbf{w}}_0, \hat{\mathbf{w}}_1, \dots$$

of local model parameters that are (hopefully) increasingly accurate approximations of the solutions $\hat{\mathbf{w}}$ to (86).

9.1 FedRelax

We now apply block-coordinate minimization [26, 36] to solve GTVMin (86).

To this end, we rewrite (86) as

$$\begin{aligned} \hat{\mathbf{w}} \in \arg \min_{\mathbf{w} \in \mathcal{W}} \underbrace{\sum_{i \in \mathcal{V}} f^{(i)}(\mathbf{w})}_{:= f^{(\text{GTV})}(\mathbf{w})} \\ \text{with } f^{(i)}(\mathbf{w}) := L_i(\mathbf{w}^{(i)}) + (\lambda/2) \sum_{i' \in \mathcal{N}^{(i)}} A_{i, i'} \left\| \mathbf{w}^{(i)} - \mathbf{w}^{(i')} \right\|_2^2. \end{aligned} \quad (87)$$

The objective function in (87) decomposes into node-wise functions $f^{(i)}(\mathbf{w})$, for each $i \in \mathcal{V}$. Block-coordinate minimization exploits this structure to decouple the optimization of local model parameters $\widehat{\mathbf{w}}^{(i)}$.

We next discuss a basic variant of block-coordinate minimization to solve (87). Given the current local model parameters $\widehat{\mathbf{w}}_r$, we compute (hopefully improved) updated local model parameters $\widehat{\mathbf{w}}_{r+1}^{(i)}$ by minimizing $f^{(\text{GTV})}(\cdot)$ along $\mathbf{w}^{(i)}$ starting from $\widehat{\mathbf{w}}_r$,

$$\begin{aligned} \widehat{\mathbf{w}}_{r+1}^{(i)} &\in \operatorname{argmin}_{\mathbf{w}^{(i)} \in \mathbb{R}^d} f^{(\text{GTV})} \left(\widehat{\mathbf{w}}_r^{(1)}, \dots, \widehat{\mathbf{w}}_r^{(i-1)}, \mathbf{w}^{(i)}, \widehat{\mathbf{w}}_r^{(i+1)}, \dots \right) \\ &\stackrel{(87)}{=} \operatorname{argmin}_{\mathbf{w}^{(i)} \in \mathbb{R}^d} f^{(i)} \left(\widehat{\mathbf{w}}_r^{(1)}, \dots, \widehat{\mathbf{w}}_r^{(i-1)}, \mathbf{w}^{(i)}, \widehat{\mathbf{w}}_r^{(i+1)}, \dots \right) \\ &\stackrel{(87)}{=} \operatorname{argmin}_{\mathbf{w}^{(i)} \in \mathbb{R}^d} L_i(\mathbf{w}^{(i)}) + (\lambda/2) \sum_{i' \in \mathcal{N}^{(i)}} A_{i,i'} \left\| \mathbf{w}^{(i)} - \widehat{\mathbf{w}}_r^{(i')} \right\|_2^2. \end{aligned} \quad (88)$$

Algorithm 1 simply repeats the update (88) for a sufficient number of iterations, i.e., until a stopping criterion is met.

Non-Parametric Models. Algorithm 1 can only be used for parametric local models $\mathcal{H}^{(i)}$ such as linear regression or ANNs with a fixed number d of parameters (which is the same for all nodes $i \in \mathcal{V}$). However, we can naturally extend the scope of Algorithm 1 to non-parametric models by applying block-coordinate minimization to the non-parametric GTVMin variant (83) instead of (86). We obtain the update

$$\widehat{h}_{r+1}^{(i)} \in \operatorname{argmin}_{h \in \mathcal{H}^{(i)}} L_i(h^{(i)}) + (\lambda/2) \sum_{i' \in \mathcal{N}^{(i)}} A_{i,i'} \underbrace{d(h^{(i)}, \widehat{h}_r^{(i')})}_{\text{see (81)}}. \quad (90)$$

Replacing the update (88) with (90) in step 4 of Algorithm 1 makes it applicable to non-parametric local models. The update (90) is obtained by

Algorithm 1 FedRelax for Parametric Local models

Input: empirical graph \mathcal{G} ; local loss functions $L_i(\cdot)$, GTV parameter λ

Output: learnt local model parameters $\widehat{\mathbf{w}}^{(i)}$

Initialize: $r := 0$; $\widehat{\mathbf{w}}_0^{(i)} := \mathbf{0}$

- 1: **while** stopping criterion is not satisfied **do**
- 2: **for** all nodes $i \in \mathcal{V}$ (simultaneously) **do**
- 3: share local parameters $\widehat{\mathbf{w}}_r^{(i)}$ with all neighbours $i' \in \mathcal{N}^{(i)}$
- 4: update local model parameters via (see (88))

$$\widehat{\mathbf{w}}_{r+1}^{(i)} := \underset{\mathbf{w}^{(i)} \in \mathbb{R}^d}{\operatorname{argmin}} L_i(\mathbf{w}^{(i)}) + (\lambda/2) \sum_{i' \in \mathcal{N}^{(i)}} A_{i,i'} \left\| \mathbf{w}^{(i)} - \widehat{\mathbf{w}}_r^{(i')} \right\|_2^2 \quad (89)$$

- 5: **end for**
 - 6: $r := r + 1$
 - 7: **end while**
 - 8: $\widehat{\mathbf{w}}^{(i)} := \widehat{\mathbf{w}}_r^{(i)}$ for all nodes $i \in \mathcal{V}$
-

using the discrepancy (81) between local models within the GTV measure (82).

We obtain Algorithm 2 when using the squared error loss to measure the discrepancy (81) between local hypothesis maps $h^{(i)}$ and $h^{(i')}$ (see (84)).

Algorithm 2 FedRelax for Non-Parametric Models

Input: empirical graph \mathcal{G} with edge weights $A_{i,i'}$; local loss functions $L_i(\cdot)$; test-set $\mathcal{D}' = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m')}\}$; GTV regularization parameter λ

Initialize: $r := 0$; $\widehat{h}_0^{(i)} \equiv 0$ for all nodes $i \in \mathcal{V}$

- 1: **while** stopping criterion is not satisfied **do**
- 2: **for** all nodes $i \in \mathcal{V}$ in parallel **do**
- 3: share test-set labels $\left\{ \widehat{h}_r^{(i)}(\mathbf{x}) \right\}_{\mathbf{x} \in \mathcal{D}'(\text{test})}$, with neighbours $i' \in \mathcal{N}^{(i)}$
- 4: update hypothesis $\widehat{h}_r^{(i)}$ as follows:

$$\widehat{h}_{r+1}^{(i)} \in \underset{h^{(i)} \in \mathcal{H}^{(i)}}{\operatorname{argmin}} \left[L_i(h^{(i)}) + (\lambda/2) \sum_{i' \in \mathcal{N}^{(i)}} (A_{i,i'}/m') \sum_{r=1}^{m'} \left(h^{(i)}(\mathbf{x}^{(r)}) - \widehat{h}_r^{(i')}(\mathbf{x}^{(r)}) \right)^2 \right]. \quad (91)$$

- 5: **end for**
 - 6: $r := r + 1$
 - 7: **end while**
-

9.2 FedSGD

Consider the GTVMin instance (86) for local loss functions $L_i(\mathbf{w}^{(i)})$ being differentiable.

9.3 FedAvg

The popular federated averaging (FedAvg) method is obtained from (86) for $\lambda \rightarrow \infty$. Indeed, for sufficiently large λ , the penalty term in (86) dominates, enforcing the learnt local parameter vectors to be nearly identical. Thus, for sufficiently large λ , we can approximate (86) by

$$\begin{aligned} \hat{\mathbf{w}} &\in \arg \min_{\mathbf{w} \in \mathcal{C}} \sum_{i \in \mathcal{V}} L_i(\mathbf{w}^{(i)}) \\ \text{with } \mathcal{C} &= \{\mathbf{w} : \mathbf{w}^{(i)} = \mathbf{w}^{(i')} \text{ for any edge } \{i, i'\} \in \mathcal{E}\}. \end{aligned} \quad (92)$$

We can solve (92) using projected GD.

Let us spell out the special case of Algorithm 3 obtained when learning local linear models by minimizing the average squared error loss (4) incurred on the local dataset

$$\mathcal{D}^{(i)} := \{(\mathbf{x}^{(i,1)}, y^{(i,1)}), \dots, (\mathbf{x}^{(i,m_i)}, y^{(i,m_i)})\}.$$

The resulting Algorithm 4 performs regularized tailored linear regression for each local dataset. The regularization is implemented by the GTV (76) of the parameters $\mathbf{w}^{(i)}$ for the node-wise linear models.

9.4 Asynchronous Computation

The FL methods presented in Section 9.1 - Section 9.3 assume a perfectly synchronous mode of computation. As a case in point, consider the synchronous mode of operation required by Algorithm 1. In particular, the step 3 of Algorithm 1 has to be carried simultaneously at all nodes. Only when step 3 has been completed at all nodes, Algorithm 1 can continue with executing

Algorithm 3 FedAvg

Input: client list \mathcal{V}

Server. (available under `fljung.cs.aalto.fi`)

Initialize. $r := 0$

- 1: **while** stopping criterion is not satisfied **do**
- 2: receive local model parameters $\mathbf{w}^{(i)}$ for all clients $i \in \mathcal{V}$
- 3: update global model parameters

$$\bar{\mathbf{w}}[r] := (1/|\mathcal{V}|) \sum_{i \in \mathcal{V}} \mathbf{w}^{(i)}.$$

- 4: send new global model parameters $\bar{\mathbf{w}}[r]$ to all clients $i \in \mathcal{V}$
- 5: $r := r + 1$
- 6: **end while**

Client. (some $i \in \mathcal{V}$)

- 1: **while** stopping criterion is not satisfied **do**
- 2: receive global parameter vector $\bar{\mathbf{w}}[r]$ from server
- 3: update local parameters by RERM (29)

$$\mathbf{w}^{(i)} := \operatorname{argmin}_{\mathbf{v} \in \mathbb{R}^d} \left[L_i(\mathbf{v}) + \lambda \|\mathbf{v} - \bar{\mathbf{w}}[r]\|^2 \right].$$

- 4: send $\mathbf{w}^{(i)}$ to `fljung.cs.aalto.fi`
 - 5: **end while**
-

Algorithm 4 FedAvg for linear regression

Input: client list \mathcal{V}

Server. (available under `fljung.cs.aalto.fi`)

Initialize. $r := 0$

- 1: **while** stopping criterion is not satisfied **do**
- 2: receive local model parameters $\mathbf{w}^{(i)}$ from all clients $i \in \mathcal{V}$
- 3: update global model parameters

$$\bar{\mathbf{w}}[r] := (1/|\mathcal{V}|) \sum_{i \in \mathcal{V}} \mathbf{w}^{(i)}.$$

- 4: send updated global model parameters $\bar{\mathbf{w}}[r]$ to all clients $i \in \mathcal{V}$
- 5: $r := r + 1$
- 6: **end while**

Client. (at some node $i \in \mathcal{V}$)

- 1: **while** stopping criterion is not satisfied **do**
- 2: receive global model parameters $\bar{\mathbf{w}}[r]$ from server
- 3: update local model parameters by RERM (see (29))

$$\mathbf{w}^{(i)} := \operatorname{argmin}_{\mathbf{v} \in \mathbb{R}^d} \left[(1/m_i) \sum_{r=1}^{m_i} (\mathbf{v}^T \mathbf{x}^{(i,r)} - y^{(i,r)})^2 + \lambda \|\mathbf{v} - \bar{\mathbf{w}}[r]\|_2^2 \right].$$

- 4: send $\mathbf{w}^{(i)}$ to `fljung.cs.aalto.fi`
 - 5: **end while**
-

step 4. However, what happens if some of the nodes $i \in \mathcal{V}$ fail to complete step 3 of Algorithm 1?

9.5 Exercises

Exercise 9.1. FedRelax uses Proximity Operator Show that the coordinate minimization update (88) is equivalent to the evaluation of the proximity operator $\text{prox}_{L_i(\cdot), \rho}(\mathbf{w}')$ for a specific vector $\mathbf{w}' \in \mathbb{R}^d$ and $\rho > 0$.

Exercise 9.2. FedAvg Linear Regression Discuss how step 3 of Algorithm 4 could be computed using the `LinearRegression.fit()` method provided by the Python package `scikit-learn`.

Part III

Trustworthy Federated Learning

10 Requirements for Trustworthy AI

FL is an important subfield of AI. As part of their AI strategy, the European Commission set up the High-Level Expert Group on Artificial Intelligence(AI HLEG) in 2018. This group put forward seven key requirements for trustworthy AI [89]:

1. **human agency and oversight.** *AI systems should empower human beings, allowing them to make informed decisions and fostering their fundamental rights. At the same time, proper oversight mechanisms need to be ensured, which can be achieved through human-in-the-loop, human-on-the-loop, and human-in-command approaches*
2. **technical robustness and safety.** *AI systems need to be resilient and secure. They need to be safe, ensuring a fall back plan in case something goes wrong, as well as being accurate, reliable and reproducible. That is the only way to ensure that also unintentional harm can be minimized and prevented.*
3. **privacy and data governance.** *besides ensuring full respect for privacy and data protection, adequate data governance mechanisms must also be ensured, taking into account the quality and integrity of the data, and ensuring legitimised access to data.*
4. **transparency.** *the data, system and AI business models should be transparent. Traceability mechanisms can help achieving this. Moreover, AI systems and their decisions should be explained in a manner adapted to the stakeholder concerned. Humans need to be aware that they are*

interacting with an AI system, and must be informed of the system's capabilities and limitations.

5. diversity non-discrimination and fairness,
6. societal and environmental well-being
7. accountability.

And I will now try to integrate these key requirements using the tools that we have built up in this course and these tools mainly revolve around design choices. Specifically, it's a design choice for data, for the model, the set of hypothesis spaces and for the loss function. So, I will show you. How these design choices could be made to satisfy these key requirements.

11 Privacy Protection

A main application domain of FL is healthcare where local datasets are collected at different healthcare providers such as hospitals or laboratories [90]. Let us assume that local datasets contain data points that represent human beings. The features of data points are different bio-physical measurements and diagnosis reports. The label might be the occurrence of some disease such as diabetes.

Some of these features and labels stored in a local dataset might carry sensitive information (e.g., presence of some infection) and cannot be shared beyond the owner of the local dataset (which could be a blood lab). We refer to such properties as private and need to ensure that they are not shared beyond the owner of the local dataset. However, the (non-sensitive part of the) information contained in the local dataset might be helpful to train high-dimensional ML models for data-driven diagnosis or personalized medication.

The FL methods discussed in Section 9 allow to share the information contained in local datasets to train local models in a privacy-friendly fashion. By privacy-friendly, we mean that no information about private properties (either a feature or a label) of a data point is revealed to anybody beyond the owner (or controller) of the local dataset. As a case in point, we will study the information flow during the iterations of Algorithm 1. Sections 11.1 and 11.2 will then explain modifications of Algorithm 1 to avoid the leakage of sensitive information. Note that Algorithm 1 only applies to parametrized hypothesis spaces which contain hypothesis maps $h^{(\mathbf{w}^{(i)})}$ that are fully determined by a parameter vector $\mathbf{w}^{(i)}$.

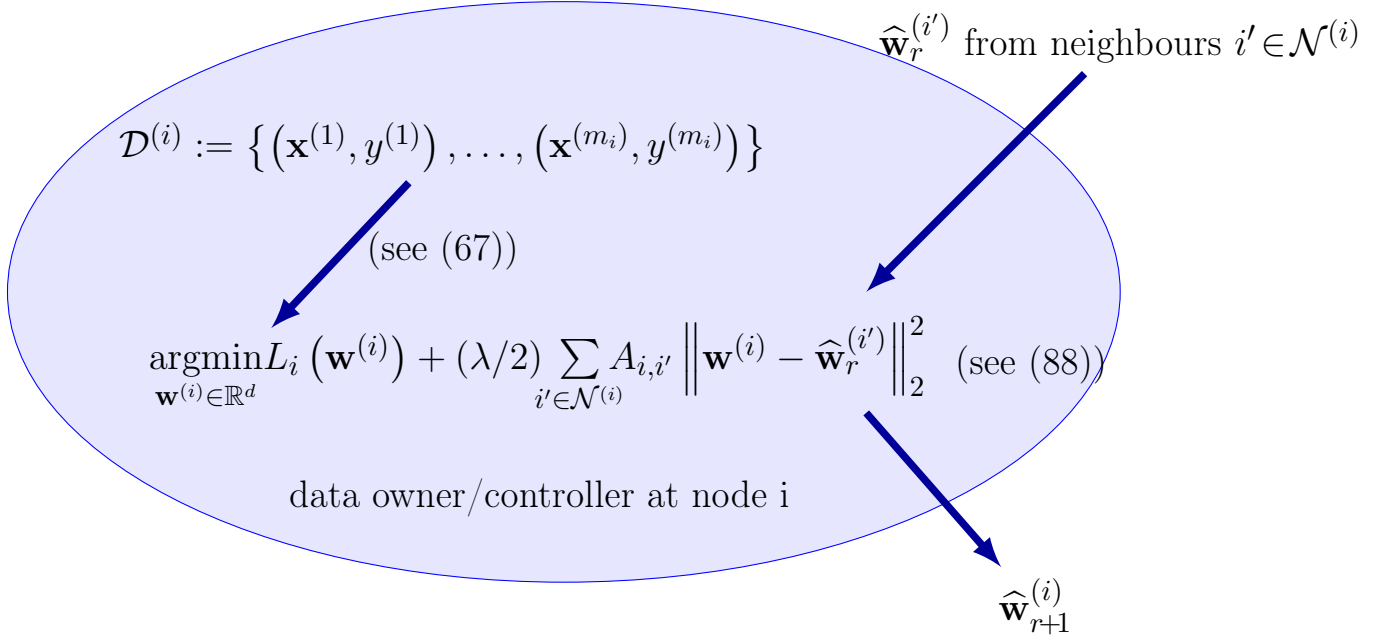


Figure 21: Data flow at node $i \in \mathcal{V}$ during a single iteration of Algorithm 1.

Figure 21 illustrates the data (and information) flow arising at node i during a single iteration of Algorithm 1. During each iteration of Algorithm 1, each node $i \in \mathcal{V}$ computes the update (88) and shares the updated local parameters $\widehat{\mathbf{w}}_{r+1}^{(i)}$ with its neighbours $\mathcal{N}^{(i)}$. It is the sharing of the updated parameters in step 3 of Algorithm 1 that might cause leakage of sensitive information.

We next discuss two different approaches to avoid this leakage of sensitive information. The first approach in Section 11.1 is to perturb (e.g., by adding noise) the update $\widehat{\mathbf{w}}_{r+1}^{(i)}$ before sharing it with the neighbours $\mathcal{N}^{(i)}$. Section 11.2 presents a complementary approach which is to directly perturb the features and labels of the data points in the local dataset $\mathcal{D}^{(i)}$.

11.1 Adding Noise

We can formalize the notion of privacy leakage and privacy protection with a probabilistic model (see Section 3.1). To this end, we interpret data points in the local dataset $\mathcal{D}^{(i)}$ as realizations of RVs. This implies, in turn, that the updates $\widehat{\mathbf{w}}_r^{(i)}$ are realizations of RVs. Moreover, a private feature $x_j^{(r)}$ becomes a realization of a RV. If the updates $\widehat{\mathbf{w}}_r^{(i)}$ could be used by an adversary to estimate or guess the value of $x_j^{(r)}$, they would leak sensitive information about private features. A quantitative measure for how much observing one RV helps to estimate another RV is mutual information (see Section 14).

To avoid leakage of sensitive information, which is carried by the private feature $x_j^{(r)}$, we need to ensure a small mutual information $I\left(x_j^{(r)}, \widehat{\mathbf{w}}_r^{(i)}\right)$. This mutual information is determined by the joint probability distribution of the local datasets $\mathcal{D}^{(i)}$, for $i \in \mathcal{V}$, and the parameters of Algorithm 1. If $I\left(x_j^{(r)}, \widehat{\mathbf{w}}_r^{(i)}\right)$ is too large, we replace the updates $\widehat{\mathbf{w}}_r^{(i)}$ in steps 3 and 4 of Algorithm 1 by the perturbed updates [91, 92]

$$\widetilde{\mathbf{w}}_r^{(i)} := \widehat{\mathbf{w}}_r^{(i)} + \sigma \boldsymbol{\varepsilon}^{(i,r)}. \quad (93)$$

Here, $\boldsymbol{\varepsilon}^{(i,r)}$ are realizations of i.i.d. RVs with common probability distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$. The parameter σ in (93) controls the amount of privacy protection offered by (93). By increasing the value of σ , we obtain an increased protection of private features, however, at the cost of deteriorating accuracy of the hypothesis obtained from $\widetilde{\mathbf{w}}_r^{(i)}$. We summarize the resulting modification of Algorithm 1 in Algorithm 5.

Algorithm 5 FedRelax with Perturbed Updates

Input: empirical graph \mathcal{G} ; local loss functions $L_i(\cdot)$, GTV parameter λ ; privacy-protection level σ

Output: learnt local model parameters $\widehat{\mathbf{w}}^{(i)}$

Initialize: $r := 0$; $\widehat{\mathbf{w}}_0^{(i)} := \mathbf{0}$

- 1: **while** stopping criterion is not satisfied **do**
 - 2: **for** all nodes $i \in \mathcal{V}$ (simultaneously) **do**
 - 3: share perturbed local parameters $\widetilde{\mathbf{w}}_r^{(i)}$ with all neighbours $i' \in \mathcal{N}^{(i)}$
 - 4: update local parameters and add noise
$$\widetilde{\mathbf{w}}_{r+1}^{(i)} := \underset{\mathbf{w}^{(i)} \in \mathbb{R}^d}{\operatorname{argmin}} \left[L_i(\mathbf{w}^{(i)}) + (\lambda/2) \sum_{i' \in \mathcal{N}^{(i)}} A_{i,i'} \left\| \mathbf{w}^{(i)} - \widetilde{\mathbf{w}}_r^{(i')} \right\|_2^2 \right] + \sigma \boldsymbol{\varepsilon}^{(i,r)}. \quad (94)$$
 - 5: **end for**
 - 6: $r := r + 1$
 - 7: **end while**
 - 8: $\widehat{\mathbf{w}}^{(i)} := \widetilde{\mathbf{w}}_r^{(i)}$ for all nodes $i \in \mathcal{V}$
-

11.2 Feature Perturbation

Let us now describe another modification of Algorithm 1 that protects a private feature $x_j^{(r)}$ of data points in the local dataset $\mathcal{D}^{(i)}$. Instead of perturbing the result of the updates 4 in Algorithm 1, we directly perturb the features of the data points in $\mathcal{D}^{(i)}$. The perturbed features result, in turn, in a perturbed local loss function $\tilde{L}_i(\cdot)$.

To perturb the features of the data points in the local dataset $\mathcal{D}^{(i)}$, we apply a feature map [1, Sec. 9.5]

$$\Phi : \mathbf{x} \mapsto \mathbf{z} = \Phi(\mathbf{x}). \quad (95)$$

A privacy-preserving feature map (95) delivers transformed features \mathbf{z} that do not allow to predict (infer) the private feature (too well). However, the new features \mathbf{z} should still allow to learn a hypothesis $h^{(\mathbf{w}^{(i)})}$ that accurately predicts the label of a data point from its privacy-preserving features \mathbf{z} . Such a hypothesis can be found via ERM (see Section 3.2), i.e., by minimizing the average loss

$$\tilde{L}_i(\mathbf{w}^{(i)}) := (1/m_i) \sum_{r=1}^{m_i} L\left((\Phi(\mathbf{x}^{(r)}), y^{(r)}), h^{(\mathbf{w}^{(i)})}\right). \quad (96)$$

We obtain Algorithm 6 from Algorithm 1 by replacing the local loss function $L_i(\cdot)$ with the perturbed version $\tilde{L}_i(\cdot)$ (96).

Examples for privacy-preserving feature maps (95) that could be used for Algorithm 6 include

- “information obfuscation” by adding noise $\Phi(\mathbf{x}) = \mathbf{x} + \text{noise}$ (see [92, Fig. 3])

Algorithm 6 FedRelax with Perturbed Features

Input: empirical graph \mathcal{G} ; feature map Φ , local dataset $\mathcal{D}^{(i)}$ for each node $i \in \mathcal{V}$, GTV parameter λ

Output: learnt local model parameters $\hat{\mathbf{w}}^{(i)}$

Initialize: $r := 0$; $\hat{\mathbf{w}}_0^{(i)} := \mathbf{0}$

- 1: **while** stopping criterion is not satisfied **do**
- 2: **for** all nodes $i \in \mathcal{V}$ (simultaneously) **do**
- 3: share local parameters $\hat{\mathbf{w}}_r^{(i)}$ with all neighbours $i' \in \mathcal{N}^{(i)}$
- 4: update local parameters

$$\begin{aligned} \hat{\mathbf{w}}_{r+1}^{(i)} := \operatorname{argmin}_{\mathbf{w}^{(i)} \in \mathbb{R}^d} & \left[(1/m_i) \sum_{r=1}^{m_i} L \left((\Phi(\mathbf{x}^{(r)}), y^{(r)}), h(\mathbf{w}^{(i)}) \right) + \right. \\ & \left. (\lambda/2) \sum_{i' \in \mathcal{N}^{(i)}} A_{i,i'} \left\| \mathbf{w}^{(i)} - \hat{\mathbf{w}}_r^{(i')} \right\|_2^2 \right] \end{aligned} \quad (97)$$

- 5: **end for**
 - 6: $r := r + 1$
 - 7: **end while**
 - 8: $\hat{\mathbf{w}}^{(i)} := \hat{\mathbf{w}}_r^{(i)}$ for all nodes $i \in \mathcal{V}$
-

- linear map $\Phi(\mathbf{x}) = \mathbf{W}\mathbf{x}$ [1, Sec. 9.5] with a matrix \mathbf{W} that is determined by the correlations between private and non-private features. Figure 22 illustrates a toy dataset for which we can find a transformation that perfectly separates the private (gender) from the non-private feature which is used to predict the food preference of a person. One challenge for the use of learnt feature maps is that they need to be shared among all nodes. Indeed, we need use the same features for computing the GTV (76) (or its non-parametric extension (81)).
- the feature map $\Phi(\mathbf{x})$ is piecewise constant such that each region contains at least a prescribed minimum number k of data points from $\mathcal{D}^{(i)}$. The resulting perturbed feature vectors $\mathbf{z}^{(r)} = \Phi(\mathbf{x}^{(r)})$ are said to provide k -anonymity [93, 94]. Figure 23 illustrates a piece-wise constant feature map that quantizes a feature that would uniquely identify a data point.

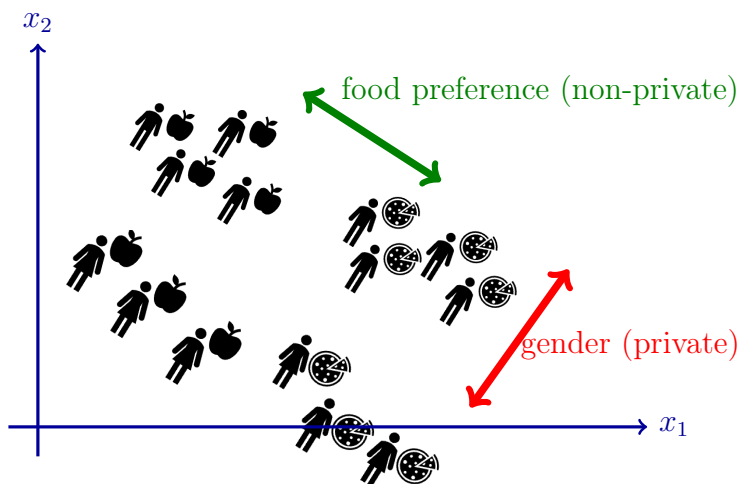


Figure 22: A toy dataset $\mathcal{D}^{(i)}$ whose data points represent persons that are characterized by two features x_1, x_2 . These features carry sensitive information (gender) and non-sensitive information (food preference) about a person. The non-sensitive information allows to predicting the quantity of interest (e.g., if a person likes pizza).

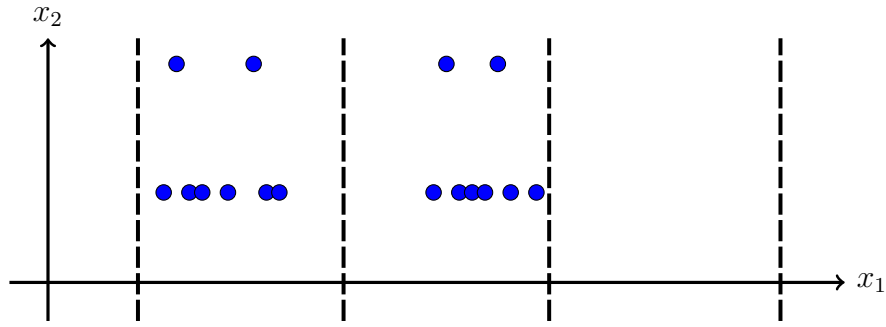


Figure 23: A toy dataset $\mathcal{D}^{(i)}$ whose data points are characterized by two features: x_1 and x_2 . The feature x_2 is private and carries sensitive information. The other feature x_1 is non-private but uniquely identifies a data point: two different data points in $\mathcal{D}^{(i)}$ have different values of x_1 . One simple approach to protect the private property x_1 is to quantize the feature x_1 resulting in a perturbed feature \tilde{x}_1 that is the same for at least k data points, thereby ensuring k -anonymity.

11.3 Exercises

Exercise 11.1. Privacy-Preserving Feature Learning (10 points)

Consider the toy dataset https://scikit-learn.org/stable/datasets/toy_dataset.html#diabetes-dataset whose data points consider patients. Each data point is characterized by $d = 10$ features $\mathbf{x} = (x_1, \dots, x_d)^T$ and a label y which measures a disease progression. The private feature x_1 characterizes the age of a person. We could simply remove this feature to obtain a privacy-preserving feature vector $\mathbf{x}'' = (x_2, \dots, x_d)^T \in \mathbb{R}^d$. This exercise explores how much better we could do by learning a privacy-preserving feature map Φ . Construct a feature map such that the transformed features do not allow to predict the age of the person. Choose an ML method from `scikit-learn` and compare the achievable performance when using the transformed features $\mathbf{x}' = \Phi(\mathbf{x})$ with using the features \mathbf{x}'' .

12 Data Poisoning

The FL algorithms of Section 9 pool the information contained in decentralized local datasets to train local (tailored) models (see Section 6.2). The benefit of information sharing, allowing to learn more accurate hypotheses, comes at the cost of increased vulnerability against data poisoning.

Data poisoning refers to the intentional manipulation (or fabrication) of local datasets to steer the training of a specific local model [95, 96]. We discuss two main flavours of data poisoning, known as a denial-of-service attack and a backdoor attack. Section 12.1 discusses how a denial-of-service attack manipulates local model training such that it performs poorly on its own local dataset [97]. Section 12.2 explains how backdoor attacks maliciously steer the training of a local model such that it performs well in general but anomalously for data points with specific features [98].

12.1 Denial-of-Service Attacks

Consider decentralized data that is represented by an empirical graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ whose nodes $i \in \mathcal{V}$ carry local datasets $\mathcal{D}^{(i)}$. We solve GTVMin to learn a hypothesis $h^{(i)}$ for some node $i \in \mathcal{V}$. Assume that some adversary controls the local datasets at a set $\mathcal{A} \subseteq \mathcal{V}$ of nodes.

12.2 Backdoor Attack

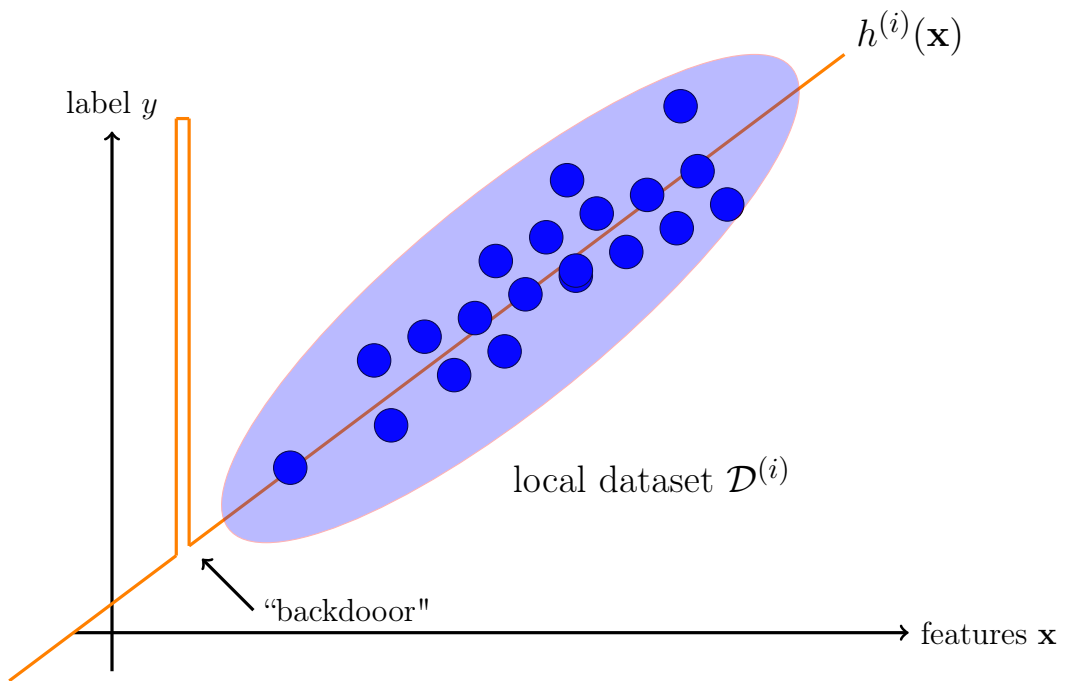


Figure 24: A backdoor attack aims at nudging the local hypothesis $h^{(i)}$ to behave in a specific way for certain feature values. These feature values can be interpreted as the key that unlocks a “backdoor”.

12.3 Exercises

Exercise 12.1. Backdoor Attack (10 points) Consider the toy dataset https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_digits.html#sklearn.datasets.load_digits whose data points are images of hand-written digits. We use the first $m' = 100$ data points as a validation set. Construct a training set for a decision tree classifier such that the trained classifier achieves an accuracy at least 0.9 on the validation set and classifies the 104th data point as any prescribed digit.

Part IV

Appendix

13 Linear Algebra and Convex Analysis

The main data structure of the FL methods discussed in part II are numeric arrays such as vectors or matrices [99]. Numeric arrays can be combined via algebraic operations such as (element-wise) addition $\mathbf{x} + \mathbf{y}$ or multiplication $\mathbf{Z}\mathbf{x}$ of compatible arrays (e.g., $\mathbf{Z} \in \mathbb{R}^{10 \times 3}$ and $\mathbf{x} \in \mathbb{R}^3$). These operations are the most important building blocks for most FL methods discussed in part II.

13.1 Convex Sets

Linear spaces are somewhat impractical as they are infinitely large. Using finite computational resources it would be useful to have a modification of linear spaces with finite extension. It turns out that convex sets pretty much fit this bill. These sets have many (algebraic and topological) properties in common with linear spaces. Moreover, many useful convex sets have a finite extension.

13.2 Convex Functions

Consider some real-valued function f whose domain is a subset of \mathbb{R}^d . The function is fully specified by its epigraph. We can then define convex functions as those functions whose epigraph is a convex set [1].

13.3 Proximal Operators

Given a closed proper convex function $f(\mathbf{x})$ with domain being a subset of \mathbb{R}^d , we define its associated convex conjugate function as [25]

$$f^*(\mathbf{x}) := \sup_{\mathbf{z} \in \mathbb{R}^d} \mathbf{x}^T \mathbf{z} - f(\mathbf{z}). \quad (98)$$

The proximity operator associated with a closed proper convex function $f(\mathbf{x})$ is defined as [100]

$$\mathbf{prox}_{f,\rho}(\mathbf{x}) := \underset{\mathbf{x}'}{\operatorname{argmin}} f(\mathbf{x}') + (\rho/2)\|\mathbf{x} - \mathbf{x}'\|_2^2 \text{ for some } \rho > 0. \quad (99)$$

Note that the minimum in (99) exists and is unique since the objective function is strongly convex [25] .

14 Information Theory

It is often useful to interpret data points as realizations of RVs. The behaviour of ML algorithms that use these data points can then be studied using the relations between these RVs. A powerful tool to study relations between RVs is the concept of mutual information [101, 102]. The mutual information between two RVs \mathbf{x}, y with joint probability distribution $p(\mathbf{x}, y)$ is defined as

$$I(\mathbf{x}, y) := \mathbb{E} \{ \log p(\mathbf{x}, y) / (p(\mathbf{x})p(y)) \} \quad (100)$$

15 Fixed Point Theory

Glossary

0/1 loss The 0/1 loss $L((\mathbf{x}, y), h)$ measures the quality of a classifier $h(\mathbf{x})$ that delivers a prediction \hat{y} (e.g., via thresholding (101)) for a label y of a data point. It is equal to 0 if the prediction is correct, i.e., $L((\mathbf{x}, y), h) = 0$ when $\hat{y} = y$. It is equal to 1 if the prediction is wrong, $L((\mathbf{x}, y), h) = 1$ when $\hat{y} \neq y$. 35, 36, 41–46, 147, 168

k -means The k -means algorithm is a hard clustering method which assigns each data points to pecisely one out of k different clusters. The method iteratively updates this assignment in order to minimize the average distance between data points in their nearest cluster mean (centre). 150, 157, 176

accuracy Consider data points characterized by features $\mathbf{x} \in \mathbf{X}$ and a categorical label y which takes on values from a finite label space \mathcal{Y} . The accuracy of a hypothesis $h : \mathcal{X} \rightarrow \mathcal{Y}$, when applied to the data points in a dataset $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$ is then defined as $1 - (1/m) \sum_{r=1}^m L((\mathbf{x}^{(r)}, y^{(r)}), h)$ using the 0/1 loss. 35, 42

activation function Each artificial neuron within an ANN consists of an activation function that maps the inputs of the neuron to a single output value. In general, an activation function is a non-linear map of the weighted sum of neuron inputs (this weighted sum is the activation of the neuron). 15

artificial intelligence Artificial intelligence aims at developing systems that behave rational in the sense of maximizing a long-term reward. 22, 129

artificial neural network An artificial neural network is a graphical (signal-flow) representation of a map from features of a data point at its input to a predicted label at its output. 15, 29–32, 57, 64, 73, 93, 115, 116, 120, 147, 154, 160, 170, 177

backdoor A backdoor (attack) is obtained by data poisoning such that a trained local model predicts well most data points but anomalously for specific feature values (which “unlock” the backdoor). 141

batch A set of randomly selected data points out of a (typically very large) dataset. 89

Bayes estimator A hypothesis h whose Bayes risk is minimal [39]. 49, 50, 52, 148

Bayes risk We use the term Bayes risk as a synonym for the risk or expected loss of a hypothesis. Some authors reserve the term Bayes risk for the risk of a hypothesis that achieves minimum risk, such a hypothesis being referred to as a Bayes estimator [39]. 148

bias Consider some unknown quantity \bar{w} , e.g., the true weight in a linear model $y = \bar{w}x + e$ relating feature and label of a data point. We might use an ML method (e.g., based on ERM) to compute an estimate \hat{w} for the \bar{w} based on a set of data points that are realizations of RVs. The (squared) bias incurred by the estimate \hat{w} is typically defined as $B^2 := (\mathbb{E}\{\hat{w}\} - \bar{w})^2$. We extend this definition to vector-valued quantities using the squared Euclidean norm $B^2 := \|\mathbb{E}\{\hat{\mathbf{w}}\} - \bar{\mathbf{w}}\|_2^2$. 14, 55, 56, 82

classification Classification is the task of determining a discrete-valued label y of a data point based solely on its features \mathbf{x} . The label y belongs to a finite set, such as $y \in \{-1, 1\}$, or $y \in \{1, \dots, 19\}$ and represents a category to which the corresponding data point belongs to. 26, 36, 39, 42–45, 47, 111, 116

classifier A classifier is a hypothesis (map) $h(\mathbf{x})$ that is used to predict a discrete-valued label. Strictly speaking, a classifier is a hypothesis $h(\mathbf{x})$ that can take only a finite number of different values. However, we are sometimes sloppy and use the term classifier also for a hypothesis that delivers a real number which is quantized, i.e., compared against a threshold, to obtain the predicted label value. For example, in a binary classification problem with label values $y \in \{-1, 1\}$, we refer to a hypothesis $h(\mathbf{x})$ as classifier if it is used to compute a predicted label according to

$$\hat{y} = 1 \text{ for } h(\mathbf{x}) \geq 0, \text{ and } \hat{y} = -1 \text{ otherwise.} \quad (101)$$

35, 36, 39, 40, 45, 46, 143, 147

cluster A cluster within an empirical graph is a subset of nodes that carry local datasets with (nearly) identical statistical distributions. Clustered FL aims at identifying the clusters and use the pooled local datasets to train a separate model for each cluster. 99

cluster A cluster is a subset of data points that are more similar to each other than to the data points outside the cluster. The notion and measure of similarity between data points is a design choice. If data

points are characterized by numeric Euclidean feature vectors it might be reasonable to define similarity between two data points using the (inverse of the) Euclidean distance between the corresponding feature vectors 94, 99, 109, 113, 115, 116

clustering Clustering methods decompose a given set of data points into few subsets, which are referred to as clusters. Each cluster consists of data points that are more similar to each other than to data points outside the cluster. Different clustering methods use different measures for the similarity between data points and different representation of clusters. The clustering method k -means uses the average feature vector (“cluster means”) of a cluster as its representative. A popular soft-clustering method based on Gaussian mixture model (GMM) represents a cluster by a multivariate normal distribution. 108, 109, 147

clustering assumption The clustering assumption postulates that data points form a (small) number of groups or clusters, within which the data points behave similar [83]. 93, 107, 115

computational aspects By computational aspects of a ML method, we mainly refer to the computational resources required for its implementation. For example, if a ML method uses iterative optimization techniques to solve ERM, then its computational aspects include (i) how many arithmetic operations are needed to implement a single iteration (gradient step) and (ii) how many iterations are needed to obtain useful model parameters. One important example for an iterative optimization technique is GD. 25, 27, 29, 30, 34–36, 51, 52, 82, 95, 97, 100,

106, 111, 165, 179

convex A set $\mathcal{C} \subseteq \mathbb{R}^d$ is convex if it contains the line segment between any two points of that set. We define a function as convex if its epigraph is a convex set [25]. 35, 36, 38, 42–44, 106, 172, 173

covariance matrix The covariance matrix of a RV $\mathbf{x} \in \mathbb{R}^d$ is defined as $\mathbb{E}\left\{(\mathbf{x} - \mathbb{E}\{\mathbf{x}\})(\mathbf{x} - \mathbb{E}\{\mathbf{x}\})^T\right\}$. 11, 67, 81, 102, 157, 168

data A (indexed) set of data points. 18, 21, 25, 61, 172

data augmentation Data augmentation methods add synthetic data points to an existing set of data points. These synthetic data points might be obtained by perturbations (adding noise) or transformations (rotations of images) of the original data points. 65–69, 71

data point A data point is any object that conveys information [101]. Data points might be students, radio signals, trees, forests, images, RVs, real numbers or proteins. We characterize data points using two types of properties. One type of property is referred to as a feature. Features are properties of a data point that can be measured or computed in an automated fashion. Another type of property is referred to as labels. The label of a data point represents a higher-level facts or quantities of interest. In contrast to features, determining the label of a data point typically requires human experts (domain experts). Roughly speaking, ML aims at predicting the label of a data point based solely on its features. 7, 11–14, 16, 17, 25–29, 31, 33–35, 37–42, 45–56, 58–68, 71–73,

79, 80, 85, 86, 88, 89, 94, 98, 99, 101, 102, 104, 109, 113, 116–118, 131–133, 135, 137–141, 143, 146–157, 159–167, 169–171, 174–180

data poisoning FL methods allow to leverage the information contained in local datasets generated by other parties to improve the training of a tailored model. Depending on how much we trust the other parties, FL can be compromised by data poisoning. Data poisoning refers to the intentional manipulation (or fabrication) of local datasets to steer the training of a specific local model [95,96]. 19, 141, 148, 153

dataset With a slight abuse of notation we use the terms “dataset” or “set of data points” to refer to an indexed list of data points $\mathbf{z}^{(1)}, \dots$. Thus, there is a first data point $\mathbf{z}^{(1)}$, a second data point $\mathbf{z}^{(2)}$ and so on. Strictly speaking a dataset is a list and not a set [103]. By using indexed lists of data points we avoid some of the challenges arising in concept of an abstract set. 12, 16, 21, 26, 39, 46, 48, 50, 51, 58, 61, 62, 67, 68, 71, 94, 113, 115, 117, 118, 135, 136, 138–141, 143, 147, 148, 154, 156, 162, 165

decision region Consider a hypothesis map h that reads in a feature vector $\mathbf{x} \in \mathbb{R}^d$ and delivers a value from a finite set \mathcal{Y} . The decision boundary induced by h is the set of vectors $\mathbf{x} \in \mathbb{R}^d$ that lie between different decision regions. More precisely, a vector \mathbf{x} belongs to the decision boundary if and only if each neighbourhood $\{\mathbf{x}' : \|\mathbf{x} - \mathbf{x}'\| \leq \varepsilon\}$, for any $\varepsilon > 0$, contains at least two vectors with different function values. 160

decision region Consider a hypothesis map h that can only take values from a finite set \mathcal{Y} . We refer to the set of features $\mathbf{x} \in \mathcal{X}$ that result in

the same output $h(\mathbf{x}) = a$ as a decision region of the hypothesis h . 15, 152, 164

decision tree A decision tree is a flow-chart like representation of a hypothesis map h . More formally, a decision tree is a directed graph which reads in the feature vector \mathbf{x} of a data point at its root node. The root node then forwards the data point to one of its children nodes based on some elementary test on the features \mathbf{x} . If the receiving children node is not a leaf node, i.e., it has itself children nodes, it represents another test. Based on the test result, the data point is further pushed to one of its neighbours. This testing and forwarding of the data point is repeated until the data point ends up in a leaf node (having no children nodes). The leaf nodes represent sets (decision regions) constituted by feature vectors \mathbf{x} that are mapped to the same function value $h(\mathbf{x})$. 15, 20, 47, 96, 105, 110, 112, 116, 143

denial-of-service attack A denial-of-service attack uses data poisoning to steer the training of a local (client) model such that it performs poorly for typical data points 141

differentiable A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is differentiable if it has a gradient $\nabla f(\mathbf{x})$ everywhere (for every $\mathbf{x} \in \mathbb{R}^d$) [5]. 15, 35, 36, 38, 42–44, 73–75, 91, 97, 122, 158, 159, 176, 177

dimensionality reduction Dimensionality reduction methods map (typically many) raw features to a (relatively small) set of new features. These methods can be used to visualize data points by learning two features that can be used as the coordinates of a depiction in a scatterplot.

effective dimension The effective dimension $d_{\text{eff}}(\mathcal{H})$ of an infinite hypothesis space \mathcal{H} is a measure of its size. Loosely speaking, the effective dimension is equal to the number of “independent” tunable parameters of the model. These parameters might be the coefficients used in a linear map or the weights and bias terms of an ANN. 55, 57, 63–65, 159

eigenvalue We refer to a number $\lambda \in \mathbb{R}$ as eigenvalue of a square matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ if there is a non-zero vector $\mathbf{x} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ such that $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$. 15, 52, 78, 80–82

empirical graph Empirical graphs represent collections of local datasets and corresponding local models [83]. An empirical graph is an undirected weighted empirical graph whose nodes carry local datasets and models. FL methods learn a local hypothesis $h^{(i)}$, for each node $i \in \mathcal{V}$, such that it incurs small loss on the local datasets. 16, 21, 69, 93–96, 99–105, 108–110, 113, 114, 141, 149, 165, 169

empirical risk The empirical risk of a given hypothesis on a given set of data points is the average loss of the hypothesis computed over all data points in that set. 14, 50–54, 61, 67, 69, 70, 73, 74, 105, 154, 169, 170, 177–179

empirical risk minimization Empirical risk minimization is the optimization problem of finding the hypothesis with minimum average loss (or empirical risk) on a given set of data points (the training set). Many ML methods are special cases of empirical risk. 21, 48, 51–54, 56, 58,

59, 61, 63–65, 68, 71, 73, 74, 79–81, 83–86, 90, 93, 135, 148, 150, 159, 168, 170, 174, 178, 179

Euclidean space The Euclidean space \mathbb{R}^d of dimension d refers to the space of all vectors $\mathbf{x} = (x_1, \dots, x_d)$, with real-valued entries $x_1, \dots, x_d \in \mathbb{R}$, whose geometry is defined by the inner product $\mathbf{x}^T \mathbf{x}' = \sum_{j=1}^d x_j x'_j$ between any two vectors $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ [5]. 97, 156, 164, 173

expectation Consider a numeric RV \mathbf{x} with probability distribution $p(\mathbf{x})$. The expectation of \mathbf{x} is defined as the integral $\mathbb{E}\{\mathbf{x}\} := \int \mathbf{x}p(\mathbf{x})$ (if this integral is well-defined!) [5, 35, 104]. 55

explainability We can define the (subjective) explainability of a ML method as the level of predictability (or lack of uncertainty) in its predictions delivered to a specific human user. Quantitative measures for the (subjective) explainability can be obtained via probabilistic models for the data fed into the ML method. In particular, we can then identify explainability of a ML method via the uncertainty or entropy of the predictions it delivers [32, 105, 106]. 28, 29, 35, 36, 162

feature A feature of a data point is one of its properties that can be measured or computed in an automated fashion. For example, if a data point is a bitmap image, then we could use the red-green-blue intensities of its pixels as features. Some widely used synonyms for the term feature are “covariate”, “explanatory variable”, “independent variable”, “input (variable)”, “predictor (variable)” or “regressor” [107–109]. However, this book makes consequent use of the term features for low-level properties

of data points that can be measured easily. 7, 11–17, 25–32, 34, 37, 38, 41, 47, 50, 52, 54–56, 58, 59, 63, 64, 66, 68, 71–73, 78–81, 84, 86, 93, 94, 109, 113, 116–118, 132, 133, 135, 137, 138, 140, 141, 147, 148, 151, 153, 156, 157, 159–161, 164–167, 170, 171, 174–176, 178, 180

feature map A map that transforms the original features of a data point into new features. The so-obtained new features might be preferable over the original features for several reasons. For example, the shape of datasets might become simpler in the new feature space, allowing to use linear models in the new features. Another reason could be that the number of new features is much smaller which is preferable in terms of avoiding overfitting. If the feature map delivers two new features, we can depict data points in a scatterplot. 15, 31, 32, 47, 80, 81, 135–137, 140

feature matrix Consider a dataset \mathcal{D} of m data points that are characterized by feature vectors $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$. It is convenient to collect these feature vectors into a feature matrix $\mathbf{X} := (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})$. 78

feature space The feature space of a given ML application or method is constituted by all potential values that the feature vector of a data point can take on. Within this book the most frequently used choice for the feature space is the Euclidean space \mathbb{R}^d with dimension d being the number of individual features of a data point. 15, 16, 26, 30, 161, 162

federated averaging (FedAvg) Federated averaging is an iterative FL algorithm that alternates between local model trainings and averaging the resulting local model parameters. Different variants of this algorithm

are obtained by different techniques for the model training. The authors of [14] consider federated averaging methods where the local model training is implemented by running several GD steps 123–125

federated learning (FL) Federated learning is an umbrella term for ML methods that train models in a collaborative fashion using decentralized data and computation. 1, 2, 18–22, 26, 29, 61, 68–70, 83, 94–98, 100, 101, 103, 105, 106, 108, 113–119, 123, 129, 131, 141, 145, 149, 152, 154, 156, 161, 180

Finnish Meteorological Institute The Finnish Meteorological Institute is a government agency responsible for gathering and reporting weather data in Finland. 104

flow-based clustering Flow-based clustering groups the nodes of an undirected graph by applying k -means clustering to node-wise feature vectors. These feature vectors are built from flows between carefully selected source and destination nodes [76]. 108

Gaussian mixture model Gaussian mixture models (GMM) are a family of probabilistic models for data points characterized by a numeric feature vector \mathbf{x} . A GMM interprets \mathbf{x} as being drawn from one out of k different multivariate normal distributions $p^{(c)} = \mathcal{N}(\boldsymbol{\mu}^{(c)}, \mathbf{C}^{(c)})$, indexed by $c = 1, \dots, k$. The probability that \mathbf{x} is drawn from the c -th multivariate normal distribution is denoted p_c . Thus, a GMM is parametrized by the probability p_c , the mean vector $\boldsymbol{\mu}^{(c)}$ and covariance matrix $\boldsymbol{\Sigma}^{(c)}$ for each $c = 1, \dots, k$. 150

Gaussian random variable A Gaussian RV $\mathbf{x} \in \mathbb{R}^d$ with a multivariate normal distribution. The special case of d corresponds to a scalar Gaussian RV [4, 41, 55]. 67

General Data Protection Regulation The General Data Protection Regulation (GDPR) is a law that has been passed by the European Union (EU) and put into effect on May 25, 2018 <https://gdpr.eu/tag/gdpr/>. The GDPR imposes obligations onto organizations anywhere, so long as they target, collect or in any other way process data related to people (i.e., personal data) in the EU. 20

generalized total variation Generalized total variation measures the changes of vector-valued node attributes of a graph. 100, 103, 105–110, 115, 121–123, 134, 136, 137, 159

gradient For a real-valued function $f : \mathbb{R}^d \rightarrow \mathbb{R} : \mathbf{w} \mapsto f(\mathbf{w})$, a vector \mathbf{a} such that $\lim_{\mathbf{w} \rightarrow \mathbf{w}'} \frac{f(\mathbf{w}) - (f(\mathbf{w}') + \mathbf{a}^T(\mathbf{w} - \mathbf{w}'))}{\|\mathbf{w} - \mathbf{w}'\|} = 0$ is referred to as the gradient of f at \mathbf{w}' . If such a vector exists it is denoted $\nabla f(\mathbf{w}')$ or $\nabla f(\mathbf{w})|_{\mathbf{w}'}$. 15, 18, 19, 44, 73–75, 77, 78, 82, 87–89, 91, 97, 153, 158, 159, 176–178

gradient descent (GD) Gradient descent is an iterative method for finding the minimum of a differentiable function $f(\mathbf{w})$. 21, 44, 52, 73, 74, 76–78, 81–84, 86, 90, 97, 123, 150, 157, 163, 164, 172, 177, 178

gradient step Given a differentiable real-valued function $f(\mathbf{w})$ and a vector \mathbf{w}' , the gradient step updates \mathbf{w}' by adding the scaled negative gradient $\nabla f(\mathbf{w}')$, $\mathbf{w}' \mapsto \mathbf{w}' - \alpha \nabla f(\mathbf{w}')$. 73–77, 79–89, 91, 114, 150, 172

gradient-based method Gradient-based methods are iterative algorithms for finding the minimum (or maximum) of a differentiable objective function of the model parameters. These algorithms construct a sequence of approximations to an optimal choice for model parameters that results in a minimum objective function value. As their name indicates, gradient-based methods use the gradients of the objective function evaluated during previous iterations to construct new (hopefully) improved model parameters. 13, 34, 36, 44, 73, 98

graph A graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is a pair that consists of a node set \mathcal{V} and an edge set \mathcal{E} . In general, a graph is specified by a map that assigns to each edge $e \in \mathcal{E}$ a pair of nodes [110]. One important family of graphs (simple undirected graphs) is obtained by identifying each edge $e \in \mathcal{E}$ with two different nodes $\{i, i'\}$. Weighted graphs also specify numeric weights A_e for each edge $e \in \mathcal{E}$. 103, 154

GTV minimization GTV minimization is an instance of RERM using the GTV of local parameter vectors as regularization term. 22, 96, 100, 105, 107–115, 119, 120, 122, 141

high-dimensional regime The high-dimensional regime of ERM is characterized by the effective dimension of the model being larger than the sample size, i.e., the number of (labeled) data points in the training set. For example, linear regression methods operate in the high-dimensional regime whenever the number d of features used to characterize data points exceeds the number of data points in the training set. Another example for ML methods that operate in the high-dimensional

regime is deep learning with large ANNs, having more tunable weights (and bias terms) than the number of data points in the training set. High-dimensional statistics is a recent main thread of probability theory that studies the behavior of ML methods in the high-dimensional regime [42, 111]. 98

hinge loss Consider a data point that is characterized by a feature vector $\mathbf{x} \in \mathbb{R}^d$ and a binary label $y \in \{-1, 1\}$. The hinge loss incurred by a specific hypothesis h is defined as

$$L((\mathbf{x}, y), h) := \max\{0, 1 - yh(\mathbf{x})\}. \quad (102)$$

A regularized variant of the hinge loss is used by the support vector machine (SVM) [112] to learn a linear classifier with maximum margin between the two classes (see Figure 25). 42–44, 160, 178

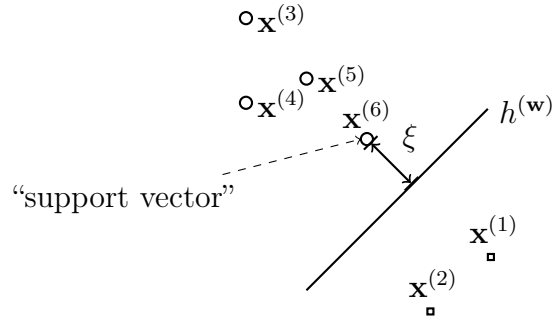


Figure 25: The SVM learns a hypothesis (or classifier) $h^{(w)}$ with minimum average soft-margin hinge loss. Minimizing this loss is equivalent to maximizing the margin ξ between the decision boundary of $h^{(w)}$ and each class of the training set.

horizontal FL Horizontal FL refers to applications with local datasets that are constituted by different data points but using the same features to characterize them [?]. One example for horizontal FL is numerical weather prediction using a network of weather (observation) stations. Each weather station measures the same parameters such as daily temperature, air pressure and precipitation. However, different weather stations measure the characteristics or features of different spatio-temporal regions (each such region being a separate data point). 113, 117, 118

Huber loss The Huber loss is a mixture of the squared error loss and the absolute value of the prediction error. 37, 38

hypothesis A map (or function) $h : \mathcal{X} \rightarrow \mathcal{Y}$ from the feature space \mathcal{X} to the label space \mathcal{Y} . Given a data point with features \mathbf{x} we use a hypothesis map h to estimate (or approximate) the label y using the predicted label $\hat{y} = h(\mathbf{x})$. ML is about learning (or finding) a hypothesis map h such that $y \approx h(\mathbf{x})$ for any data point. 7, 13, 14, 17, 25, 28–30, 33–42, 44, 46–50, 52–56, 58–67, 72, 73, 81–83, 96–99, 110, 111, 116, 122, 131, 133, 141, 147–149, 152, 154, 160, 163, 165–171, 174–179

hypothesis space Every practical ML method uses a specific hypothesis space (or model) \mathcal{H} . The hypothesis space of a ML method is a subset of all possible maps from the feature space to label space. The design choice of the hypothesis space should take into account available computational resources and statistical aspects. If the computational infrastructure allows for efficient matrix operations and we expect a linear relation between features and label, a reasonable first candidate for the hypothesis

space is a linear model. 13, 17, 25, 28–34, 46–48, 51–54, 58, 59, 62, 63, 66, 73, 93, 96, 110, 111, 154, 161, 164–166, 168–170, 174, 178, 180

i.i.d. It can be useful to interpret data points $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}$ as realizations of independent and identically distributed RVs with a common probability distribution. If these RVs are continuous, their joint probability density function (pdf) is $p(\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}) = \prod_{r=1}^m p(\mathbf{z}^{(r)})$ with $p(\mathbf{z})$ being the common marginal pdf of the underlying RVs. 19, 42, 48, 50, 54, 55, 59, 62, 63, 67, 72, 81, 87, 101, 133, 162, 163, 167, 169, 175, 176

i.i.d. assumption The i.i.d. assumption interprets data points of a dataset as the realizations of i.i.d. RVs. 48, 49, 51, 52, 54, 89, 99, 100, 169, 170, 175

interpretability A synonym for explainability [31]. 30, 165

kernel Consider data points characterized by features $\mathbf{x} \in \mathcal{X}$ with values in some arbitrary feature space. A kernel is a map that assigns each pair of features $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ a real number. This number is interpreted as a measure for the similarity between \mathbf{x} and \mathbf{x}' . For more details about kernels and the resulting kernel methods, we refer to the literature [112, 113]. 31

Kullback-Leibler divergence The Kullback–Leibler divergence is a quantitative measure for how much one probability distribution is different from another probability distribution [101]. 101

label A higher level fact or quantity of interest associated with a data point. If a data point is an image, its label might be the fact that it shows a cat

(or not). Some widely used synonyms for the term label are "response variable", "output variable" or "target" [107–109]. 7, 11, 13, 25–29, 34, 37, 38, 40, 41, 50–56, 59, 64, 66, 68, 71, 72, 75, 86, 94, 117, 132, 135, 140, 147, 149, 151, 160, 161, 163–167, 170, 171, 174, 176, 178

label space Consider a ML application that involves data points characterized by features and labels. The label space is constituted by all potential values that the label of a data point can take on. Regression methods, aiming at predicting numeric labels, often use the label space $\mathcal{Y} = \mathbb{R}$. Binary classification methods use a label space that consists of two different elements, e.g., $\mathcal{Y} = \{-1, 1\}$, $\mathcal{Y} = \{0, 1\}$ or $\mathcal{Y} = \{\text{"cat image"}, \text{"no cat image"}\}$ 26, 30, 39, 41, 45–47, 147, 161, 166

law of large numbers The law of large numbers refers to the convergence of the average of an increasing (large) number of i.i.d. RVs to the mean (or expectation) of their common probability distribution. Different instances of the law of large numbers are obtained using different notions of convergence. 42, 50, 51

learning rate Consider an iterative method for finding or learning a good choice for a hypothesis. Such an iterative method repeats similar computational (update) steps that adjust or modify the current choice for the hypothesis to obtain an improved hypothesis. A prime example for such an iterative learning method is GD and its variants. We refer by learning rate to any parameter of an iterative learning method that controls the extent by which the current hypothesis might be modified or improved in each iteration. A prime example for such a parameter

is the step size used in GD. Some authors use the term learning rate mostly as a synonym for the step size of (a variant of) GD 13, 76–80, 84, 85, 88, 90

learning task A learning task consists of a specific choice for a collection of data points (e.g., all images stored in a particular database), their features and labels. 14, 26, 69, 115

least absolute shrinkage and selection operator (Lasso) The least absolute shrinkage and selection operator (Lasso) is an instance of structural risk minimization (SRM) for learning the weights \mathbf{w} of a linear map $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$. The Lasso minimizes the sum consisting of an average squared error loss (as in linear regression) and the scaled ℓ_1 norm of the weight vector \mathbf{w} . 66, 98, 174

linear classifier A classifier $h(\mathbf{x})$ maps the feature vector $\mathbf{x} \in \mathbb{R}^d$ of a data point to a predicted label $\hat{y} \in \mathcal{Y}$ out of a finite set of label values \mathcal{Y} . We can characterize such a classifier equivalently by the decision regions \mathcal{R}_a , for every possible label value $a \in \mathcal{Y}$. Linear classifiers are such that the boundaries between the regions \mathcal{R}_a are hyperplanes in the Euclidean space \mathbb{R}^d . 160

linear model This book uses the term linear model in a very specific sense. In particular, a linear model is a hypothesis space which consists of all linear maps,

$$\mathcal{H}^{(d)} := \{h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} : \mathbf{w} \in \mathbb{R}^d\}. \quad (103)$$

Note that (103) defines an entire family of hypothesis spaces, which is parametrized by the number d of features that are linearly combined

to form the prediction $h(\mathbf{x})$. The design choice of d is guided by computational aspects (smaller d means less computation), statistical aspects (increasing d might reduce prediction error) and interpretability (a linear model using few carefully chosen features might be considered interpretable). 29–32, 34, 52, 54–57, 63, 68, 71, 72, 93, 96–98, 123, 148, 162

linear regression Linear regression aims at learning a linear hypothesis map to predict a numeric label based on numeric features of a data point. The quality of a linear hypothesis map is measured using the average squared error loss incurred on a set of labeled data points (which we refer to as training set). 52, 58, 67, 71, 73, 74, 81, 82, 84, 86, 98, 120, 123, 125, 159, 164

local dataset The concept of a local dataset is in-between the concept of a data point and a dataset. Similar to the basic notion of a dataset, also a local dataset consists of several individual data points which are characterized by features and label. In contrast to a single dataset used in basic ML methods, a local dataset is also related to other local datasets via different notions of similarities. These similarities might arising from probabilistic models or communication infrastructure and are encoded in the edges of an empirical graph. 1, 2, 16, 18–20, 22, 69, 93–105, 108, 113–119, 123, 131–133, 135, 141, 152, 154, 161, 165, 169, 180

local model Consider a collections of local datasets that are assigned to the nodes of an empirical graph. A local model $\mathcal{H}^{(i)}$ is a hypothesis space

that is assigned to a node $i \in \mathcal{V}$. Different nodes might be assigned different hypothesis spaces, i.e., in general $\mathcal{H}^{(i)} \neq \mathcal{H}^{(i')}$ for different nodes $i, i' \in \mathcal{V}$. 1, 2, 16, 19, 20, 61, 68, 83, 93, 94, 96, 97, 99, 100, 105–113, 115, 116, 119–122, 131, 152, 154, 169

logistic loss Consider a data point that is characterized by the features \mathbf{x} and a binary label $y \in \{-1, 1\}$. We use a real-valued hypothesis h to predict the label y solely from the features \mathbf{x} . The logistic loss incurred by a specific hypothesis h is defined as

$$L((\mathbf{x}, y), h) := \log(1 + \exp(-yh(\mathbf{x}))). \quad (104)$$

. Carefully note that the expression (104) for the logistic loss applies only for the label space $\mathcal{Y} = \{-1, 1\}$ and using the thresholding rule (101). 35, 36, 44–46, 73, 166, 168

logistic regression Logistic regression learns a linear hypothesis map (classifier) $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ to predict a binary label y based on numeric features \mathbf{x} of a data point. The quality of a linear hypothesis map is measured using its average logistic loss on some labeled data points (the training set). 44

loss With a slight abuse of language, we use the term loss either for loss function itself or for its value for a specific pair of data point and hypothesis. 14, 17, 33–38, 41–46, 48–50, 61, 66, 67, 69, 81, 96, 98, 99, 107, 116, 121, 134, 147, 148, 154, 160, 167, 168, 170, 174–178

loss function A loss function is a map

$$L : \mathcal{X} \times \mathcal{Y} \times \mathcal{H} \rightarrow \mathbb{R}_+ : ((\mathbf{x}, y), h) \mapsto L((\mathbf{x}, y), h)$$

which assigns a pair consisting of a data point, with features \mathbf{x} and label y , and a hypothesis $h \in \mathcal{H}$ the non-negative real number $L((\mathbf{x}, y), h)$. The loss value $L((\mathbf{x}, y), h)$ quantifies the discrepancy between the true label y and the predicted label $h(\mathbf{x})$. Smaller (closer to zero) values $L((\mathbf{x}, y), h)$ mean a smaller discrepancy between predicted label and true label of a data point. Figure 26 depicts a loss function for a given data point, with features \mathbf{x} and label y , as a function of the hypothesis $h \in \mathcal{H}$. 21, 25, 33–38, 40–45, 50–52, 58, 61, 64, 66, 67, 73, 96–100,

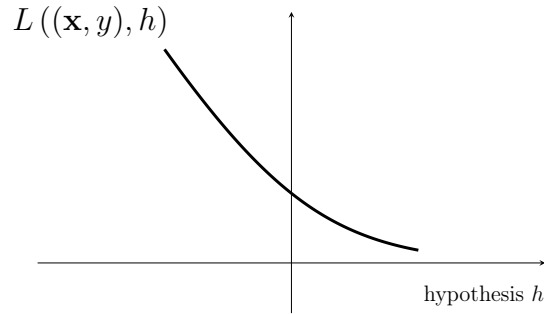


Figure 26: Some loss function $L((\mathbf{x}, y), h)$ for a fixed data point, with feature vector \mathbf{x} and label y , and varying hypothesis h . ML methods try to find (learn) a hypothesis that incurs minimum loss.

106–111, 113, 115, 116, 122, 135, 166–168, 175

maximum likelihood Consider data points $\mathcal{D} = \{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ that are interpreted as realizations of i.i.d. RVs with a common probability distribution $p(\mathbf{z}; \mathbf{w})$ which depends on a parameter vector $\mathbf{w} \in \mathcal{W} \subseteq \mathbb{R}^n$. Maximum likelihood methods aim at finding a parameter vector \mathbf{w} such that the probability (density) $p(\mathcal{D}; \mathbf{w}) = \prod_{r=1}^m p(\mathbf{z}^{(r)}; \mathbf{w})$ of observing the data is maximized. Thus, the maximum likelihood estimator is

obtained as a solution to the optimization problem $\max_{\mathbf{w} \in \mathcal{W}} p(\mathcal{D}; \mathbf{w})$.

35, 38

mean The expectation $\mathbb{E}\{\mathbf{x}\}$ of a numeric RV \mathbf{x} . 91

metric A metric refers to a loss function that is used solely for the final performance evaluation of a learnt hypothesis. The metric is typically a loss function that has a “natural” interpretation (such as the 0/1 loss) but is not a good choice to guide the learning process, e.g., via ERM. For ERM, we typically prefer loss functions that depend smoothly on the (parameters of the) hypothesis. Examples for such smooth loss functions include the squared error loss (105) and the logistic loss (104).

36

model We use the term model as a synonym for hypothesis space 13, 21, 25, 28, 30, 33, 34, 61, 63–67, 71, 73, 152, 154, 159, 161, 170

model parameters Model parameters are numbers that select a hypothesis map out of a hypothesis space. 34, 43, 69, 70, 73–75, 86, 87, 89, 93, 94, 97–100, 103, 105–109, 114–116, 119–121, 124, 125, 136, 150, 159, 172, 178

multivariate normal distribution The multivariate normal distribution $\mathcal{N}(\mathbf{m}, \mathbf{C})$ is an important family of probability distributions for a continuous RV $\mathbf{x} \in \mathbb{R}^d$ [?, 41, 114]. This family is parametrized by the mean \mathbf{m} and covariance matrix \mathbf{C} of \mathbf{x} . If the covariance matrix is invertible, the probability distribution of \mathbf{x} is

$$p(\mathbf{x}) \propto \exp \left(- (1/2)(\mathbf{x} - \mathbf{m})^T \mathbf{C}^{-1}(\mathbf{x} - \mathbf{m}) \right).$$

11, 54, 101, 150, 157, 158

networked data Networked data consists of local datasets that are related by some notion of pair-wise similarity. We represent networked data using an empirical graph whose nodes carry local datasets and an edge indicates a similarity between the connected nodes. 1, 18, 95, 108

networked federated learning Networked federated learning refers to methods that learn personalized models in a distributed fashion from local datasets that are related by an intrinsic network structure. 109

networked model A networked model over an empirical graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ assigns a local model (hypothesis space) to each node $i \in \mathcal{V}$ of the empirical graph \mathcal{G} . 96, 97, 105

non-smooth We refer to a function as non-smooth if it is not smooth [54]. 35

objective function An objective function is a map that assigns each possible value of an optimization variable, such as the parameters \mathbf{w} of a hypothesis $h^{(\mathbf{w})}$, to an objective value $f(\mathbf{w})$. The objective value $f(\mathbf{w})$ could be the risk or the empirical risk of a hypothesis $h^{(\mathbf{w})}$. 51, 52, 55, 73, 74, 76–78, 89, 90

outlier Many ML methods are motivated by the i.i.d. assumption which interprets data points as realizations of i.i.d. RVs with a common probability distribution. The i.i.d. assumption is useful for applications where the statistical properties of the data generation process are

stationary (time-invariant). However, in some applications the data consists of a majority of “regular” data points that conform with an i.i.d. assumption and a small number of data points that have fundamentally different statistical properties compared to the regular data points. We refer to a data point that substantially deviates from the statistical properties of the majority of data points as an outlier. Different methods for outlier detection use different measures for this deviation. 35, 38

overfitting Consider a ML method that uses ERM to learn a hypothesis with minimum empirical risk on a given training set. Such a method is “overfitting” the training set if it learns hypothesis with small empirical risk on the training set but unacceptably large loss outside the training set. 61, 63, 64, 83, 156

parameters The parameters of a ML model are tunable (learnable or adjustable) quantities that allow to choose between different hypothesis maps. For example, the linear model $\mathcal{H} := \{h : h(x) = w_1x + w_2\}$ consists of all hypothesis maps $h(x) = w_1x + w_2$ with a particular choice for the parameters w_1, w_2 . Another example of parameters are the weights assigned to the connections of an ANN. 73, 121, 132, 134

polynomial regression Polynomial regression aims at learning a polynomial hypothesis map to predict a numeric label based on numeric features of a data point. For data points characterized by a single numeric features, polynomial regression uses the hypothesis space $\mathcal{H}_d^{(\text{poly})} := \{h(x) = \sum_{j=0}^{d-1} x^j w_j\}$. The quality of a polynomial hypothesis map is measured using the average squared error loss incurred on a set

of labeled data points (which we refer to as training set). 63

positive semi-definite A symmetric matrix $\mathbf{Q} = \mathbf{Q}^T \in \mathbb{R}^{d \times d}$ is referred to as positive semi-definite if $\mathbf{x}^T \mathbf{Q} \mathbf{x} \geq 0$ for every vector $\mathbf{x} \in \mathbb{R}^d$. 10, 15, 78, 90

prediction A prediction is an estimate or approximation for some quantity of interest. ML revolves around learning or finding a hypothesis map h that reads in the features \mathbf{x} of a data point and delivers a prediction $\hat{y} := h(\mathbf{x})$ for its label y . 7, 14, 25, 26, 28, 30, 33, 34, 37, 39, 64, 75, 110, 147, 165, 175

probabilistic model A probabilistic model interprets data points as realizations of RVs with a joint probability distribution. This joint probability distribution typically involves parameters which have to be manually chosen (=design choice) or learnt via statistical inference methods [39]. 35, 48, 49, 54, 56, 58, 87, 95, 101, 133, 155, 157, 165, 177

probability We assign a probability value, typically chosen in the interval $[0, 1]$, to each event that might occur in a random experiment [4, 35, 104, 115]. 72

probability density function (pdf) The probability density function (pdf) $p(x)$ of a real-valued RV $x \in \mathbb{R}$ is a particular representation of its probability distribution. If the pdf exists, it can be used to compute the probability that x takes on a value from a (measurable) set $\mathcal{B} \subseteq \mathbb{R}$ via $p(x \in \mathcal{B}) = \int_{\mathcal{B}} p(x') dx'$ [4, Ch. 3]. The pdf of a vector-valued RV $\mathbf{x} \in \mathbb{R}^d$ (if it exists) allows to compute the probability that \mathbf{x} falls into

a (measurable) region \mathcal{R} via $p(\mathbf{x} \in \mathcal{R}) = \int_{\mathcal{R}} p(\mathbf{x}') dx'_1 \dots dx'_d$ [4, Ch. 3].

162

probability distribution The data generated in some ML applications can be reasonably well modelled as realizations of a RV. The overall statistical properties (or intrinsic structure) of such data are then governed by the probability distribution of this RV. We use the term probability distribution in a highly informal manner and mean the collection of probabilities assigned to different values or value ranges of a RV. The probability distribution of a binary RV $y \in \{0, 1\}$ is fully specified by the probabilities $p(y = 0)$ and $p(y = 1) (= 1 - p(y = 0))$. The probability distribution of a real-valued RV $x \in \mathbb{R}$ might be specified by a probability density function $p(x)$ such that $p(x \in [a, b]) \approx p(a)|b - a|$. In the most general case, a probability distribution is defined by a probability measure [35, 41]. 11, 35, 42, 48–53, 59, 63, 72, 81, 87, 89, 95, 99, 101, 102, 133, 146, 155, 162, 163, 167–169, 171, 175, 177

projected GD Projected GD extends basic GD for unconstrained optimization to handle constraints on the optimization variable (model parameters). A single iteration of projected GD consists of first taking a gradient step and then projecting the result back into a constrain set.

123

proximity operator Given a convex function and a vector \mathbf{x} , we define its proximity operator as

$$\mathbf{prox}_{L_i(\cdot), 2\lambda}(\mathbf{w}'') := \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} f(\mathbf{w}) + (\rho/2) \|\mathbf{w} - \mathbf{w}'\|_2^2 \text{ with } \rho > 0.$$

Convex functions for which the proximity operator can be computed efficiently are sometimes referred to as “proximable” or “simple” [116].
69, 71, 127, 146

quadratic function A quadratic function $f(\mathbf{w})$, reading in a vector $\mathbf{w} \in \mathbb{R}^d$ as its argument, is such that

$$f(\mathbf{w}) = \mathbf{w}^T \mathbf{Q} \mathbf{w} + \mathbf{q}^T \mathbf{w} + a,$$

with some matrix $\mathbf{Q} \in \mathbb{R}^{d \times d}$, vector $\mathbf{q} \in \mathbb{R}^d$ and scalar $a \in \mathbb{R}$. 73, 91

random variable (RV) A random variable is a mapping from a probability space \mathcal{P} to a value space [35]. The probability space, whose element are elementary events, is equipped with a probability measure that assigns a probability to subsets of \mathcal{P} . A binary random variable maps elementary events to a set containing two different value, such as $\{-1, 1\}$ or $\{\text{cat}, \text{no cat}\}$. A real-valued random variable maps elementary events to real numbers \mathbb{R} . A vector-valued random variable maps elementary events to the Euclidean space \mathbb{R}^d . Probability theory uses the concept of measurable spaces to rigorously define and study the properties of (large) collections of random variables [35, 41]. 11, 14, 42, 48, 50, 51, 54, 55, 58, 59, 62, 72, 81, 87, 89, 91, 95, 101, 133, 146, 148, 151, 155, 158, 162, 163, 167–169, 171–173, 175, 176, 180

realization Consider a RV x which maps each element (outcome, or elementary event) $\omega \in \mathcal{P}$ of a probability space \mathcal{P} to an element a of a measurable space \mathcal{N} [5, 35, 104]. A realization of x is any element $a' \in \mathcal{N}$

such that there is an element $\omega' \in \mathcal{P}$ with $x(\omega') = a'$. 14, 42, 48, 50, 51, 54, 55, 58, 59, 62, 63, 67, 72, 87, 89, 91, 95, 101, 133, 171, 175

regression Regression problems revolve around the problem of predicting a numeric label solely from the features of a data point. 26

regularization Regularization techniques modify the ERM principle such that the learnt hypothesis performs well (generalizes) beyond the training set. One specific implementation of regularization is to add a penalty or regularization term to the objective function of ERM (which is the average loss on the training set). This regularization term can be interpreted as an estimate for the increase in the expected loss (risk) compared to the average loss on the training set. 2, 15, 20, 21, 61, 65, 66, 68, 69, 71, 81, 82, 93, 98, 100, 107, 111, 113, 115, 116, 122, 123, 159, 174, 175, 177

regularized empirical risk minimization Synonym for SRM. 68, 69, 71, 108, 111, 124, 125, 159

regularizer A regularizer assigns each hypothesis h from a hypothesis space \mathcal{H} a quantitative measure $\mathcal{R}\{h\}$ for how much its prediction error on a training set might differ from its prediction errors on data points outside the training set. Ridge regression uses the regularizer $\mathcal{R}\{h\} := \|\mathbf{w}\|_2^2$ for linear hypothesis maps $h^{(\mathbf{w})}(\mathbf{x}) := \mathbf{w}^T \mathbf{x}$ [1, Ch. 3]. The Lasso uses the regularizer $\mathcal{R}\{h\} := \|\mathbf{w}\|_1$ for linear hypothesis maps $h^{(\mathbf{w})}(\mathbf{x}) := \mathbf{w}^T \mathbf{x}$ [1, Ch. 3]. 66, 67

ridge regression Ridge regression learns the parameter (or weight) vector \mathbf{w}

of a linear hypothesis map $h^{(\mathbf{w})}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$. The quality of a particular choice for the parameter vector \mathbf{w} is measured by the sum of two components. The first component is the average squared error loss incurred by $h^{(\mathbf{w})}$ on a set of labeled data points (the training set). The second component is the scaled squared Euclidean norm $\lambda \|\mathbf{w}\|_2^2$ with a regularization parameter $\lambda > 0$. It can be shown that the effect of adding $\lambda \|\mathbf{w}\|_2^2$ to the average squared error loss is equivalent to replacing the original data points by an ensemble of realizations of a RV centered around these data points. 66, 67, 71, 81, 82, 98, 174

risk Consider a hypothesis h that is used to predict the label y of a data point based on its features \mathbf{x} . We measure the quality of a particular prediction using a loss function $L((\mathbf{x}, y), h)$. If we interpret data points as the realizations of i.i.d. RVs, also the $L((\mathbf{x}, y), h)$ becomes the realization of a RV. Using such an i.i.d. assumption allows to define the risk of a hypothesis as the expected loss $\mathbb{E}\{L((\mathbf{x}, y), h)\}$. Note that the risk of h depends on both, the specific choice for the loss function and the probability distribution of the data points. 48–51, 53–56, 59, 61, 64, 65, 74, 89, 148, 169, 174, 179

sample covariance matrix The sample covariance matrix $\hat{\Sigma} \in \mathbb{R}^{d \times d}$ for a given set of feature vectors $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d$ is defined as

$$\hat{\Sigma} = (1/m) \sum_{r=1}^m (\mathbf{x}^{(r)} - \hat{\mathbf{m}})(\mathbf{x}^{(r)} - \hat{\mathbf{m}})^T.$$

Here, we used the sample mean $\hat{\mathbf{m}}$. 102

sample mean The sample mean $\mathbf{m} \in \mathbb{R}^{d \times d}$ for a given set of feature vectors

$\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d$ is defined as

$$\mathbf{m} = (1/m) \sum_{r=1}^m \mathbf{x}^{(r)}.$$

102, 175

sample size The number of individual data points contained in a dataset that is obtained from realizations of i.i.d. RVs. 42, 46, 55, 159

scatterplot A visualization technique that depicts data points by markers in a two-dimensional plane. 153, 156

smooth We refer to a real-valued function as smooth if it is differentiable and its gradient is continuous [37, 54]. In particular, a differentiable function $f(\mathbf{w})$ is referred to as β -smooth if the gradient $\nabla f(\mathbf{w})$ is Lipschitz continuous with Lipschitz constant β , i.e.,

$$\|\nabla f(\mathbf{w}) - \nabla f(\mathbf{w}')\| \leq \beta \|\mathbf{w} - \mathbf{w}'\|.$$

169

spectral clustering Spectral clustering groups the nodes of an undirected graph by applying k -means clustering to node-wise feature vectors. These feature vectors are built from the eigenvectors the graph Laplacian matrix [63, 76]. 108

squared error loss The squared error loss measures the prediction error of a hypothesis h when predicting a numeric label $y \in \mathbb{R}$ from the features \mathbf{x} of a data point. It is defined as

$$L((\mathbf{x}, y), h) := \left(y - \underbrace{h(\mathbf{x})}_{=\hat{y}}\right)^2. \quad (105)$$

37–41, 58, 66, 67, 71, 111, 122, 123, 165, 168, 170, 175

statistical aspects By statistical aspects of a ML method, we refer to (properties of) the probability distribution of its output given a probabilistic model for the data fed into the method. 25, 27, 29, 30, 34, 36, 42, 51, 52, 95, 97, 100, 106, 111, 165, 179

stochastic gradient descent Stochastic GD is obtained from GD by replacing the gradient of the objective function by a noisy (or stochastic) estimate. 74, 87–89, 91

stopping criterion Many ML methods use iterative algorithms that construct a sequence of model parameters (such as the weights of a linear map or the weights of an ANN) that (hopefully) converge to an optimal choice for the model parameters. In practice, given finite computational resources, we need to stop iterating after a finite number of times. A stopping criterion is any well-defined condition required for stopping iterating. 76–78, 84, 120–122, 124, 125, 134, 136

strongly convex A continuously differentiable real-valued function $f(\mathbf{x})$ is strongly convex with coefficient σ if $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + (\sigma/2) \|\mathbf{y} - \mathbf{x}\|_2^2$ [54], [36, Sec. B.1.1.]. 74

structural risk minimization Structural risk minimization is the problem of finding the hypothesis that optimally balances the average loss (or empirical risk) on a training set with a regularization term. The regularization term penalizes a hypothesis that is not robust against (small) perturbations of the data points in the training set. 164, 174

subgradient For a real-valued function $f : \mathbb{R}^d \rightarrow \mathbb{R} : \mathbf{w} \mapsto f(\mathbf{w})$, a vector \mathbf{a} such that $f(\mathbf{w}) \geq f(\mathbf{w}') + (\mathbf{w} - \mathbf{w}')^T \mathbf{a}$ is referred to as a subgradient of f at \mathbf{w}' [117, 118]. 44

subgradient descent Subgradient descent is a generalization of GD that is obtained by using sub-gradients (instead of gradients) to construct local approximations of an objective function. The objective function might be the empirical risk $\widehat{L}(h^{(\mathbf{w})} | \mathcal{D})$ viewed as a function of the model parameters \mathbf{w} that select a hypothesis $h^{(\mathbf{w})} \in \mathcal{H}$. 44

support vector machine A binary classification method for learning a linear hypothesis map that maximally separates data points the two classes in the feature space (“maximum margin”). Maximizing this separation is equivalent to minimizing a regularized variant of the hinge loss (102). 160

training error The average loss of a hypothesis when predicting the labels of data points in a training set. We sometimes refer by training error also the minimum average loss incurred on the training set by any hypothesis out of a hypothesis space. 14, 53, 61, 63–65, 178

training set A set of data points that is used in ERM to learn a hypothesis \hat{h} . The average loss of \hat{h} on the training set is referred to as the training error. The comparison between training error and validation error of \hat{h} allows to diagnose ML methods and informs how to improve them (e.g., using a different hypothesis space or collecting more data points). 13, 14, 38, 40, 48, 53–56, 58, 59, 61, 63–67, 69, 78, 86, 88, 89, 93, 99, 143, 154, 159, 160, 165, 166, 170, 171, 174, 175, 177–179

trustworthiness Beside the computational aspects and statistical aspects, a third main design aspect for ML methods is their trustworthiness. The European Union has put forward seven key requirements for trustworthy AI systems (that typically build on ML methods) [89]: **human agency and oversight, technical robustness and safety, privacy and data governance, transparency, diversity non-discrimination and fairness, societal and environmental well-being, accountability.** 27, 28

underfitting Consider a ML method that uses ERM to learn a hypothesis with minimum empirical risk on a given training set. Such a method is “underfitting” the training set if it is not able to learn a hypothesis with sufficiently small empirical risk on the training set. If a method is underfitting it will typically also not be able to learn a hypothesis with a small risk. 64

validation Consider a hypothesis \hat{h} that has been learn via ERM on some training set \mathcal{D} . Validation refers to the practice of trying out a hypothesis \hat{h} on a validation set that consists of data points that are not contained in the training set \mathcal{D} . 77, 107

validation error Consider a hypothesis \hat{h} which is obtained by ERM on a training set. The average loss of \hat{h} on a validation set, which is different from the training set, is referred to as the validation error. 14, 178

validation set A set of data points that has not been used as training set in ERM to train a hypothesis \hat{h} . The average loss of \hat{h} on the validation

set is referred to as the validation error and used to diagnose the ML method. The comparison between training and validation error informs adaptations of the ML method (such as using a different hypothesis space). 14, 143, 179

Vapnik–Chervonenkis (VC) dimension The VC dimension of an infinite hypothesis space is a widely-used measure for its size. We refer to [119] for a precise definition of VC dimension as well as a discussion of its basic properties and use in ML. 63

variance The variance of a real-valued RV x is defined as the expectation $\mathbb{E}\{(x - \mathbb{E}\{x\})^2\}$ of the squared difference x and its expectation $\mathbb{E}\{x\}$. We extend this definition to vector-valued RVs \mathbf{x} as $\mathbb{E}\{\|\mathbf{x} - \mathbb{E}\{\mathbf{x}\}\|_2^2\}$. 14, 54, 55, 59

vertical FL Vertical FL refers to applications with local datasets that are constituted by the same data points but characterizing them with different features [87]. For example, different healthcare providers might all contain information about the same population of patients. However, different healthcare providers collect different measurements (blood values, electrocardiography, lung x-ray) for the same patients. 113, 116

References

- [1] A. Jung, *Machine Learning: The Basics*, 1st ed. Springer Singapore, Feb. 2022.
- [2] W. Rudin, *Real and Complex Analysis*, 3rd ed. New York: McGraw-Hill, 1987.
- [3] G. H. Golub and C. F. Van Loan, *Matrix Computations*, 3rd ed. Baltimore, MD: Johns Hopkins University Press, 1996.
- [4] D. Bertsekas and J. Tsitsiklis, *Introduction to Probability*, 2nd ed. Athena Scientific, 2008.
- [5] W. Rudin, *Principles of Mathematical Analysis*, 3rd ed. New York: McGraw-Hill, 1976.
- [6] S. Cui, A. Hero, Z.-Q. Luo, and J. Moura, Eds., *Big Data over Networks*. Cambridge Univ. Press, 2016.
- [7] M. Wollschlaeger, T. Sauter, and J. Jasperneite, “The future of industrial communication: Automation networks in the era of the internet of things and industry 4.0,” *IEEE Industrial Electronics Magazine*, vol. 11, no. 1, pp. 17–27, 2017.
- [8] M. Satyanarayanan, “The emergence of edge computing,” *Computer*, vol. 50, no. 1, pp. 30–39, Jan. 2017. [Online]. Available: <https://doi.org/10.1109/MC.2017.9>

- [9] H. Ates, A. Yetisen, F. Güder, and C. Dincer, “Wearable devices for the detection of covid-19,” *Nature Electronics*, vol. 4, no. 1, pp. 13–14, 2021. [Online]. Available: <https://doi.org/10.1038/s41928-020-00533-1>
- [10] H. Boyes, B. Hallaq, J. Cunningham, and T. Watson, “The industrial internet of things (iiot): An analysis framework,” *Computers in Industry*, vol. 101, pp. 1–12, 2018. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0166361517307285>
- [11] M. E. J. Newman, *Networks: An Introduction*. Oxford Univ. Press, 2010.
- [12] A. Barabási, N. Gulbahce, and J. Loscalzo, “Network medicine: a network-based approach to human disease,” *Nature Reviews Genetics*, vol. 12, no. 56, 2011.
- [13] K. Grantz, H. Meredith, D. Cummings, C. Metcalf, B. Grenfell, J. Giles, S. Mehta, S. Solomon, A. Labrique, N. Kishore, C. Buckee, and A. Wesolowski, “The use of mobile phone data to inform analysis of COVID-19 pandemic epidemiology,” *Nature Communications*, vol. 11, no. 1, p. 4961, 2020. [Online]. Available: <https://doi.org/10.1038/s41467-020-18190-5>
- [14] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, ser. Proceedings of Machine Learning Research, A. Singh and J. Zhu, Eds., vol. 54. Fort Lauderdale, FL,

- USA: PMLR, 20–22 Apr 2017, pp. 1273–1282. [Online]. Available: <http://proceedings.mlr.press/v54/mcmahan17a.html>
- [15] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, “Federated learning: Challenges, methods, and future directions,” *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 50–60, May 2020.
 - [16] Y. Cheng, Y. Liu, T. Chen, and Q. Yang, “Federated learning for privacy-preserving ai,” *Communications of the ACM*, vol. 63, no. 12, pp. 33–36, Dec. 2020.
 - [17] N. Agarwal, A. Suresh, F. Yu, S. Kumar, and H. McMahan, “cpSGD: Communication-efficient and differentially-private distributed sgd,” in *Proc. Neural Inf. Proc. Syst. (NIPS)*, 2018.
 - [18] V. Smith, C.-K. Chiang, M. Sanjabi, and A. Talwalkar, “Federated Multi-Task Learning,” in *Advances in Neural Information Processing Systems*, vol. 30, 2017. [Online]. Available: <https://proceedings.neurips.cc/paper/2017/file/6211080fa89981f66b1a0c9d55c61d0f-Paper.pdf>
 - [19] J. You, J. Wu, X. Jin, and M. Chowdhury, “Ship compute or ship data? why not both?” in *18th USENIX Symposium on Networked Systems Design and Implementation (NSDI 21)*. USENIX Association, April 2021, pp. 633–651. [Online]. Available: <https://www.usenix.org/conference/nsdi21/presentation/you>
 - [20] D. Tse and P. Viswanath, *Fundamentals of Wireless Communication*. Cambridge University Press, 2005.

- [21] T. Yang, G. Andrew, H. Eichner, H. Sun, W. Li, N. Kong, D. Ramage, and F. Beaufays, “Applied federated learning: Improving google keyboard query suggestions,” 2018. [Online]. Available: <https://arxiv.org/abs/1812.02903>
- [22] A. Ghosh, J. Chung, D. Yin, and K. Ramchandran, “An efficient framework for clustered federated learning,” in *34th Conference on Neural Information Processing Systems (NeurIPS 2020)*, Vancouver, Canada, 2020.
- [23] F. Sattler, K. Müller, and W. Samek, “Clustered federated learning: Model-agnostic distributed multitask optimization under privacy constraints,” *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [24] F. Pedregosa, “Scikit-learn: Machine learning in python,” *Journal of Machine Learning Research*, vol. 12, no. 85, pp. 2825–2830, 2011. [Online]. Available: <http://jmlr.org/papers/v12/pedregosa11a.html>
- [25] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, UK: Cambridge Univ. Press, 2004.
- [26] D. P. Bertsekas and J. Tsitsiklis, *Parallel and Distributed Computation: Numerical Methods*. Athena Scientific, 2015.
- [27] M. Conti and D. Moretti, “System level analysis of the bluetooth standard,” in *Design, Automation and Test in Europe*, 2005, pp. 118–123 Vol. 3.

- [28] M. Ribeiro, S. Singh, and C. Guestrin, ““Why should i trust you?”: Explaining the predictions of any classifier,” in *Proc. 22nd ACM SIGKDD*, Aug. 2016, pp. 1135–1144.
- [29] R. Raina, A. Madhavan, and A. Ng, “Large-scale deep unsupervised learning using graphics processors,” in *Proceedings of the 26th Annual International Conference on Machine Learning*, ser. ICML ’09. New York, NY, USA: Association for Computing Machinery, 2009, pp. 873–880. [Online]. Available: <https://doi.org/10.1145/1553374.1553486>
- [30] X. Lin, Y. Rivenson, N. Yardimci, M. Veli, Y. Luo, M. Jarrahi, and A. Ozcan, “All-optical machine learning using diffractive deep neural networks,” *Science*, vol. 361, no. 6406, pp. 1004–1008, 2018. [Online]. Available: <https://www.science.org/doi/abs/10.1126/science.aat8084>
- [31] C. Molnar, *Interpretable Machine Learning - A Guide for Making Black Box Models Explainable*. [online] Available: <https://christophm.github.io/interpretable-ml-book/>, 2019.
- [32] A. Jung and P. Nardelli, “An information-theoretic approach to personalized explainable machine learning,” *IEEE Sig. Proc. Lett.*, vol. 27, pp. 825–829, 2020.
- [33] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.
- [34] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, ser. Springer Series in Statistics. New York, NY, USA: Springer, 2001.

- [35] P. Billingsley, *Probability and Measure*, 3rd ed. New York: Wiley, 1995.
- [36] D. P. Bertsekas, *Convex Optimization Algorithms*. Athena Scientific, 2015.
- [37] S. Bubeck, “Convex optimization. algorithms and complexity.” in *Foundations and Trends in Machine Learning*. Now Publishers, 2015, vol. 8.
- [38] A. Nedic and A. Ozdaglar, “Distributed subgradient methods for multi-agent optimization,” *IEEE Transactions on Automatic Control*, vol. 54, no. 1, pp. 48–61, 2009.
- [39] E. L. Lehmann and G. Casella, *Theory of Point Estimation*, 2nd ed. New York: Springer, 1998.
- [40] M. J. Wainwright and M. I. Jordan, *Graphical Models, Exponential Families, and Variational Inference*, ser. Foundations and Trends in Machine Learning. Hanover, MA: Now Publishers, 2008, vol. 1, no. 1–2.
- [41] R. Gray, *Probability, Random Processes, and Ergodic Properties*, 2nd ed. New York: Springer, 2009.
- [42] M. Wainwright, *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge: Cambridge University Press, 2019.
- [43] J. Cho and H. White, “Generalized runs tests for the iid hypothesis,” *Journal of Econometrics*, vol. 162, no. 2, pp. 326–344, 2011.
[Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0304407611000285>

- [44] P. J. Huber, *Robust Statistics*. New York: Wiley, 1981.
- [45] T. W. Anderson, *An Introduction to Multivariate Statistical Analysis*. Wiley, 2003.
- [46] L. Breiman and D. Freedman, “How many variables should be entered in a regression equation?” *Journal of the American Statistical Association*, vol. 78, no. 381, pp. 131–136, 1983.
- [47] J. Mourtada, “Exact minimax risk for linear least squares, and the lower tail of sample covariance matrices,” *The Annals of Statistics*, vol. 50, no. 4, pp. 2157 – 2178, 2022. [Online]. Available: <https://doi.org/10.1214/22-AOS2181>
- [48] A. V. Oppenheim, R. W. Schaffer, and J. R. Buck, *Discrete-Time Signal Processing*, 2nd ed. Englewood Cliffs, NJ: Prentice Hall, 1998.
- [49] N. Tran, O. Abramenko, and A. Jung, “On the sample complexity of graphical model selection from non-stationary samples,” *IEEE Transactions on Signal Processing*, vol. 68, pp. 17–32, 2020.
- [50] N. Parikh and S. Boyd, “Proximal algorithms,” *Foundations and Trends in Optimization*, vol. 1, no. 3, pp. 123–231, 2013.
- [51] G. Strang, *Introduction to Linear Algebra*, 5th ed. Wellesley-Cambridge Press, MA, 2016.
- [52] G. Golub and C. van Loan, “An analysis of the total least squares problem,” *SIAM J. Numerical Analysis*, vol. 17, no. 6, pp. 883–893, Dec. 1980.

- [53] Y. L. T. Schaul, X. Zhang, “No more pesky learning rates,” in *Proc. of the 30th International Conference on Machine Learning, PMLR 28(3)*, vol. 28, Atlanta, Georgia, June 2013, pp. 343–351.
- [54] Y. Nesterov, *Introductory lectures on convex optimization*, ser. Applied Optimization. Kluwer Academic Publishers, Boston, MA, 2004, vol. 87, a basic course. [Online]. Available: <http://dx.doi.org/10.1007/978-1-4419-8853-9>
- [55] A. Papoulis and S. U. Pillai, *Probability, Random Variables, and Stochastic Processes*, 4th ed. New York: Mc-Graw Hill, 2002.
- [56] A. Jung, “A fixed-point of view on gradient methods for big data,” *Frontiers in Applied Mathematics and Statistics*, vol. 3, 2017. [Online]. Available: <https://www.frontiersin.org/article/10.3389/fams.2017.00018>
- [57] A. Rakhlin, O. Shamir, and K. Sridharan, “Making gradient descent optimal for strongly convex stochastic optimization,” in *Proceedings of the 29 th International Conference on Machine Learning*, Edinburgh, Scotland, UK,, 2012.
- [58] A. Jung, “Networked exponential families for big data over networks,” *IEEE Access*, vol. 8, pp. 202 897–202 909, 2020.
- [59] H. Ambos, N. Tran, and A. Jung, “The logistic network lasso,” *arXiv*, 2018.
- [60] —, “Classifying big data over networks via the logistic network lasso,” in *Proc. 52nd Asilomar Conf. Signals, Systems, Computers*, Oct./Nov. 2018.

- [61] A. Jung and N. Tran, “Localized linear regression in networked data,” *IEEE Sig. Proc. Lett.*, vol. 26, no. 7, Jul. 2019.
- [62] P. J. Brockwell and R. A. Davis, *Time Series: Theory and Methods*. Springer New York, 1991.
- [63] U. von Luxburg, “A tutorial on spectral clustering,” *Statistics and Computing*, vol. 17, no. 4, pp. 395–416, Dec. 2007.
- [64] A. Y. Ng, M. I. Jordan, and Y. Weiss, “On spectral clustering: Analysis and an algorithm,” in *Adv. Neur. Inf. Proc. Syst.*, 2001.
- [65] V. Kalofolias, “How to learn a graph from smooth signals,” in *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, ser. Proceedings of Machine Learning Research, A. Gretton and C. C. Robert, Eds., vol. 51. Cadiz, Spain: PMLR, 09–11 May 2016, pp. 920–929.
- [66] X. Dong, D. Thanou, M. Rabbat, and P. Frossard, “Learning graphs from data: A signal representation perspective,” *IEEE Signal Processing Magazine*, vol. 2019, no. 3, May 2019.
- [67] L. Stanković, D. Mandić, M. Daković, M. Brajović, B. Scalzo, S. Li, and A. Constantinides, “Data analytics on graphs part III: Machine learning on graphs, from graph topology to applications,” *Foundations and Trends® in Machine Learning*, vol. 13, no. 4, pp. 332–530, 2020. [Online]. Available: <http://dx.doi.org/10.1561/22000000078-3>
- [68] D. Hallac, J. Leskovec, and S. Boyd, “Network lasso: Clustering and optimization in large graphs,” in *Proc. SIGKDD*, 2015, pp. 387–396.

- [69] Y. SarcheshmehPour, M. Leinonen, and A. Jung, “Federated learning from big data over networks,” in *Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, preprint: <https://arxiv.org/pdf/2010.14159.pdf>, 2021.
- [70] A. Jung, N. Tran, and A. Mara, “When is Network Lasso Accurate?” *Front. Appl. Math. Stat.*, vol. 3, Jan. 2018.
- [71] B. Nadler, N. Srebro, and X. Zhou, “Statistical analysis of semi-supervised learning: The limit of infinite unlabelled data,” in *Advances in Neural Information Processing Systems 22*, 2009, pp. 1330–1338.
- [72] F. Bach, R. Jenatton, J. Mairal, and G. Obozinski, “Optimization with sparsity-inducing penalties,” *Found. Trends Mach. Learn.*, vol. 4, no. 1, pp. 1–106, Jan. 2012.
- [73] L. Jacob, J.-P. Vert, and F. Bach, “Clustered multi-task learning: A convex formulation,” in *Advances in Neural Information Processing Systems*, D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, Eds., vol. 21. Curran Associates, Inc., 2009. [Online]. Available: <https://proceedings.neurips.cc/paper/2008/file/fccb3cdc9acc14a6e70a12f74560c026-Paper.pdf>
- [74] T. Evgeniou, C. Micchelli, and M. Pontil, “Learning multiple tasks with kernel methods,” *Journal of Machine Learning Research*, vol. 6, no. 21, pp. 615–637, 2005. [Online]. Available: <http://jmlr.org/papers/v6/evgeniou05a.html>

- [75] L. Z. Y. SarcheshmehPour, Y. Tian and A. Jung, “Networked federated learning,” *arXiv e-prints*, 2022.
- [76] Y. SarcheshmehPour, Y. Tian, L. Zhang, and A. Jung, “Flow-based clustering and spectral clustering: A comparison,” in *2021 55th Asilomar Conference on Signals, Systems, and Computers*, 2021, pp. 1292–1296.
- [77] D. Sun, K.-C. Toh, and Y. Yuan, “Convex clustering: Model, theoretical guarantee and efficient algorithm,” *Journal of Machine Learning Research*, vol. 22, no. 9, pp. 1–32, 2021. [Online]. Available: <http://jmlr.org/papers/v22/18-694.html>
- [78] K. Pelckmans, J. D. Brabanter, J. Suykens, and B. D. Moor, “Convex clustering shrinkage,” in *PASCAL Workshop on Statistics and Optimization of Clustering Workshop*, 2005.
- [79] C. G. Atkeson, A. W. Moore, and S. Schaal, “Locally weighted learning,” *Artif. Intell. Rev.*, vol. 11, no. 1–5, pp. 11–73, feb 1997. [Online]. Available: <https://doi.org/10.1023/A:1006559212014>
- [80] V. Mothukuri, R. M. Parizi, S. Pouriye, Y. Huang, A. Dehghantanha, and G. Srivastava, “A survey on security and privacy of federated learning,” *Future Generation Computer Systems*, vol. 115, pp. 619–640, 2021. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0167739X20329848>
- [81] Q. Yang, Y. Liu, Y. Cheng, Y. Kang, T. Chen, and H. Yu, *Federated Learning*, 1st ed. Springer, 2022.

- [82] S. R. Pokhrel and J. Choi, “Federated learning with blockchain for autonomous vehicles: Analysis and design challenges,” *IEEE Transactions on Communications*, vol. 68, no. 8, pp. 4734–4746, 2020.
- [83] O. Chapelle, B. Schölkopf, and A. Zien, Eds., *Semi-Supervised Learning*. Cambridge, Massachusetts: The MIT Press, 2006.
- [84] A. Ghosh, J. Chung, D. Yin, and K. Ramchandran, “An efficient framework for clustered federated learning,” *IEEE Transactions on Information Theory*, vol. 68, no. 12, pp. 8076–8091, 2022.
- [85] H. Ludwig and N. Baracaldo, Eds., *Federated Learning: A Comprehensive Overview of Methods and Applications*. Springer, 2022.
- [86] A. Shamsian, A. Navon, E. Fetaya, and G. Chechik, “Personalized federated learning using hypernetworks,” in *Proceedings of the 38th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, M. Meila and T. Zhang, Eds., vol. 139. PMLR, 18–24 Jul 2021, pp. 9489–9502. [Online]. Available: <https://proceedings.mlr.press/v139/shamsian21a.html>
- [87] Q. Yang, Y. Liu, Y. Cheng, Y. Kang, T. Chen, and H. Yu, *Vertical Federated Learning*. Cham: Springer International Publishing, 2020, pp. 69–81. [Online]. Available: https://doi.org/10.1007/978-3-031-01585-4_5
- [88] ———, *Horizontal Federated Learning*. Cham: Springer International Publishing, 2020, pp. 49–67. [Online]. Available: https://doi.org/10.1007/978-3-031-01585-4_4

- [89] E. Commission, C. Directorate-General for Communications Networks, and Technology, *The Assessment List for Trustworthy Artificial Intelligence (ALTAI) for self assessment*. Publications Office, 2020.
- [90] M. J. Sheller, B. Edwards, G. A. Reina, J. Martin, S. Pati, A. Kotrotsou, M. Milchenko, W. Xu, D. Marcus, R. R. Colen, and S. Bakas, “Federated learning in medicine: facilitating multi-institutional collaborations without sharing patient data,” *Scientific Reports*, vol. 10, no. 1, p. 12598, 2020. [Online]. Available: <https://doi.org/10.1038/s41598-020-69250-1>
- [91] K. Chaudhuri, C. Monteleoni, and A. Sarwate, “Differentially private empirical risk minimization,” *J. Mach. Learn. Res.*, vol. 12, pp. 1069–1109, Mar. 2011.
- [92] H. Hsu, N. Martinez, M. Bertran, G. Sapiro, and F. P. Calmon, “A survey on statistical, information, and estimation—theoretic views on privacy,” *IEEE BITS the Information Theory Magazine*, vol. 1, no. 1, pp. 45–56, 2021.
- [93] P. Samarati, “Protecting respondents identities in microdata release,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 13, no. 6, pp. 1010–1027, 2001.
- [94] L. Sweeney, “ k -anonymity: a model for protecting privacy,” *International Journal on Uncertainty, Fuzziness and Knowledge-based Systems*, vol. 10, no. 5, pp. 557–570, 2002.

- [95] X. Liu, H. Li, G. Xu, Z. Chen, X. Huang, and R. Lu, “Privacy-enhanced federated learning against poisoning adversaries,” *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 4574–4588, 2021.
- [96] J. Zhang, B. Chen, X. Cheng, H. T. T. Binh, and S. Yu, “PoisonGAN: Generative poisoning attacks against federated learning in edge computing systems,” *IEEE Internet of Things Journal*, vol. 8, no. 5, pp. 3310–3322, 2021.
- [97] N. M. Müller, S. Roschmann, and K. Böttinger, “Defending Against Adversarial Denial-of-Service Data Poisoning Attacks,” *arXiv e-prints*, p. arXiv:2104.06744, April 2021.
- [98] H. Wang, K. Sreenivasan, S. Rajput, H. Vishwakarma, S. Agarwal, J.-y. Sohn, K. Lee, and D. Papailiopoulos, “Attack of the tails: Yes, you really can backdoor federated learning,” in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 16 070–16 084. [Online]. Available: <https://proceedings.neurips.cc/paper/2020/file/b8ffa41d4e492f0fad2f13e29e1762eb-Paper.pdf>
- [99] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. F. del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant, “Array programming with numpy,”

Nature, vol. 585, no. 7825, pp. 357–362, 2020. [Online]. Available: <https://doi.org/10.1038/s41586-020-2649-2>

- [100] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, *Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers*. Hanover, MA: Now Publishers, 2010, vol. 3, no. 1.
- [101] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. New Jersey: Wiley, 2006.
- [102] C. E. Shannon, “Communication in the presence of noise,” 1948.
- [103] P. Halmos, *Naive set theory*. Springer-Verlag, 1974.
- [104] P. R. Halmos, *Measure Theory*. New York: Springer, 1974.
- [105] A. Jung, “Explainable empirical risk minimization,” *submitted to IEEE Sig. Proc. Letters (preprint: <https://arxiv.org/pdf/2009.01492.pdf>)*, 2020.
- [106] J. Chen, L. Song, M. Wainwright, and M. Jordan, “Learning to explain: An information-theoretic perspective on model interpretation,” in *Proc. 35th Int. Conf. on Mach. Learning*, Stockholm, Sweden, 2018.
- [107] D. Gujarati and D. Porter, *Basic Econometrics*. McGraw Hill, 2009.
- [108] Y. Dodge, *The Oxford Dictionary of Statistical Terms*. Oxford University Press, 2003.

- [109] B. Everitt, *Cambridge Dictionary of Statistics*. Cambridge University Press, 2002.
- [110] R. T. Rockafellar, *Network Flows and Monotropic Optimization*. Athena Scientific, Jul. 1998.
- [111] P. Bühlmann and S. van de Geer, *Statistics for High-Dimensional Data*. New York: Springer, 2011.
- [112] C. Lampert, “Kernel methods in computer vision,” *Foundations and Trends in Computer Graphics and Vision*, 2009.
- [113] B. Schölkopf and A. Smola, *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. Cambridge, MA, USA: MIT Press, Dec. 2002.
- [114] A. Lapidoth, *A Foundation in Digital Communication*. New York: Cambridge University Press, 2009.
- [115] O. Kallenberg, *Foundations of modern probability*. New York: Springer, 1997.
- [116] L. Condat, “A primal–dual splitting method for convex optimization involving lipschitzian, proximable and linear composite terms,” *Journal of Opt. Th. and App.*, vol. 158, no. 2, pp. 460–479, Aug. 2013.
- [117] D. P. Bertsekas, A. Nedic, and A. E. Ozdaglar, *Convex Analysis and Optimization*. Athena Scientific, 2003.
- [118] D. P. Bertsekas, *Nonlinear Programming*, 2nd ed. Belmont, MA: Athena Scientific, June 1999.

- [119] S. Shalev-Shwartz and S. Ben-David, *Understanding Machine Learning*
– *from Theory to Algorithms*. Cambridge University Press, 2014.