

Machine Learning and Neural Networks III (MATH3431)

Epiphany term

Georgios P. Karagiannis

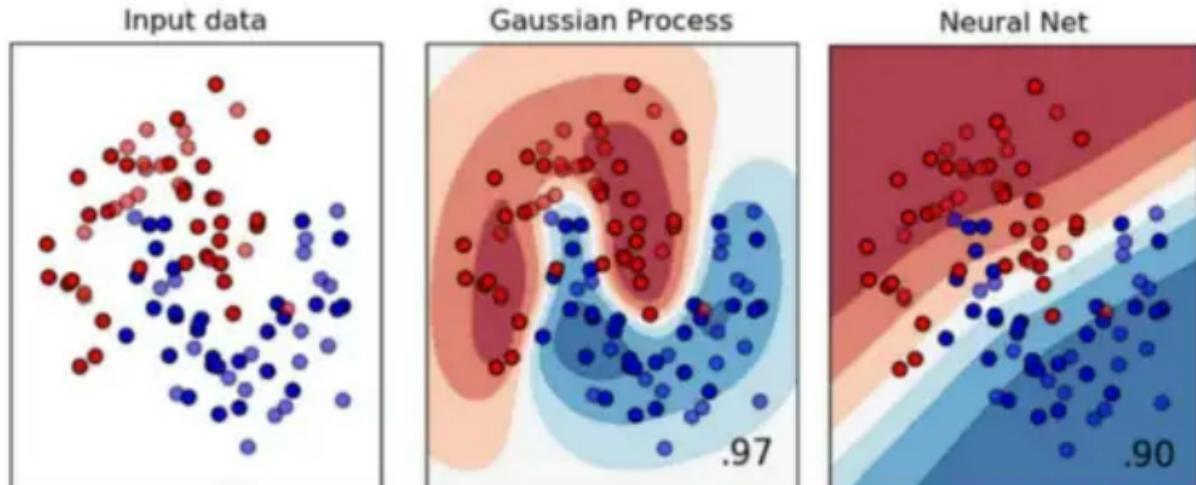
georgios.karagiannis@durham.ac.uk

Department of Mathematical Sciences (Office MCS3088)
Durham University
Stockton Road Durham DH1 3LE UK

2025/05/03 at 06:58:17

Concepts

- Convex learning problems
- Stochastic learning
- Support vector machines
- Artificial neural networks
- Kernel methods, trick
- Gaussian process regression



READING LIST

These lecture Handouts have been derived based on the reading list below.

Main texts:

- Bishop, C. M. (2006). Pattern recognition and machine learning. New York: Springer.
 - It is a classical textbook in machine learning (ML) methods. It discusses all the concepts introduced in the course (not necessarily in the same depth). It is one of the main textbooks in the module. The level on difficulty is easy.
 - Students who wish to have a textbook covering traditional concepts in machine learning are suggested to get a copy of this textbook. It is available online from the Microsoft's website <https://www.microsoft.com/en-us/research/publication/pattern-recognition-machine-learning/>
- Shalev-Shwartz, S., & Ben-David, S. (2014). Understanding machine learning: From theory to algorithms. Cambridge university press.
 - It has several elements of theory about machine learning algorithms. It is one of the main textbooks in the module. The level on difficulty is advanced as it requires moderate knowledge of maths.
- Bishop, C. M. (1995). Neural networks for pattern recognition. Oxford university press.
 - It is a classical textbook about ‘traditional’ artificial neural networks (ANN). It is very comprehensive (compared to others) and it goes deep enough for the module although it may be a bit outdated. It is one of the main textbooks in the module for ANN. The level on difficulty is moderate.

Supplementary textbooks:

- Vapnik, V. (1999). The nature of statistical learning theory. Springer science & business media.
 - Important textbook in the statistical machine learning theory. To have have an in deep understanding about statistical learning, one has to read it together with the textbook of Shalev-Shwartz, S., & Ben-David, S. (2014)
- Devroye, L., Györfi, L., & Lugosi, G. (2013). A probabilistic theory of pattern recognition (Vol. 31). Springer Science & Business Media.
 - Important textbook in the theoretical aspects about machine learning algorithms. The level on difficulty is advanced as it requires moderate knowledge of probability.
- Theodoridis, S. (2015). Machine learning: a Bayesian and optimization perspective. Academic press.
 - A textbook in general machine learning methods. It includes a large collection of different machine learning methods. It does not go too deep into the theoretical details of the method.
- Murphy, K. P. (2012). Machine learning: a probabilistic perspective. MIT press.
 - A popular textbook in general machine learning methods. It discusses all the concepts introduced in the module. It focuses more on the probabilistic/Bayesian framework but not with great detail. It can be used as a comparison textbook

for brief reading about ML methods just to see another perspective than that in (Bishop, C. M., 2006). The level on difficulty is easy.

- Murphy, K. P. (2022). Probabilistic machine learning: an introduction. MIT press.
 - A textbook in general machine learning methods. It covers a smaller number of ML concepts than (Murphy, K. P., 2012) but it contains more fancy/popular topics such as deep learning ideas. It is suggested to be used in the same manner as (Murphy, K. P., 2012). The level on difficulty is easy.
- Bishop, C. M., & Bishop, H. (2024). Deep learning: foundations and concepts. New York: Springer.
 - A textbook in machine learning methods focusing on Neural Networks philosophy and deep learning. Essentially, this is an update of the earlier book “Bishop, C. M. (1995). Neural networks for pattern recognition. Oxford university press.” where: the concepts about Neural networks, stochastic gradient descent have been extended; new deep learning concepts have been added; and the kernel methods are left out.
- Ripley, B. D. (2007). Pattern recognition and neural networks. Cambridge university press.
 - A classical textbook in artificial neural networks (ANN) that also covers other machine learning concepts. It contains interesting theory about ANN.
 - It is suggested to be used as a supplementary reading for neural networks as it contains a few interesting theoretical results. The level on difficulty is moderate.
- Williams, C. K., & Rasmussen, C. E. (2006). Gaussian processes for machine learning (Vol. 2, No. 3, p. 4). Cambridge, MA: MIT press.
 - A classic book in Gaussian process regression (GPR) that covers the material we will discuss in the course about GPR. It can be used as a companion textbook with that of (Bishop, C. M., 2006). The level on difficulty is easy.

Preparatory textbooks:

- Strang, G. (2019). Linear algebra and learning from data. Wellesley-Cambridge Press.
 - Material regarding linear algebra concepts.
- Boyd, S. P., & Vandenberghe, L. (2004). Convex optimization. Cambridge university press.
 - Material regarding Convex optimization, quadratic optimization, Lagrange duality.

CONTENTS

- (1) Lecture notes 1: Machine learning –A recap on: definitions, notation, and formalism
- (2) Lecture notes 2: Elements of convex learning problems
- (3) Lecture notes 3: Learnability, and stability
- (4) Lecture notes 4: Gradient descent
- (5) Lecture notes 5: Stochastic gradient descent
- (6) Lecture notes 6: Variations of stochastic gradient descent
- (7) Lecture notes 7: Bayesian Learning via Stochastic gradient and Stochastic gradient Langevin dynamics
- (8) Lecture notes 8: Support Vector Machines
- (9) Lecture notes 9: Kernel methods
- (10) Lecture notes 10: Multi-class classification
- (11) Lecture notes 11: Artificial neural networks
- (12) Lecture notes 12: Gaussian process regression

Lecture notes 1: Machine learning -A recap on: definitions, notation, and formalism

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim. To get some definitions and set-up about the learning procedure; essentially to formalize what introduced in term 1.

Reading list & references:

- Bishop, C. M. (2006). Pattern recognition and machine learning. New York: Springer.
– Ch. 1 Introduction
- Shalev-Shwartz, S., & Ben-David, S. (2014). Understanding machine learning: From theory to algorithms. Cambridge university press.
– Ch. 1 Introduction

1. GENERAL INTRODUCTIONS AND DEFINITIONS

Pattern recognition is the automated discovery of patterns and regularities in data $z \in \mathcal{Z}$. **Machine learning (ML)** are statistical procedures for building and understanding probabilistic methods that 'learn'. **ML algorithms** build a (probabilistic/deterministic) model able to make predictions or decisions with minimum human interference and can be used for pattern recognition. **Learning** (or training, estimation, fitting) is called the procedure where the ML model is tuned. **Training data** (or observations, sample data set, examples) is a set of observables $\{z_i \in \mathcal{Z}\}$ used to tune the parameters of the ML model. By \mathcal{Z} we denote the examples (or observables) domain. **Test set** is a set of available examples/observables $\{z'_i\}$ (different than the training data) used to verify the performance of the ML model for a given a measure of success. **Measure of success** (or performance) is a quantity that indicates how bad the corresponding ML model or Algorithm performs (eg quantifies the failure/error), and can also be used for comparisons among different ML models; eg, **Risk function** or **Empirical Risk Function**. Two main problems in ML are the supervised learning (we will focus on this here) and the unsupervised learning.

Supervised learning problems involve applications where the training data $z \in \mathcal{Z}$ comprise examples of the input vectors $x \in \mathcal{X}$ along with their corresponding target vectors $y \in \mathcal{Y}$; i.e. $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$. By \mathcal{X} we denote the inputs (or instances) domain, and by \mathcal{Y} we denote the target domain. **Classification problems** are those which aim to assign each input vector x to one of a finite number of discrete categories of y . **Regression problems** are those where the output y consists of one or more continuous variables. All in all, the learner wishes to discover an unknown pattern (i.e. functional relationship) between components $x \in \mathcal{X}$ that serves as inputs and components $y \in \mathcal{Y}$ that act as outputs; i.e. $x \mapsto y$. Hence, \mathcal{X} is the input domain, and \mathcal{Y} is the output (or target) domain. The goal of learning is to discover a function which predicts (or help us make decisions about) $y \in \mathcal{Y}$ from $x \in \mathcal{X}$.

Unsupervised learning problems involve applications where the training data $z \in \mathcal{Z}$ consist of a set of input vectors $x \in \mathcal{X}$ without any corresponding target values ; i.e. $\mathcal{Z} = \mathcal{X}$. In clustering the goal is to discover groups of similar examples within the data of it is to discover groups of similar examples within the data.

2. NOTATION & DEFINITIONS IN LEARNING

Note 1. The aim is to learn the mapping $x \rightarrow y$ where $x \in \mathcal{X}$ and $y \in \mathcal{Y}$ with purpose to predict $y \in \mathcal{Y}$ from $x \in \mathcal{X}$. Denote $z = (x, y)^\top$ and $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$.

Definition 2. The learner's output is a function, $h : \mathcal{X} \rightarrow \mathcal{Y}$ which predicts $y \in \mathcal{Y}$ from $x \in \mathcal{X}$. It is also called hypothesis, prediction rule, predictor, or classifier.

Notation 3. We often denote the set of hypothesis as \mathcal{H} ; i.e. $h \in \mathcal{H}$.

Example 4. (Linear Regression)¹ Consider the regression problem where the goal is to learn the mapping $x \rightarrow y$ where $x \in \mathcal{X} \subseteq \mathbb{R}^d$ and $y \in \mathcal{Y} \subseteq \mathbb{R}$. A hypothesis is a linear function $h : \mathcal{X} \rightarrow \mathcal{Y}$ (that learner wishes to learn) with $h(x) = \langle w, x \rangle$ approximating the mapping $x \rightarrow y$. The hypothesis set is $\mathcal{H} = \{x \rightarrow \langle w, x \rangle : w \in \mathbb{R}^d\}$.

Example 5. (Binary Classification) Consider the classification problem where the goal is to learn the mapping $x \rightarrow y$ where $x \in \mathcal{X} \subseteq \mathbb{R}^d$ and $y \in \mathcal{Y} = \{-1, +1\}$. A hypothesis can be a function $h : \mathcal{X} \rightarrow \mathcal{Y}$ with $h(x) = \text{sign}(\langle w, x \rangle)$ approximating the mapping $x \rightarrow y$. The hypothesis set $\mathcal{H} = \{x \rightarrow \text{sign}(\langle w, x \rangle) : w \in \mathbb{R}^d\}$.

Definition 6. (Loss function) Given any set of hypothesis \mathcal{H} and some domain $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, a loss function $\ell(\cdot)$ is any function $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}_+$. Loss function $\ell(h, z)$ for $h \in \mathcal{H}$ and $z \in \mathcal{Z}$ is specified according to the purpose the machine learning algorithm. It reflects how the “error” is quantified for a given hypothesis h and a given example z . The rule is “the greater the error the greater the value of the loss”.

Example 7. (Cont. Example 4) In regression problems $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ and $\mathcal{Y} \subset \mathbb{R}$ is uncountable, a potential loss function is

$$\ell_{\text{sq}}(h, (x, y)) = (h(x) - y)^2$$

Example 8. (Cont. Example 5) In binary classification problems with hypothesis $h : \mathcal{X} \rightarrow \mathcal{Y}$ where $\mathcal{Y} = \{0, 1\}$ is discrete, a loss function can be

$$\ell_{0-1}(h, (x, y)) = 1(h(x) \neq y),$$

Definition 9. A learning problem with hypothesis class \mathcal{H} , examples domain \mathcal{Z} , and loss function ℓ may be denoted with a triplet $(\mathcal{H}, \mathcal{Z}, \ell)$.

Example 10. The standard multiple linear regression problem with regressors $x \in \mathcal{X} \subseteq \mathbb{R}^d$ and response $y \in \mathcal{Y} \subseteq \mathbb{R}$, is a learning problem with examples domain $\mathcal{Z} = \mathcal{X} \times \mathcal{Y} = \mathbb{R}^{d+1}$, hypothesis class $\mathcal{H} = \{x \rightarrow \langle w, x \rangle : w \in \mathbb{R}^d\}$, and loss function $\ell_{\text{sq}}(h, (x, y)) = (h(x) - y)^2$.

¹ $\langle w, x \rangle = w^\top x$

Example 11. The binary classification regression problem with regressors $x \in \mathcal{X} \subseteq \mathbb{R}^d$ and response $y \in \mathcal{Y} = \{-1, +1\}$, is a learning problem with examples domain $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, hypothesis class $\mathcal{H} = \{x \rightarrow \text{sign}(\langle w, x \rangle) : w \in \mathbb{R}^d\}$, and loss function $\ell_{0-1}(h, (x, y)) = 1(h(x) \neq y) = 1(\langle w, x \rangle y < 0)$.

Definition 12. Data generation model $g(\cdot)$ is the probability distribution over $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, unknown to the learner which describes the probabilistic law that realizations $z = (x, y) \in \mathcal{Z}$ are realized.

Conventions...

Note 13. If the Hypothesis set \mathcal{H} is a known parametric family of functions; i.e. $\mathcal{H} = \{h_w(\cdot) ; w \in \mathcal{W}\}$ parameterized by unknown $w \in \mathcal{W}$, then we can equivalently consider the convention $\mathcal{H} = \{w \in \mathcal{W}\} = \mathcal{W}$ keeping in mind that the learner's output is restricted to formula $h_w(\cdot)$.

Example 14. Because it involves only linear functions as predictors $h_w(x) = \langle w, x \rangle$, we could consider a hypothesis class $\mathcal{H} = \{w \in \mathbb{R}^d\} = \mathbb{R}^d$ and loss function loss $\ell(w, (x, y)) = (\langle w, x \rangle - y)^2$ for simplicity.

Example 15. Because it involves only linear functions as predictors $h_w(x) = \text{sign}(\langle w, x \rangle)$, we could consider a hypothesis class $\mathcal{H} = \{w \in \mathbb{R}^d\} = \mathbb{R}^d$ and loss function loss $\ell(w, (x, y)) = 1(\langle x, w \rangle y < 0)$ for simplicity.

2.1. Learning...

Definition 16. (Risk function) The risk function $R_g(h)$ of h is the expected loss of the hypothesis $h \in \mathcal{H}$, w.r.t. the data generation model (which is a probability distribution) g over domain Z ; i.e.

$$(2.1) \quad R_g(h) = \mathbb{E}_{z \sim g}(\ell(h, z))$$

Definition 17. (Risk minimization learning) The Risk minimization (RM) learning paradigm computes the optimal predictor h^* as the minimizer of the risk $R_g(h)$ of h ; i.e.

$$(2.2) \quad h^* = \arg \min_h (R_g(h))$$

Example 18. (Cont. Ex. 7) The risk function is $R_g(h_w) = \mathbb{E}_{z \sim g}(h_w(x) - y)^2 = \mathbb{E}_{z \sim g}(\langle w, x \rangle - y)^2$, and it measures the quality of the hypothesis function $h : \mathcal{X} \rightarrow \mathcal{Y}$, (or equiv. the validity of the class of hypotheses \mathcal{H}) against the data generating model g , as the expected square difference between the predicted values form h and the true target values y at every x .

Example 19. (Cont. Ex.) The risk function is $R_g(h_w) = \mathbb{E}_{z \sim g}(1(h_w(x) \neq y)) = \Pr_{z \sim g}(\langle x, w \rangle y < 0)$.

Note 20. Computing the risk minimizer may be practically challenging due to the integration w.r.t. the unknown data generation model g involved in the expectation (2.1). Sub-optimally, instead of the Risk function in (2.2), one use the Empirical risk function as a Monte Carlo approximation of it.

Definition 21. Training data set \mathcal{S} of size m is any finite sequence of pairs $(z_i = (x_i, y_i) ; i = 1, \dots, m)$, called examples, in $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$; i.e. $\mathcal{S} = \{(x_i, y_i) ; i = 1, \dots, m\}$. This is the information that the learner has assess.

Definition 22. (Empirical risk function) The Empirical Risk Function (ERF) $\hat{R}_S(h)$ of h is the expectation of loss of h over a given sample $S = (z_1, \dots, z_m) \in \mathcal{Z}^m$; i.e.

$$\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, z_i).$$

Definition 23. (Empirical risk minimization learning) The Empirical risk minimization (ERM) learning paradigm computes the optimal predictor h^* as the minimizer of the ERF $\hat{R}_S(h)$ of h ; i.e.

$$(2.3) \quad h^* = \arg \min_h (\hat{R}_S(h))$$

Example 24. (Cont. Example 18) Given given sample $S = \{(x_i, y_i); i = 1, \dots, m\}$ the empirical risk function is $\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i)^2 = \frac{1}{m} \sum_{i=1}^m (\langle x_i^\top w - y_i \rangle^2)$.

Example 25. (Cont. Example 8) Given given sample $S = \{(x_i, y_i); i = 1, \dots, m\}$ the empirical risk function is $\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}(h(x_i) \neq y_i) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}(\langle x_i, w \rangle y_i < 0)$.

Definition 26. We denote as $\mathfrak{A}(\mathcal{S})$ the hypothesis (outcome) that a learning algorithm \mathfrak{A} returns given training sample S .

3. REGULARIZED LOSS MINIMIZATION (RLM) LEARNING

Note 27. Consider a learning problem $(\mathcal{H}, \mathcal{Z}, \ell)$ with training data set \mathcal{S} .

Definition 28. (Regularized loss minimization learning) Regularized loss minimization learning paradigm computes the optimal rule $h^* \in \mathcal{H}$ by jointly minimizing the Risk function $R_g(h)$ (or Empirical Risk Function $\hat{R}_S(h)$) and a regularization function $J : \mathcal{H} \rightarrow \mathbb{R}$ where $J(h)$ increases in value with the complexity of the hypothesis $h \in \mathcal{H}$. Formally, the Regularized loss minimization rule h^* is

$$h^* = \arg \min_{h \in \mathcal{H}} (R_g(h) + J(h)), \text{ given the Risk function}$$

$$h^* = \arg \min_{h \in \mathcal{H}} (\hat{R}_S(h) + J(h)), \text{ given the Empirical risk function}$$

3.1. ℓ_0 , ℓ_1 , and ℓ_2 norm regularization problems.

Note 29. Consider a parameterized hypothesis $h_w(\cdot)$ (e.g. $h_w(x) = \eta(\langle x, w \rangle)$) with $\eta(\cdot)$ a function and $w \in \mathbb{R}^d$. Consider the hypothesis class in the simplified form $\mathcal{H} = \mathcal{W} = \{w \in \mathbb{R}^d\} = \mathbb{R}^d$ of Note 13. Assume we want the vector w to be sparse, in the sense that “sparser” means more zero elements.

3.1.1. ℓ_0 norm regularization.

Note 30. The problem of minimizing the empirical risk subject to a budget of k features can be written as

$$\underset{w \in \mathcal{H}}{\text{minimize}} R(w)$$

$$\text{subject to } \|w\|_0 \leq k_0$$

where $\|w\|_0 = \#\{j : w_j \neq 0\}$. By Lagrange multiplies we can re-write it as

$$w^* = \arg \min_{w \in \mathcal{H}} (R(w) + \lambda_0 \|w\|_0)$$

for some $\lambda_0 \geq 0$ (k_0 dependent). Hence, $h^*(x) = h_{w^*}(x)$.

Example 31. (Cont. Example 14) The ℓ_0 -RLM rule in Example 14 becomes

$$w^* = \arg \min_w \left(\mathbb{E}_{z \sim g} (\langle w, x_i \rangle - y)^2 + \lambda_0 \|w\|_0 \right)$$

using the Risk function, and

$$w^* = \arg \min_w \left(\frac{1}{m} \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2 + \lambda_0 \|w\|_0 \right)$$

using the Empirical risk function.

3.1.2. ℓ_1 norm regularization (LASSO).

Note 32. Solving ℓ_0 norm optimization problem is computationally hard. To simplify computations, we can replace $J(\cdot) = \|\cdot\|_0$ with $J(\cdot) = \|\cdot\|_1$ i.e.

$$\begin{aligned} & \underset{w \in \mathcal{H}}{\text{minimize}} R(w) \\ & \text{subject to } \|w\|_1 \leq k_1 \end{aligned}$$

where $\|w\|_1 = \sum_j |w_j|$. By Lagrange multiplies we can re-write it as

$$w^* = \arg \min_w (R(w) + \lambda_1 \|w\|_1)$$

for some $\lambda_1 \geq 0$ (k_1 dependent). Hence, $h^*(x) = h_{w^*}(x)$.

Example 33. (Cont. Example 14) The ℓ_1 -RLM rule in Example 14 becomes

$$w^* = \arg \min_w \left(\mathbb{E}_{z \sim g} (\langle w, x \rangle - y)^2 + \lambda_1 \sum_{j=1}^d |w_j| \right)$$

using the Risk function , and

$$w^* = \arg \min_w \left(\frac{1}{m} \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2 + \lambda_1 \sum_{j=1}^d |w_j| \right)$$

using the Empirical risk function.

3.1.3. ℓ_2 norm regularization (Ridge).

Note 34. To further simplify computations, we can consider $J(\cdot) = \|\cdot\|_2$ i.e.

$$\begin{aligned} & \underset{w}{\text{minimize}} R(w) \\ & \text{subject to } \|w\|_2^2 \leq k_2 \end{aligned}$$

where $\|w\|_2^2 = \sum_j w_j^2$. By Lagrange multipliers we can re-write it as

$$w^* = \arg \min_w \left(R(w) + \lambda_2 \|w\|_2^2 \right)$$

for some $\lambda_1 \geq 0$ (k_1 dependent). Hence, $h^*(x) = h_{w^*}(x)$.

Example 35. (Cont. Example 14) The ℓ_2 -RLM rule in Example 14 becomes

$$w^* = \arg \min_w \left(\mathbb{E}_{z \sim g} (\langle w, x \rangle - y)^2 + \lambda_2 \sum_{j=1}^d w_j^2 \right)$$

using the Risk function, and

$$w^* = \arg \min_w \left(\frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i)^2 + \lambda_2 \sum_{j=1}^d w_j^2 \right)$$

using the Empirical risk function.

Example 36. (Cont. Example 15) The ℓ_2 -RLM rule in Example 15 becomes

$$w^* = \arg \min_w \left(\Pr_{z \sim g} (\langle x, w \rangle y < 0) + \lambda_2 \sum_{j=1}^d w_j^2 \right)$$

using the Risk function, and

$$w^* = \arg \min_w \left(\frac{1}{m} \sum_{i=1}^m \mathbb{1}(\langle x_i, w \rangle y_i < 0) + \lambda_2 \sum_{j=1}^d w_j^2 \right)$$

using the Empirical risk function.

APPENDIX A. FURTHER EXAMPLES

Example 37. Consider a learning problem where the true data generation distribution (unknown to the learner) is $g(z)$, the statistical model (known to the learner) is given by a sampling distribution $f_\theta(y) := f(y|\theta)$ labeled by an unknown parameter θ . The goal is to learn θ . If we assume loss function

$$\ell(\theta, z) = \log\left(\frac{g(z)}{f_\theta(z)}\right)$$

then the risk is

$$(A.1) \quad R_g(\theta) = E_{z \sim g}\left(\log\left(\frac{g(z)}{f_\theta(z)}\right)\right) = E_{z \sim g}(\log(g(z))) - E_{z \sim g}(\log(f_\theta(z)))$$

whose minimizer is

$$\theta^* = \arg \min_{\forall \theta} (R_g(\theta)) = \arg \min_{\forall \theta} (E_{z \sim g}(-\log(f_\theta(z))))$$

as the first term in (A.1) is constant. Note that in the Maximum Likelihood Estimation technique the MLE θ_{MLE} is the minimizer

$$\theta_{MLE} = \arg \min_{\theta} \left(\frac{1}{m} \sum_{i=1}^m (-\log(f_\theta(z_i))) \right)$$

where $S = \{z_1, \dots, z_m\}$ is an IID sample from g . Hence, MLE θ_{MLE} can be considered as the minimizer of the empirical risk $R_S(\theta) = \frac{1}{m} \sum_{i=1}^m (-\log(f_\theta(z_i)))$.

Lecture notes 2: Elements of convex learning problems

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim. To introduce convex learning problems that can be used as a framework for the analysis of stochastic gradient related learning algorithms.

Reading list & references:

- Shalev-Shwartz, S., & Ben-David, S. (2014). Understanding machine learning: From theory to algorithms. Cambridge university press.
 - Ch. 12 Convex Learning Problems

Further reading

- Bishop, C. M. (2006). Pattern recognition and machine learning. New York: Springer.

1. MOTIVATIONS

Note 1. We introduce learning problems associated with convexity, Lipschitzness and smoothness to be able to analyze and understand the machine learning (ML) tools we will discuss later on (eg stochastic gradient descent, SVM). We discuss ways allowing to address non-convex ML problems (eg, Artificial neural networks, Gaussian process regression) in the convex setting.

2. CONVEXITY

Note 2. Convexity is a central concept in learning, e.g. least squares, as it often considers a Euclidean distance as a Risk function to be minimized.

Definition 3. A set C is convex if for any $u, v \in C$ and for any $\alpha \in [0, 1]$ we have that $\alpha u + (1 - \alpha) v \in C$.

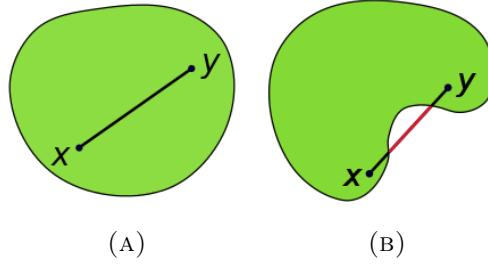


FIGURE 2.1. (2.1a) is a Convex set ; (2.1b) is a non-convex set

Example 4. For instance \mathbb{R}^d for $d \geq 1$ is a convex set.

Definition 5. Let C be a convex set. A function $f : C \rightarrow \mathbb{R}$ is convex function if for any $u, v \in C$ and for any $\alpha \in [0, 1]$

$$f(\alpha u + (1 - \alpha)v) \leq \alpha f(u) + (1 - \alpha)f(v)$$

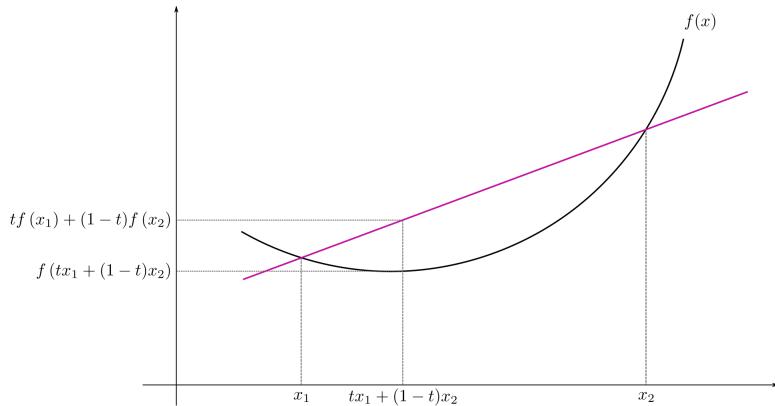


FIGURE 2.2. A convex function

Example 6. The function $f : \mathbb{R}^d \rightarrow \mathbb{R}_+$ with $f(x) = \|x\|_2^2$ is convex function. For any $u, v \in C$ and for any $\alpha \in [0, 1]$ it is

$$\begin{aligned} & \|\alpha u + (1 - \alpha)v\|_2^2 - \alpha \|u\|_2^2 - (1 - \alpha)\|v\|_2^2 \\ &= (\alpha u + (1 - \alpha)v)^\top (\alpha u + (1 - \alpha)v) - \alpha \|u\|_2^2 - (1 - \alpha)\|v\|_2^2 \\ &= \dots = -\alpha(1 - \alpha)\|u - v\|_2^2 \leq 0 \end{aligned}$$

Note 7. Every local minimum of a convex function is the global minimum.

Note 8. Let $f : C \rightarrow \mathbb{R}$ be convex function. The tangent of f at $w \in C$ is below f , namely

$$\forall u \in C \quad f(u) \geq f(w) + \langle \nabla f(w), u - w \rangle$$

Proposition 9. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ such that $f(w) = g(\langle w, x \rangle + y)$ for some $x \in \mathbb{R}^d$, $y \in \mathbb{R}$. If g is convex function then f is convex function.

Proof. See Exercise 1 in the Exercise sheet. \square

Example 10. Consider the regression problem with regressor $x \in \mathbb{R}^d$, and response $y \in \mathbb{R}$ and predictor rule $h(x) = \langle w, x \rangle$. The loss function $\ell(w, (x, y)) = (\langle w, x \rangle - y)^2$ is convex because $g(a) = a^2$ is convex and Proposition 9.

Note 11. Let $f_j : \mathbb{R}^d \rightarrow \mathbb{R}$ convex functions for $j = 1, \dots, r$. Then:

- (1) $g(x) = \max_{\forall j} (f_j(x))$ is a convex function
- (2) $g(x) = \sum_{j=1}^r w_j f_j(x)$ is a convex function where $w_j > 0$

Proof.

- (1) For any $u, v \in \mathbb{R}^d$ and for any $\alpha \in [0, 1]$

$$\begin{aligned} g(\alpha u + (1 - \alpha)v) &= \max_{\forall j} (f_j(\alpha u + (1 - \alpha)v)) \\ &\leq \max_{\forall j} (\alpha f_j(u) + (1 - \alpha)f_j(v)) && (f_j \text{ is convex}) \\ &\leq \alpha \max_{\forall j} (f_j(u)) + (1 - \alpha) \max_{\forall j} (f_j(v)) && (\max(\cdot) \text{ is convex}) \\ &\leq \alpha g(u) + (1 - \alpha)g(v) \end{aligned}$$

- (2) For any $u, v \in \mathbb{R}^d$ and for any $\alpha \in [0, 1]$

$$\begin{aligned} g(\alpha u + (1 - \alpha)v) &= \sum_{j=1}^r w_j f_j(\alpha u + (1 - \alpha)v) \\ &\leq \alpha \sum_{j=1}^r w_j f_j(u) + (1 - \alpha) \sum_{j=1}^r w_j f_j(v) && (f_j \text{ is convex}) \\ &\leq \alpha g(u) + (1 - \alpha)g(v) \end{aligned}$$

\square

Example 12. $g(x) = \|x\|_1$ is convex according to Note 11, since $|x_j| = \max(-x_j, x_j)$ and $\|x\|_1 = \sum_j |x_j|$.

3. STRONG CONVEXITY

Note 13. Strong convexity is a central concept in Ridge regularization, as it makes a convex loss function strongly convex by adding a shrinkage term.

Definition 14. (Strongly convex functions) A function f is λ -strongly convex function is for all w, u , and $\alpha \in (0, 1)$ we have

$$(3.1) \quad f(\alpha w + (1 - \alpha) u) \leq \alpha f(w) + (1 - \alpha) f(u) - \frac{\lambda}{2} \alpha (1 - \alpha) \|w - u\|^2$$

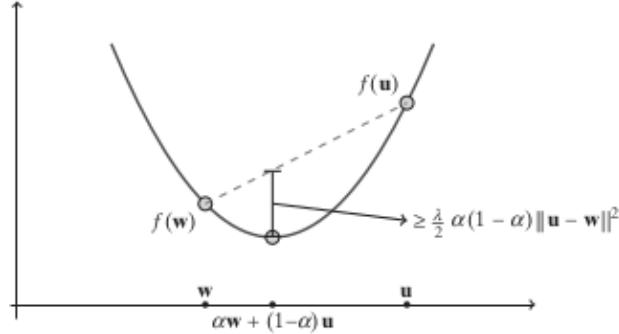


FIGURE 3.1. Strongly convex function

Note 15. The following can be checked from the definition by substitution

- (1) The function $f(w) = \lambda \|w\|^2$ is 2λ -strongly convex
- (2) If f is λ -strongly convex and g is convex then $f + g$ is λ -strongly convex

4. LIPSCHITZNESS

Definition 16. Let $C \subseteq \mathbb{R}^d$. Function $f : C \rightarrow \mathbb{R}^k$ is ρ -Lipschitz over C if for every $w_1, w_2 \in C$ we have that

$$(4.1) \quad \|f(w_1) - f(w_2)\| \leq \rho \|w_1 - w_2\|. \quad \text{Lipschitz condition}$$

Note 17. That means: a Lipschitz function $f(x)$ cannot change too drastically wrt x .

Example 18. Consider the function $f : \mathbb{R} \rightarrow \mathbb{R}_+$ with $f(x) = x^2$.

- (1) f is not a ρ -Lipschitz in \mathbb{R} .
- (2) f is a ρ -Lipschitz in $C = \{x \in \mathbb{R} : |x| < \rho/2\}$.

Solution.

- (1) For $x_1 = 0$ and $x_2 = 1 + \rho$, it is

$$|f(x_2) - f(x_1)| = (1 + \rho)^2 > \rho(1 + \rho) = \rho|x_2 - x_1|$$

- (2) It is

$$|f(x_2) - f(x_1)| = |x_2^2 - x_1^2| = |(x_2 + x_1)(x_2 - x_1)| \leq 2\rho/2(x_2 - x_1) = \rho|x_2 - x_1|$$

Note 19. Let functions g_1 be ρ_1 -Lipschitz and g_2 be ρ_2 -Lipschitz. Then f with $f(x) = g_1(g_2(x))$ is $\rho_1\rho_2$ -Lipschitz.

[See Exercise 2 from the exercise sheet]

Example 20. Let functions g be ρ -Lipschitz. Then f with $f(x) = g(\langle v, x \rangle + b)$ is $(\rho|v|)$ -Lipschitz.

Solution. It is

$$\begin{aligned} |f(w_1) - f(w_2)| &= |g(\langle v, w_1 \rangle + b) - g(\langle v, w_2 \rangle + b)| \leq \rho |\langle v, w_1 \rangle + b - \langle v, w_2 \rangle - b| \\ &\leq \rho |v^\top w_1 - v^\top w_2| \leq \rho |v| |w_1 - w_2| \end{aligned}$$

Note 21. So, given Examples 18 and 20, in the linear regression setting using loss $\ell(w, z = (x, y)) = (w^\top x - y)^2$, the loss function is β -Lipschitz for a given $z = (x, y)$ and bounded $\|w\| < \rho$.

5. SMOOTHNESS

Definition 22. A differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is β -smooth if its gradient is β -Lipschitz; namely for all $v, w \in \mathbb{R}^d$

$$(5.1) \quad \|\nabla f(w_1) - \nabla f(w_2)\| \leq \beta \|w_1 - w_2\|.$$

Note 23. Function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is β -smooth iff

$$(5.2) \quad f(v) \leq f(w) + \langle \nabla f(w), v - w \rangle + \frac{\beta}{2} \|v - w\|^2$$

Note 24. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ with $f(w) = g(\langle w, x \rangle + y)$ $x \in \mathbb{R}^d$ and $y \in \mathbb{R}$. Let $g : \mathbb{R} \rightarrow \mathbb{R}$ be a β -smooth function. Then f is a $(\beta \|x\|^2)$ -smooth.

[See Exercise 3 from the Exercise sheet]

Example 25. Let $f(w) = (\langle w, x \rangle + y)^2$ for $w \in \mathbb{R}^d$ and $y \in \mathbb{R}$. Then f is $(2\|x\|^2)$ -smooth.

To show that we set $f(w) = g(\langle w, x \rangle + y)$ where $g(a) = a^2$. g is 2-smooth since

$$\|g'(w_1) - g'(w_2)\| = \|2w_1 - 2w_2\| \leq 2 \|w_1 - w_2\|.$$

Hence from Theorem 24, f is $(2\|x\|^2)$ -smooth.

Example 26. Consider the regression problem with predictor rule $h(x) = \langle w, x \rangle$, loss function $\ell(w, (x, y)) = (\langle w, x \rangle - y)^2$, feature $x \in \mathbb{R}^d$, and target $y \in \mathbb{R}$. Then $\ell(w, \cdot)$ is $(2\|x\|^2)$ -smooth.

[Follows from Example 25.]

6. CONVEX LEARNING PROBLEMS

Definition 27. Convex learning problem is a learning problem $(\mathcal{H}, \mathcal{Z}, \ell)$ that the hypothesis class \mathcal{H} is a convex set, and the loss function $\ell(\cdot, z)$ is a convex function for each example $z \in \mathcal{Z}$.

Example 28. Consider the regression problem with predictor rule $h(x) = \langle w, x \rangle$, loss function $\ell(w, (x, y)) = (\langle w, x \rangle - y)^2$, feature $x \in \mathbb{R}^d$, and target $y \in \mathbb{R}$. This imposes a convex learning problem due to Examples 4 and 11.

Definition 29. Convex-Lipschitz-Bounded Learning Problem $(\mathcal{H}, \mathcal{Z}, \ell)$ with parameters ρ , and B , is called the learning problem whose the hypothesis class \mathcal{H} is a convex set, for all $w \in \mathcal{H}$ it is $\|w\| \leq B$, and the loss function $\ell(\cdot, z)$ is convex and ρ -Lipschitz function for all $z \in \mathcal{Z}$.

Example 30. Consider the regression problem with predictor rule $h(x) = \langle w, x \rangle$, loss function $\ell(w, (x, y)) = (\langle w, x \rangle - y)^2$, feature $x \in \mathbb{R}^d$, and target $y \in \mathbb{R}$. This imposes a Convex-Lipschitz-Bounded Learning Problem if $\mathcal{H} = \{w \in \mathbb{R}^d : \|w\| \leq B\}$ due to Examples 11, and 18(2).

Definition 31. Convex-Smooth-Bounded Learning Problem $(\mathcal{H}, \mathcal{Z}, \ell)$ with parameters β , and B , is called the learning problem whose the hypothesis class \mathcal{H} is a convex set, for all $w \in \mathcal{H}$ it is $\|w\| \leq B$, and the loss function $\ell(\cdot, z)$ is convex, nonnegative, and β -smooth function for all $z \in \mathcal{Z}$.

Example 32. Consider the regression problem with predictor rule $h(x) = \langle w, x \rangle$, loss function $\ell(w, (x, y)) = (\langle w, x \rangle - y)^2$, feature $x \in \mathbb{R}^d$, and target $y \in \mathbb{R}$. This imposes a Convex-Smooth-Bounded Learning Problem if $\mathcal{H} = \{w \in \mathbb{R}^d : \|w\| \leq B\}$ due to Examples 11, and 26.

Note 33. (Empirical risk minimization learning problem) If ℓ is a convex loss function and the class \mathcal{H} is convex, then the ERM $_{\mathcal{H}}$ problem, of minimizing the empirical risk $\hat{R}_{\mathcal{S}}(w)$ over \mathcal{H} , is a convex optimization problem

Proof. The Empirical risk minimization (ERM) problem $(\mathcal{H}, \mathcal{Z}, \ell)$ is

$$w^* = \arg \min_{w \in \mathcal{H}} \left\{ \hat{R}_{\mathcal{S}}(w) \right\}$$

given a sample $\mathcal{S} = \{z_1, \dots, z_m\}$ for $\hat{R}_{\mathcal{S}}(w) = \frac{1}{m} \sum_{i=1}^m \ell(w, z_i)$. $\hat{R}_{\mathcal{S}}(w)$ is a convex function from Note (11). Hence ERM rule is a problem of minimizing a convex function subject to the constraint that the solution should be in a convex set. \square

Example 34. Multiple linear regression with predictor rule $h(x) = \langle w, x \rangle$, loss function $\ell(w, (x, y)) = (\langle w, x \rangle - y)^2$, feature $x \in \mathbb{R}^d$, and target $y \in \mathbb{R}$ where

$$w^* = \arg \min_w \text{E} (\langle w, x \rangle - y)^2$$

or

$$w^{**} = \arg \min_w \frac{1}{m} \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2$$

is a convex learning problem –from Note 33.

Note 35. (Ridge / ℓ_2 -RLM problem) If ℓ is a convex loss function w.r.t. w , the class \mathcal{H} is convex, and $J(\cdot; \lambda) = \lambda \|\cdot\|_2^2$ with $\lambda > 0$ then $\hat{R}_{\mathcal{S}}(w) + \lambda \|w\|_2^2$ is a 2λ -strongly convex function, and hence the ℓ_2 -RLM problem

$$w^* = \arg \min_{w \in \mathcal{H}} \left\{ \hat{R}_{\mathcal{S}}(w) + \lambda \|w\|_2^2 \right\}$$

is a strongly convex optimization problem (i.e. the learning rule is the minimizer of a strongly convex function over a convex set).

Proof. $\hat{R}_{\mathcal{S}}(\cdot)$ is a convex function from Note 33, $\lambda \|\cdot\|_2^2$ is 2λ -strongly convex, hence $\hat{R}_{\mathcal{S}}(w) + \lambda \|w\|_2^2$ is a 2λ -strongly convex function. Hence the above Ridge RLM problem is a strongly convex optimization problem. \square

7. NON-CONVEX LEARNING PROBLEMS (SURROGATE TREATMENT)

Note 36. A learning problem may involve non-convex loss function $\ell(w, z)$ w.r.t. w which implies a non-convex risk function $R_g(w)$. A suitable treatment to analyze the associated learning tools within the convex setting would be to upper bound the non-convex loss function $\ell(w, z)$ by a convex surrogate loss function $\tilde{\ell}(w, z)$ for all w , and use $\tilde{\ell}(w, z)$ instead of $\ell(w, z)$.

Example 37. Consider the binary classification problem with inputs $x \in \mathcal{X}$, outputs $y \in \{-1, +1\}$; we need to learn $w \in \mathcal{H}$ from hypothesis class $\mathcal{H} \subset \mathbb{R}^d$ with respect to the loss

$$\ell(w, (x, y)) = 1_{(y\langle w, x \rangle \leq 0)}$$

with $y \in \mathbb{R}$, and $x \in \mathbb{R}^d$. Here $\ell(\cdot)$ is non-convex. A convex surrogate loss function can be

$$\tilde{\ell}(w, (x, y)) = \max(0, 1 - y\langle w, x \rangle)$$

which is convex (Example 11) wrt w . Note that:

- $\tilde{\ell}(w, (x, y))$ is convex wrt w ; because $\max(\cdot)$ is convex
- $\ell(w, (x, y)) \leq \tilde{\ell}(w, (x, y))$ for all $w \in \mathcal{H}$

Then we can compute

$$\tilde{w}_* = \arg \min_{\forall x} (\tilde{R}_g(w)) = \arg \min_{\forall x} (\mathbb{E}_{(x,y) \sim g} (\max(0, 1 - y\langle w, x \rangle)))$$

instead of

$$w_* = \arg \min_{\forall x} (R_g(w)) = \arg \min_{\forall x} (\mathbb{E}_{(x,y) \sim g} (1_{(y\langle w, x \rangle \leq 0)}))$$

Of course by using the surrogate loss instead of the actual one, we introduce some approximation error in the produced output $\tilde{w}_* \neq w_*$.

Note 38. (Intuitions...) Using a convex surrogate loss function instead the convex one, facilitates computations but introduces extra error to the solution. If $R_g(\cdot)$ is the risk under the non-convex loss, $\tilde{R}_g(\cdot)$ is the risk under the convex surrogate loss, and \tilde{w}_{alg} is the output of the learning algorithm under $\tilde{R}_g(\cdot)$ then we have the upper bound

$$R_g(\tilde{w}_{\text{alg}}) \leq \underbrace{\min_{w \in \mathcal{H}} (R_g(w))}_{\text{I}} + \underbrace{\left(\min_{w \in \mathcal{H}} (\tilde{R}_g(w)) - \min_{w \in \mathcal{H}} (R_g(w)) \right)}_{\text{II}} + \underbrace{\epsilon}_{\text{III}}$$

where term I is the approximation error measuring how well the hypothesis class performs on the generating model, term II is the optimization error due to the use of surrogate loss instead of the actual non-convex one, and term III is the estimation error due to the use of a training set and not the whole generation model.

Lecture notes 3: Learnability and stability in learning problems

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim. To introduce concepts PAC, fitting vs stability trade off, stability, and their implementation in regularization problems and convex problems.

Reading list & references:

- Shalev-Shwartz, S., & Ben-David, S. (2014). Understanding machine learning: From theory to algorithms. Cambridge university press.
– Ch. 2, 3, 13
- Vapnik, V. (2000). The nature of statistical learning theory. Springer science & business media.

1. LEARNABLE PROBABLY APPROXIMATELY CORRECT (PAC) LEARNING

Note 1. We formally define the broad learning problem we will work on.

Note 2. Learning algorithms \mathfrak{A} use training data sets \mathcal{S} which may miss characteristics of the unknown data-generating process g . Essentially, approximations (aka errors) are inevitable; some characteristics of data generating process g will be missed even if we use a very representative training set \mathcal{S} . We cannot hope that the learning algorithm will find a hypothesis whose error is smaller than the minimal possible error.

Note 3. The PAC learning problem requires no prior assumptions about the data-generating process g . It requires that the learning algorithm \mathfrak{A} will find a predictor $\mathfrak{A}(\mathcal{S})$ whose error is not much larger than the best possible error of a predictor in some given benchmark hypothesis class. So essentially, in practice, the researcher's effort falls on the hypothesis class \mathcal{H} .

Definition 4. (Agnostic PAC Learnability for General Loss Functions) A hypothesis class \mathcal{H} is agnostic PAC learnable with respect to a domain \mathcal{Z} and a loss function $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}_+$, if there exist a function $m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$ and a learning algorithm \mathfrak{A} with the following property:

- for every $\epsilon \in (0, 1)$, $\delta \in (0, 1)$, and distribution g over \mathcal{Z} , when running algorithm \mathfrak{A} given training set $\mathcal{S}_m = \{z_1, \dots, z_m\}$ with $z_i \stackrel{\text{IID}}{\sim} g$ for $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ then \mathfrak{A} returns

$\mathfrak{A}(\mathcal{S}) \in \mathcal{H}$ such as

$$(1.1) \quad \Pr_{\mathcal{S} \sim g} \left(R_g(\mathfrak{A}(\mathcal{S}_m)) - \min_{h' \in \mathcal{H}} (R_g(h')) \leq \epsilon \right) \geq 1 - \delta$$

Note 5. About (1.1): The accuracy parameter ϵ determines how far the output rule $\mathfrak{A}(\mathcal{S}_m)$ of \mathfrak{A} can be from the optimal rule (this corresponds to the ‘approximately correct’ part of “PAC”). The confidence parameter δ indicates how likely the classifier is to meet that accuracy requirement (corresponds to the “probably” part of “PAC”).

Note 6. It may be easier to work with expectations by using Markov inequality in (1.1).

Proposition 7. *Let \mathfrak{A} be a learning algorithm that guarantees the following:*

- If $m \geq m_{\mathcal{H}}(\epsilon)$ then for every distribution g , it is

$$E_{\mathcal{S} \sim g}(R_g(\mathfrak{A}(\mathcal{S}_m))) \leq \min_{h' \in \mathcal{H}} (R_g(h')) + \epsilon$$

Then \mathfrak{A} satisfies the PAC guarantee in Definition 4:

- for every $\delta \in (0, 1)$, if $m \geq m_{\mathcal{H}}(\epsilon\delta)$ then

$$\Pr_{\mathcal{S} \sim g} \left(R_g(\mathfrak{A}(\mathcal{S}_m)) - \min_{h' \in \mathcal{H}} (R_g(h')) \leq \epsilon \right) \geq 1 - \delta$$

Proof. Let $\xi = R_g(\mathfrak{A}(\mathcal{S}_m)) - \min_{h' \in \mathcal{H}} (R_g(h'))$. From Markov’s inequality $\Pr(\xi \geq E(\xi)/\delta) \leq \frac{1}{E(\xi)/\delta} E(\xi) = \delta$. Namely, $\Pr(\xi \leq E(\xi)/\delta) = 1 - \delta$. But it is given that if $m \geq m_{\mathcal{H}}((\epsilon\delta))$ for every distribution g it is $E(\xi) \leq (\epsilon\delta)$. So by substitution $\Pr(\xi \leq (\epsilon\delta)/\delta) \geq 1 - \delta$ implies $\Pr(\xi \leq \epsilon) \geq 1 - \delta$. Now substitute back ξ and we conclude the proof. \square

2. ANALYSIS OF THE RISK BASED ON THE TRADE-OFF FITTING VS STABILITY

Note 8. Let $R^* = \min_{\forall h} (R(h))$ be an ideal/optimal (hence minimum) Risk, and $\mathfrak{A}(\mathcal{S})$ the learning rule from a learning algorithm \mathfrak{A} trained against dataset \mathcal{S} . The Risk of a learning algorithm \mathfrak{A} can be decomposed as

$$(2.1) \quad R_g(\mathfrak{A}(\mathcal{S})) - R^* = \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) - R^* + \underbrace{R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))}_{\text{over-fitting}}$$

Note 9. Over-fitting can be represented by $R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))$. However, $\hat{R}_{\mathcal{S}}(\cdot)$ is a random variable and, here, for our computational convenience, we focus on its expectation w.r.t. $\mathcal{S} \sim g$. Hence, we provide the following (arguable) definition for over-fitting on which we base our analysis.

Definition 10. For a learning algorithm \mathfrak{A} , as a measure of over-fitting we consider the expected difference between true Risk and empirical Risk

$$(2.2) \quad E_{\mathcal{S} \sim g} \left(R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right)$$

Definition 11. We say that learning algorithm \mathfrak{A} suffers from over-fitting when (2.2) is ‘too’ large.

Note 12. The expected Risk of a learning algorithm $\mathfrak{A}(\mathcal{S})$ can be decomposed as

$$(2.3) \quad \mathbb{E}_{\mathcal{S} \sim g}(R_g(\mathfrak{A}(\mathcal{S}))) = \underbrace{\mathbb{E}_{\mathcal{S} \sim g}(\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})))}_{(I)} + \underbrace{\mathbb{E}_{\mathcal{S} \sim g}(R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})))}_{(II)}$$

by applying expectations in (2.1) and ignoring R^* . (I) indicates how well $\mathfrak{A}(\mathcal{S})$ fits the training set \mathcal{S} , and (II) indicates the discrepancy between the true and empirical risks of $\mathfrak{A}(\mathcal{S})$. In Section 3, we argue that the over-fitting term (II) is directly related to a certain type of stability of $\mathfrak{A}(\mathcal{S})$.

Note 13. The following result connects the need for upper bounding (and minimizing this bound) the Expected Risk in (2.3) and PAC learning (Definition 4).

Note 14. Hence, we aim to design a learning algorithm $\mathfrak{A}(\mathcal{S})$ that both fits the training set and is stable; i.e, keep both (I) and (II) in (2.3) small. As seen later, in certain learning problems, there may be a trade-off between empirical risk term (I) and (II) in (2.3). Consequently, we aim to upper bound (2.3) by upper bounding (I) and (II) individually and decide which term we should benefit against the other to achieve a certain total bound.

3. STABILITY AND ITS ASSOCIATION WITH OVER-FITTING

Notation 15. Consider a learning problem $(\mathcal{H}, \mathcal{Z}, \ell)$. Let $\mathcal{S} = \{z_1, \dots, z_m\}$ be a training sample, and \mathfrak{A} be a learning algorithm with output $\mathfrak{A}(\mathcal{S})$.

Note 16. It is reasonable to consider that a learning algorithm \mathfrak{A} can be stable if a small change of the algorithm input does not change the algorithm output much. Mathematically formalizing this, we can say that if $\mathcal{S}^{(i)} = \{z_1, \dots, z_{i-1}, z', z_{i+1}, \dots, z_m\}$ is another training dataset equal to \mathcal{S} but the i th element being replaced by another independent example $z' \sim g$, then a good learning algorithm \mathfrak{A} would produce a small value of

$$\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \geq 0$$

Definition 17. A learning algorithm \mathfrak{A} is **on-average-replace-one-stable** with rate a decreasing function $\epsilon : \mathbb{N} \rightarrow \mathbb{R}$ if for every distribution g

$$(3.1) \quad \mathbb{E}_{\substack{\mathcal{S} \sim g \\ z' \sim g, i \sim U\{1, \dots, m\}}} (\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i)) \leq \epsilon(m)$$

Note 18. The following result demonstrates the direct association of stability and over-fitting as defined in Definitions 11 & 17.

Theorem 19. For any learning algorithm \mathfrak{A} it is

$$(3.2) \quad E_{\mathcal{S} \sim g} \left(R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right) = E_{\substack{\mathcal{S} \sim g, z' \sim g \\ i \sim U\{1, \dots, m\}}} (\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i))$$

where g is a distribution, $\mathcal{S} = \{z_1, \dots, z_m\}$ and $\mathcal{S}^{(i)} = \{z_1, \dots, z_{i-1}, z', z_{i+1}, \dots, z_m\}$ are training datasets with $z', z_1, \dots, z_m \stackrel{\text{iid}}{\sim} g$.

Proof. As $z', z_1, \dots, z_m \stackrel{\text{iid}}{\sim} g$, then for every i

$$E_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S}))) = E_{\substack{\mathcal{S} \sim g \\ z' \sim g}} (\ell(\mathfrak{A}(\mathcal{S}), z')) = E_{\substack{\mathcal{S} \sim g \\ z' \sim g}} (\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i))$$

and

$$E_{\mathcal{S} \sim g} (\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))) = E_{\substack{\mathcal{S} \sim g \\ i \sim U\{1, \dots, m\}}} (\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i))$$

□

Note 20. Hence, we desire to design learning algorithms corresponding to as small as possible (3.2).

4. IMPLEMENTATION IN REGULARIZED LOSS LEARNING PROBLEMS

Definition 21. A Regularized Loss Minimization (RLM) learning problem $(\mathcal{H}, \mathcal{Z}, \ell)$ with a regularization function $J : \mathcal{H} \rightarrow \mathbb{R}$ aims at a RLM learning rule that is the minimizer

$$(4.1) \quad h^* = \arg \min_{h \in \mathcal{H}} \left(\hat{R}_{\mathcal{S}}(h) + J(h) \right)$$

where $\hat{R}_{\mathcal{S}}(\cdot) = \frac{1}{m} \sum_{i=1}^m \ell(\cdot, z_i)$, and $\mathcal{S} = \{z_1, \dots, z_m\} \subset$ an IID training sample.

Remark 22. The motivation for considering the regularization function J in (4.1) is to: (1.) control complexity and (2.) improve stability; as we will see later.

Note 23. Here, we make our example more specific and narrow it to the Ridge RLM learning problem (could have been LASSO etc.).

Definition 24. The Ridge RLM learning problem $(\mathcal{H}, \mathcal{Z}, \ell)$, with $\mathcal{H} = \mathcal{W} \subset \mathbb{R}^d$, uses regularization function $J(w; \lambda) = \lambda \|w\|_2^2$ with $\lambda > 0$, $w \in \mathcal{W}$ and produces learning rule

$$(4.2) \quad \mathfrak{A}(\mathcal{S}) = \arg \min_{w \in \mathcal{W}} \left(\hat{R}_{\mathcal{S}}(w) + \lambda \|w\|_2^2 \right)$$

Note 25. Recall (Term 1) how the regularization function in Ridge RLM learning problem $(\mathcal{H}, \mathcal{Z}, \ell)$ penalizes complexity. Essentially, it implies a sequence of hypothesis $\mathcal{H}_1 \subset \mathcal{H}_2 \subset \dots$ with $\mathcal{H}_i = \{w \in \mathbb{R}^d : \|w\|_2 < i\}$ due to duality of the corresponding minimization problem.

Note 26. Below, we will try to analyze the behavior of Ridge RLM learning rule (4.2) w.r.t. the Risk decomposition (2.3). In particular, to upper bounded w.r.t. the shrinkage term λ , training sample size m , and other characteristics.

4.1. Bounding the empirical risk (I) in (2.3).

Note 27. From (4.2), we have

$$\begin{aligned}\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) &\leq \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) + \lambda \|\mathfrak{A}(\mathcal{S})\|_2^2 \\ &\leq \hat{R}_{\mathcal{S}}(w) + \lambda \|w\|_2^2; \quad \forall w \in \mathcal{W}\end{aligned}$$

and by taking expectations w.r.t. \mathcal{S} , it is

$$(4.3) \quad E_{\mathcal{S} \sim g} \left(\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right) \leq R_g(w) + \lambda \|w\|_2^2; \quad \forall w \in \mathcal{W}$$

because $E_{\mathcal{S} \sim g} \left(\hat{R}_{\mathcal{S}}(\cdot) \right) = \frac{1}{m} \sum_{i=1}^m E_{\mathcal{S} \sim g} (\ell(\cdot, z_i)) = R_g(\cdot)$.

Note 28. From (4.3), we observe that part (I) in the expected risk decomposition (2.3), (aka the upper bound of the expected empirical risk) increases with the regularization term $\lambda > 0$. (!!!) –Although anticipated, it did not start well.

4.2. Bounding the empirical risk (II) in (2.3).

Note 29. As we will see to bound (II) in (2.3) we will have to impose an additional condition on the behavior of loss function apart from convexity; e.g. Lipschitzness.

Assumption 30. *The loss function $\ell(\cdot, z)$ in (4.2) is convex for any $z \in \mathcal{Z}$.*

Note 31. Let $\tilde{R}_{\mathcal{S}}(w) = \hat{R}_{\mathcal{S}}(w) + \lambda \|w\|_2^2$. $\tilde{R}_{\mathcal{S}}(\cdot)$ is 2λ -strongly convex as the sum of a convex function $\hat{R}_{\mathcal{S}}(\cdot)$ (Assumption 30) and a 2λ -strongly convex function $J(\cdot; \lambda) = \lambda \|\cdot\|_2^2$ (results directly from Definition 14 in Lecture notes 2).

Fact 32. *(To be used in Note 33) If f is λ -strongly convex and u is a minimizer of f then for any w*

$$f(w) - f(u) \geq \frac{\lambda}{2} \|w - u\|^2$$

Note 33. Let $\mathfrak{A}(\mathcal{S})$ be the Ridge learning algorithm output minimizing (4.2). Because $\tilde{R}_{\mathcal{S}}(\cdot)$ is 2λ -strongly convex and $\mathfrak{A}(\mathcal{S})$ is its minimizer, according to Fact 32, for $\mathfrak{A}(\mathcal{S})$ and any $w \in \mathcal{W}$, it is

$$(4.4) \quad \tilde{R}_{\mathcal{S}}(w) - \tilde{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \geq \lambda \|w - \mathfrak{A}(\mathcal{S})\|^2, \quad \forall w \in \mathcal{W}$$

Note 34. Also, for any $w, u \in \mathcal{W}$, it is

$$\begin{aligned}\tilde{R}_{\mathcal{S}}(w) - \tilde{R}_{\mathcal{S}}(u) &= \left(\hat{R}_{\mathcal{S}}(w) + \lambda \|w\|_2^2 \right) - \left(\hat{R}_{\mathcal{S}}(u) + \lambda \|u\|_2^2 \right) \\ &= \left(\hat{R}_{\mathcal{S}^{(i)}}(w) + \lambda \|w\|_2^2 \right) - \left(\hat{R}_{\mathcal{S}^{(i)}}(u) + \lambda \|u\|_2^2 \right) \\ &\quad + \frac{\ell(w, z_i) - \ell(u, z_i)}{m} + \frac{\ell(u, z') - \ell(w, z')}{m} \\ &= \tilde{R}_{\mathcal{S}^{(i)}}(w) - \tilde{R}_{\mathcal{S}^{(i)}}(u) + \frac{\ell(w, z_i) - \ell(u, z_i)}{m} + \frac{\ell(u, z') - \ell(w, z')}{m}\end{aligned}$$

Choosing $w = \mathfrak{A}(\mathcal{S}^{(i)})$ and $u = \mathfrak{A}(\mathcal{S})$, and the fact that $\tilde{R}_{\mathcal{S}^{(i)}}(\mathfrak{A}(\mathcal{S}^{(i)})) \leq \tilde{R}_{\mathcal{S}^{(i)}}(\mathfrak{A}(\mathcal{S}))$, it is

$$(4.5) \quad \tilde{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}^{(i)})) - \tilde{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \leq \frac{\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i)}{m} + \frac{\ell(\mathfrak{A}(\mathcal{S}), z') - \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z')}{m}$$

Note 35. Then (4.4) and (4.5) imply

$$(4.6) \quad \lambda \|\mathfrak{A}(\mathcal{S}^{(i)}) - \mathfrak{A}(\mathcal{S})\|^2 \leq \frac{\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i)}{m} + \frac{\ell(\mathfrak{A}(\mathcal{S}), z') - \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z')}{m}$$

Note 36. Now that we brought it in form (4.6), we can impose/use an additional assumption on the loss to bound it. In this example we use Lipschitzness.

Assumption 37. The loss function $\ell(\cdot, z)$ in (4.2) is convex and ρ -Lipschitz for any $z \in \mathcal{Z}$.

Note 38. Given ρ -Lipschitzness in Assumption 37, it is

$$(4.7) \quad \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \leq \rho \|\mathfrak{A}(\mathcal{S}^{(i)}) - \mathfrak{A}(\mathcal{S})\|$$

$$(4.8) \quad \ell(\mathfrak{A}(\mathcal{S}), z') - \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z') \leq \rho \|\mathfrak{A}(\mathcal{S}^{(i)}) - \mathfrak{A}(\mathcal{S})\|$$

and hence plugging 4.7 and (4.8) in (4.6) yields

$$(4.9) \quad \|\mathfrak{A}(\mathcal{S}^{(i)}) - \mathfrak{A}(\mathcal{S})\| \leq 2 \frac{\rho}{\lambda m}$$

Note 39. Plugging (4.9) in (4.7) yields

$$\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \leq 2 \frac{\rho^2}{\lambda m}$$

Note 40. Using Theorem 19, we get an upper bound for the stability / over-fitting

$$\mathbb{E}_{\mathcal{S} \sim g} \left(R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{D}}(\mathfrak{A}(\mathcal{S})) \right) = \mathbb{E}_{\substack{\mathcal{S} \sim g, z' \sim g \\ i \sim \text{U}\{1, \dots, m\}}} (\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i)) \leq 2 \frac{\rho^2}{\lambda m}$$

Note 41. After this saga, the researcher could come to the conclusion that: A Ridge RLM learning problem $(\mathcal{H}, \mathcal{Z}, \ell)$ with the loss function which is convex and ρ -Lipschitz, regularizer $J(\cdot; \lambda) = \lambda \|\cdot\|^2$, $\lambda > 0$, and learning rule trained against an iid sample $\mathcal{S} = \{z_i\}_{i=1}^m$ from g is on-average-replace-one-stable with rate $\epsilon(m) = 2\frac{\rho^2}{\lambda m}$; i.e.

$$(4.10) \quad \mathbb{E}_{\mathcal{S} \sim g} \left(R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right) \leq 2 \frac{\rho^2}{\lambda m}$$

Note 42. From (4.10), we see that stability improves (and over-fitting decreases) as the shrinkage parameter λ increases.

4.3. Bounding the Risk (2.3).

Note 43. Given the bounds (4.3) and (4.10), the decomposition of the expected Risk in (2.3) yields that: A Ridge RLM learning problem $(\mathcal{H}, \mathcal{Z}, \ell)$ with the loss function which is convex and ρ -Lipschitz, regularizer $J(\cdot; \lambda) = \lambda \|\cdot\|^2$, $\lambda > 0$, and learning rule trained against an iid sample $\mathcal{S} = \{z_i\}_{i=1}^m$ from g has Expected Risk bound

$$(4.11) \quad \mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S}))) \leq \underbrace{R_g(w) + \lambda \|w\|_2^2}_{(I)} + \underbrace{2 \frac{\rho^2}{\lambda m}}_{(II)}; \quad \forall w \in \mathcal{W}$$

Note 44. From (4.11), we see that there is a trade-off between Empirical Risk (I) and stability/overfitting (II) with regards the regularization parameter λ . It is desirable to use the optimal $\lambda > 0$ corresponding to the smallest bound in (4.11); it has to both fit the training data well (but perhaps not too well) and be very stable to different training data from the same g (but perhaps not too stable)!

Assumption 45. Assume that the learning problem is additionally B -bounded by $B > 0$; i.e. $\mathcal{H} = \{w \in \mathbb{R}^d : \|w\|_2 \leq B\}$.

Note 46. If additionally the learning problem is B -bounded then the upper bound in (4.11) becomes

$$(4.12) \quad \mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S}))) \leq \underbrace{\min_{w \in \mathcal{H}} R_g(w) + \lambda B^2}_{(I)} + \underbrace{2 \frac{\rho^2}{\lambda m}}_{(II)}$$

Note 47. Furthermore, if we choose a shrinkage parameter $\lambda_m = \sqrt{\frac{2}{m} \frac{\rho}{B}}$ in (4.2) the upper bound in (4.11) becomes

$$(4.13) \quad \mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S}))) \leq \min_{w \in \mathcal{H}} R_g(w) + \rho B \sqrt{\frac{8}{m}}$$

4.4. Deriving PAC guarantees.

Note 48. In particular, if the size m of the training sample \mathcal{S} is

$$m \geq \frac{\lambda_m^2 B^2 8}{\epsilon^2} = \frac{8\rho^2 B^2}{\epsilon^2}$$

for some $\epsilon > 0$ then for every distribution g , the upper bound in (4.11) becomes

$$(4.14) \quad \mathbb{E}_{\mathcal{S} \sim g}(R_g(\mathfrak{A}(\mathcal{S}))) \leq \min_{w \in \mathcal{H}} R_g(w) + \epsilon$$

and hence we can have a PAC guarantee that is summarized in the following.

Summary 49. Let $(\mathcal{H}, \mathcal{Z}, \ell)$ be a learning problem convex-Lipschitz-bounded with parameters ρ and B , and let m be the training data size. The associated Ridge Regularized Loss Minimization rule with regularization function $\lambda_m = \sqrt{\frac{2}{m} \frac{\rho}{B}}$ satisfies

$$\mathbb{E}_{\mathcal{S} \sim g}(R_g(\mathfrak{A}(\mathcal{S}))) \leq \min_{w \in \mathcal{H}} R_g(w) + \lambda_m B \sqrt{\frac{8}{m}}$$

Moreover, if $m \geq \frac{8\rho^2 B^2}{\epsilon^2}$, $\epsilon > 0$ then for every distribution g we can get a PAC -guarantee

$$\mathbb{E}_{\mathcal{S} \sim g}(R_g(\mathfrak{A}(\mathcal{S}))) \leq \min_{w \in \mathcal{H}} R_g(w) + \epsilon$$

Lecture notes 4: Gradient descent

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim. To introduce gradient descent, its motivation, description, practical tricks, analysis in the convex scenario, and implementation.

Reading list & references:

- (1) Shalev-Shwartz, S., & Ben-David, S. (2014). Understanding machine learning: From theory to algorithms. Cambridge university press.
 - Ch. 14.1 Gradient Descent

1. MOTIVATIONS

Problem 1. Consider a learning problem $(\mathcal{H}, \mathcal{Z}, \ell)$. Learning may involve the computation of the minimizer $h^* \in \mathcal{H}$, where \mathcal{H} is a class of hypotheses, of the empirical risk function (ERF) $\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n \ell(h, z_i)$ given a finite sample $\{z_i; i = 1, \dots, n\}$ generated from the data generating model $g(\cdot)$ and using loss $\ell(\cdot)$; that is

$$(1.1) \quad h^* = \arg \min_h \left(\hat{R}(h) \right) = \arg \min_h \left(\frac{1}{n} \sum_{i=1}^n \ell(h, z_i) \right)$$

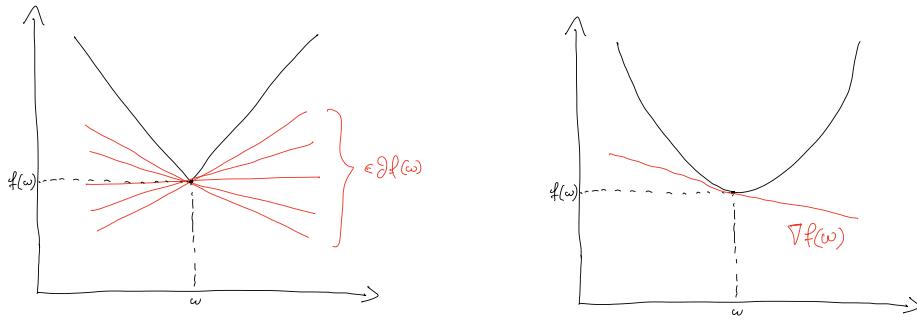
or its corresponding regularized version. If analytical minimization of (1.1) is impossible or impractical, numerical procedures can be applied; eg Gradient Descent (GD) algorithms. Such approaches introduce numerical errors in the solution.

2. SUBGRADIENT

Definition 2. Vector v is called subgradient of a function $f : S \rightarrow \mathbb{R}$ at $w \in S$ if

$$(2.1) \quad \forall u \in S, \quad f(u) \geq f(w) + \langle u - w, v \rangle$$

Note 3. There may be more than one subgradients of a function at a specific point. As seen by (2.1), subgradients are the slopes of all the lines passing through the point $(w, f(w))$ and been under the function $f(\cdot)$.



(A) subgradients satisfying (2.1)
in the non-differentiable case

(B) gradient satisfying the
equality in (2.1) in the differen-
tiable case

Definition 4. The set of subgradients of function $f : S \rightarrow \mathbb{R}$ at $w \in S$ is denoted by $\partial f(w)$.

Fact 5. Properties of subgradient sets useful for the subgradient construction

- (1) If function $f : S \rightarrow \mathbb{R}$ is differentiable at w then the only subgradient of f at w is the gradient $\nabla f(w)$, and (2.1) is equality; i.e. $\partial f(w) = \{\nabla f(w)\}$.
- (2) for constants α, β and convex function $f(\cdot)$, it is

$$\partial(\alpha f(w) + \beta) = \alpha(\partial f(w)) = \{\alpha v : v \in \partial f(w)\}$$

- (3) for convex functions $f(\cdot)$ and $g(\cdot)$, it is

$$\partial(f(w) + g(w)) = \partial f(w) + \partial g(w) = \{v + u : v \in \partial f(w), \text{ and } u \in \partial g(w)\}$$

Example 6. Consider the function $f : \mathbb{R} \rightarrow \mathbb{R}_+$ with $f(w) = |w| = \begin{cases} w & w \geq 0 \\ -w & w < 0 \end{cases}$. Find the set of subgradients $\partial f(w)$ for each $w \in \mathbb{R}$.

Solution: Using Fact 5, it is $\partial f(w) = 1$ for $w > 0$ and $\partial f(w) = -1$ for $w < 0$ as f is differentiable for $x \neq 0$. At $x = 0$, f is not differentiable; hence from condition (2.1) it is

$$\forall u \in \mathbb{R}, \quad |u| \geq |0| + (u - 0)v$$

which is satisfied for $v \in [-1, 1]$. Hence,

$$\partial f(w) = \begin{cases} \{-1\} & , w < 0 \\ [-1, 1] & , w = 0 \\ \{1\} & , w > 0 \end{cases}$$

3. DESCRIPTION OF GRADIENT DESCENT

Problem 7. For the sake of notation simplicity and generalization, we will present Gradient Descent (GD) in the following minimization problem

$$(3.1) \quad w^* = \arg \min_{w \in \mathcal{H}} (f(w))$$

where $f : \mathcal{H} \rightarrow \mathbb{R}$, and $w \in \mathcal{H} \subseteq \mathbb{R}^d$; $f(\cdot)$ is the function to be minimized, e.g., $f(\cdot)$ can be an empirical risk function $\hat{R}(\cdot)$.

Algorithm 8. Gradient Descent (GD) algorithm with learning rate $\eta_t > 0$ for the solution of the minimization problem (3.1)

For $t = 1, 2, 3, \dots$ iterate:

(1) compute

$$(3.2) \quad w^{(t+1)} = w^{(t)} - \eta_t v_t; \text{ where } v_t \in \partial f(w^{(t)})$$

(2) terminate if a termination criterion is satisfied, e.g.

If $t \geq T_{\max}$ then STOP

Note 9. (Intuition) Assume f is differentiable. GD produces a chain $\{w^{(t)}\}$ that drifts towards a minimum w^* . The chain is directed towards the opposite direction of that of the gradient $\nabla f(\cdot)$ and at a rate controlled by the learning rate η_t .

Note 10. (More intuition) Assume f is differentiable. Consider the (1st order) Taylor polynomial for the approximation of $f(w)$ in a small area around u (i.e. $\|v - u\| = \text{small}$)

$$f(u) \approx P(u) = f(w) + \langle u - w, \nabla f(w) \rangle$$

Assuming convexity for f , it is

$$(3.3) \quad f(u) \geq \underbrace{f(w) + \langle u - w, \nabla f(w) \rangle}_{=P(u;w)}$$

meaning that P lower bounds f . Hence we could design an updating mechanism producing $w^{(t+1)}$ which is nearby $w^{(t)}$ (small steps) and which minimize the linear approximation $P(w)$ of $f(w)$ at $w^{(t)}$

$$(3.4) \quad P(w; w^{(t)}) = f(w^{(t)}) + \langle w - w^{(t)}, \nabla f(w^{(t)}) \rangle.$$

while hoping that this mechanism would push the produced chain $\{w^{(t)}\}$ towards the minimum because of (3.3). Hence we could recursively minimize the linear approximation (3.4) and the

distance between the current state $w^{(t)}$ and the next w value to produce $w^{(t+1)}$; namely

$$\begin{aligned}
 (3.5) \quad w^{(t+1)} &= \arg \min_{\forall w} \left(\frac{1}{2} \|w - w^{(t)}\|^2 + \eta P(w; w^{(t)}) \right) \\
 &= \arg \min_{\forall w} \left(\frac{1}{2} \|w - w^{(t)}\|^2 + \eta \left(f(w^{(t)}) + \langle w - w^{(t)}, \nabla f(w^{(t+1)}) \rangle \right) \right) \\
 &= w^{(t)} - \eta \nabla f(w^{(t)})
 \end{aligned}$$

where parameter $\eta > 0$ controls the trade off in (3.5).

Note 11. GD output (e.g. in Note 14) converges to a local minimum, $w_{\text{GD}}^{(T)} \rightarrow w_*$ (in some sense), under different sets of regularity conditions (some are weaker other stronger). Section 4 has a brief analysis.

Note 12. The parameter η_t is called learning rate (or step size, gain). It determines the size of the steps GD takes to reach a (local) minimum. $\{\eta_t\}$ is a non-negative sequence and it is chosen by the practitioner. In principle, regularity conditions (Note 11) often imply restrictions on the decay of $\{\eta_t\}$ which guide the practitioner to parametrize it properly. Some popular choices of learning rate η_t are:

- (1) constant; $\eta_t = \eta$, for where $\eta > 0$ is a small value. The rationale is that GD chain $\{w_t\}$ performs constant small steps towards the (local) minimum w_* and then oscillate around it.
 - (2) decreasing and converging to zero; $\eta_t \downarrow$ with $\lim_{t \rightarrow \infty} \eta_t = 0$. E.g. $\eta_t = (\frac{C}{t})^\varsigma$ where $\varsigma \in [0.5, 1]$ and $C > 0$. The rationale is that GD algorithm starts by performing larger steps (controlled by C) at the beginning to explore the area for discovering possible minima. Also it reduces the size of those steps with the iterations (controloed by ς) such that eventually when the chain $\{w_t\}$ is close to a possible minimum w_* value to converge and do not overshoot.
 - (3) decreasing and converging to a tiny value τ_* ; $\eta_t \downarrow$ with $\lim_{t \rightarrow \infty} \eta_t = \tau_*$. E.g. $\eta_t = (\frac{C}{t})^\varsigma + \tau_*$ with $\varsigma \in (0.5, 1]$, $C > 0$, and $\tau_* \approx 0$. Same as previously, but the algorithm aims at oscillating around the detected local minimum.
 - (4) constant until an iteration T_0 and then decreasing; Eg $\eta_t = \left(\frac{C}{\max(t, T_0)}\right)^\varsigma$ with $\varsigma \in [0.5, 1]$ and $C > 0$, and $T_0 < T$. The rationale is that at the first stage of the iterations (when $t \leq T_0$) the algorithm may need a constant large steps for a significant number of iterations T_0 in order to explore the domain; and hence in order for the chain $\{w_t\}$ to reach the area around the (local) minimum w_* . In the second stage, hoping that the chain $\{w_t\}$ may be in close proximity to the (local) minimum w_* the algorithm progressively performs smaller steps to converge towards the minimum w_* . The first stage ($t \leq T_0$) is called burn-in; the values $\{w_t\}$ produced during the burn-in ($t \leq T_0$) are often discarded/ignored from the output of the GD algorithm.
- Parameters $C, \varsigma, \tau_*, T_0$ may be chosen based on pilot runs against a small fraction of the training data set.

Note 13. There are several practical termination criteria that can be used in GD Algorithm 8(step 2). They aim to terminate the recursion in practice. Some popular termination criteria are

- (1) terminate when the gradient is sufficiently close to zero; i.e. if $\|\nabla f(w^{(t)})\| \leq \epsilon$, for some pre-specified tiny $\epsilon > 0$ then STOP
- (2) terminate when the chain $w^{(t)}$ does not change; i.e. if $\|w^{(t+1)} - w^{(t)}\| \leq \epsilon \|w^{(t)}\|$ for some pre-specified tiny $\epsilon > 0$ then STOP
- (3) terminate when a pre-specified number of iterations T is performed; i.e. if $t \geq T$ then STOP

Here (1) may be deceive if the chain is in a flat area, (2) may be deceived if the learning rate become too small, (3) is obviously a last resort.

Note 14. Given T iterations of GD algorithm, the output of GD can be (but not a exclusively),

- (1) the average (after discarding the first few iterations of $w^{(t)}$ for stability reasons)

$$(3.6) \quad w_{\text{GD}}^{(T)} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$$

- (2) or the best value discovered

$$w_{\text{GD}}^{(T)} = \arg \min_{\forall w_t} \left(f(w^{(t)}) \right)$$

- (3) or the last value discovered

$$w_{\text{GD}}^{(T)} = w^{(T)}$$

4. ANALYSIS OF GRADIENT DESCENT (ALGORITHM 8)

Note 15. We analyze the behavior of GD addressing the minimization Problem 7 under Assumptions 16.

Assumption 16. *For the sake of the analysis of the GD, let us consider:*

- (1) *The function $f(\cdot)$ is convex and Lipschitz*
- (2) *GD has a constant learning rate $\eta_t = \eta$,*
- (3) *GD output $w_{\text{GD}}^{(T)} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$*

Lemma 17. *Let $\{v_t; t = 1, \dots, T\}$ be a sequence of vectors. Any algorithm with $w^{(1)} = 0$ and $w^{(t+1)} = w^{(t)} - \eta v_t$ for $t = 1, \dots, T$ satisfies*

$$(4.1) \quad \sum_{t=1}^T \langle w^{(t)} - w^*, v_t \rangle \leq \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2$$

Proof. Omitted; see the Exercise 12 in the Exercise sheet 1. □

No need to memorize

Note 18. To find an upper bound of the GD error, we try to bound the error $f(w_{\text{GD}}^{(T)}) - f(w^*)$ with purpose to use Lemma 17.

Note 19. Consider the minimization problem (3.1). Given Assumptions 16, the error can be bounded as

$$(4.2) \quad f(w_{\text{GD}}^{(T)}) - f(w^*) \leq \frac{\|w^*\|^2}{2\eta T} + \frac{\eta}{2} \frac{1}{T} \sum_{t=1}^T \|v_t\|^2$$

where $v_t \in \partial f(w^{(t)})$. If $f(\cdot)$ is differentiable then $v_t = \nabla f(w^{(t)})$

Proof. It is¹

$$\begin{aligned}
f(w_{\text{GD}}^{(T)}) - f(w^*) &= f\left(\frac{1}{T} \sum_{t=1}^T w_t\right) - f(w^*) \\
(4.3) \quad &\leq \frac{1}{T} \sum_{t=1}^T (f(w_t) - f(w^*)) \quad (\text{by Jensen's inequality}) \\
&\leq \frac{1}{T} \sum_{t=1}^T \langle w^{(t)} - w^*, v_t \rangle \quad (\text{by convexity of } f(\cdot)) \\
(4.4) \quad &\leq \frac{\|w^*\|^2}{2\eta T} + \frac{\eta}{2} \frac{1}{T} \sum_{t=1}^T \|v_t\|^2 \quad (\text{by proceeding Lemma})
\end{aligned}$$

□

Note 20. Note 19 shows that it is important to bound the (sub-)gradient in (4.2) in a meaningful manner. Lemma 21 shows that if we assume Lipschitzness as well, the (sub-)gradient has the so-called self-bounded behavior.

Lemma 21. ² $f : S \rightarrow \mathbb{R}$ is ρ -Lipschitz over an open convex set S if and only if for all $w \in S$ and $v \in \partial f(w)$ it is $\|v\| \leq \rho$.

Solution: \implies Let $f : S \rightarrow \mathbb{R}$ be ρ -Lipschitz over convex set S , $w \in S$ and $v \in \partial f(w)$.

- Since S is open we get that there exist $\epsilon > 0$ such as $u := w + \epsilon \frac{v}{\|v\|}$ where $u \in S$. So $\langle u - w, v \rangle = \epsilon \|v\|$ and $\|u - w\| = \epsilon$.
 - From the subgradient definition we get

$$f(u) - f(w) \geq \langle u - w, v \rangle = \epsilon \|v\|$$

– From the Lipschitzness of $f(\cdot)$ we get

$$f(u) - f(w) \leq \rho \|u - w\| = \rho \epsilon$$

Therefore $\|v\| \leq \rho$.

Proof. For \Leftarrow see Exercise 4 in the Exercise sheet. □

Note 22. The following summarizes Note 19 and Lemma 21 with respect to the GD algorithm under Assumption 16.

Note 23. Let $f(\cdot)$ be a convex and ρ -Lipschitz function. If we run GD algorithm of f with learning rate $\eta > 0$ for T steps the output $w_{\text{GD}}^{(T)} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$ satisfies

$$f(w_{\text{GD}}^{(T)}) - f(w^*) \leq \frac{\|w^*\|^2}{2\eta T} + \frac{\eta}{2} \frac{1}{T} \sum_{t=1}^T \|v_t\|^2$$

where $\|v_t\| \leq \rho$, $v_t \in \partial f(w^{(t)})$. If $f(\cdot)$ is differentiable then $v_t = \nabla f(w^{(t)})$.

¹Jensen's inequality for convex $f(\cdot)$ is $f(\mathbf{E}(x)) \leq \mathbf{E}(f(x))$

²If this was a Homework there would be a Hint:

• If S is open there exist $\epsilon > 0$ such as $u = w + \epsilon \frac{v}{\|v\|}$ such as $u \in S$

Proof. Straightforward from Lemma 17 and Note 19. \square

Note 24. The following shows that given a learning rate depending on the iteration t , we can reduce the upper bound of the error as well as find the number of required iterations to achieve convergence.

Note 25. (Cont Prop. 23) Let $f(\cdot)$ be a convex and ρ -Lipschitz function, and let $\mathcal{H} = \{w \in \mathbb{R} : \|w\| \leq B\}$. Assume we run GD algorithm of $f(\cdot)$ with learning rate $\eta_t = \sqrt{\frac{B^2}{\rho^2 T}}$ for T steps, and output $w_{\text{GD}}^{(T)} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$. Then

(1) upper bound on the sub-optimality is

$$(4.5) \quad f\left(w_{\text{GD}}^{(T)}\right) - f(w^*) \leq \frac{B\rho}{\sqrt{T}}$$

(2) a given level off accuracy ε such that $f\left(w_{\text{GD}}^{(T)}\right) - f(w^*) \leq \varepsilon$ can be achieved after T iterations

$$T \geq \frac{B^2 \rho^2}{\varepsilon^2}.$$

Proof. Part 1 is a simple substitution from Proposition 23, and part 2 is implied from part 1. \square

The result on Note 25 heavily relies on setting suitable values for B and ρ which is rather a difficult task to be done in very complicated learning problems (e.g., learning a neural network).

5. EXAMPLES ³

Example 26. Consider the simple Normal linear regression problem where the dataset $\{z_i = (y_i, x_i)\}_{i=1}^n \in \mathcal{S}$ is generated from a Normal data generating model

$$(5.1) \quad \begin{pmatrix} y_i \\ x_i \end{pmatrix} \stackrel{\text{iid}}{\sim} N \left(\begin{pmatrix} \mu_y \\ \mu_x \end{pmatrix}, \begin{bmatrix} \sigma_y^2 & \rho\sqrt{\sigma_y^2 \sigma_x^2} \\ \rho\sqrt{\sigma_y^2 \sigma_x^2} & \sigma_x^2 \end{bmatrix} \right)$$

for $i = 1, \dots, n$. Consider a hypothesis space \mathcal{H} of linear functions $h : \mathbb{R}^2 \rightarrow \mathbb{R}$ with $h(w) = w_1 + w_2 x$. The exact solution (which we pretend we do not know) is given as

$$(5.2) \quad \begin{pmatrix} w_1^* \\ w_2^* \end{pmatrix} = \begin{pmatrix} \mu_y - \rho \frac{\sigma_y}{\sigma_x} \mu_x \\ \rho \frac{\sigma_y}{\sigma_x} \end{pmatrix}.$$

To learn the optimal $w^* = (w_1^*, w_2^*)^\top$, we consider a loss $\ell(w, z_i = (x_i, y_i)^\top) = (y_i - [w_1 + w_2 x_i])^2$, which leads to the minimization problem

$$w^* = \arg \min_w \hat{R}_{\mathcal{S}}(w) = \arg \min_w \left(\frac{1}{n} \sum_{i=1}^n (y_i - w_1 - w_2 x_i)^2 \right)$$

The GD Algorithm 8 with learning rate η is

For $t = 1, 2, 3, \dots$ iterate:

³Code is available in https://github.com/georgios-stats/Machine_Learning_and_Neural_Networks_III_Epiphanys_2025/tree/main/Lecture_notes/code/example_GD.R

(1) compute

$$(5.3) \quad w^{(t+1)} = w^{(t)} - \eta v_t, \quad \text{where } v_t = \begin{pmatrix} 2w_1^{(t)} + 2w_2^{(t)}\bar{x} - 2\bar{y} \\ 2w_1^{(t)}\bar{x} + 2w_2^{(t)}\bar{x^2} - 2y^\top x \end{pmatrix}$$

(2) terminate if a termination criterion is satisfied, e.g.

If $t \geq T$ then STOP

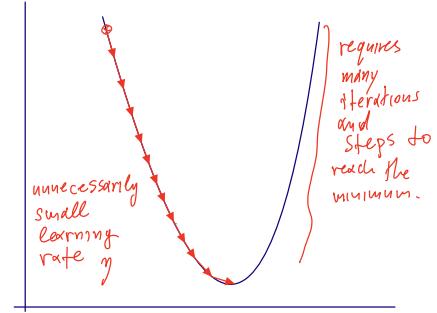
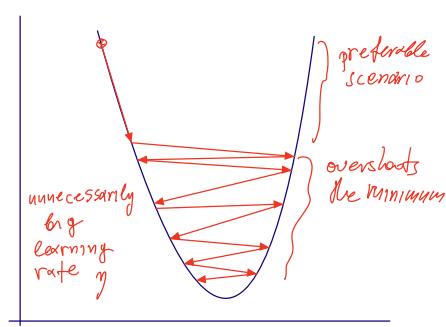
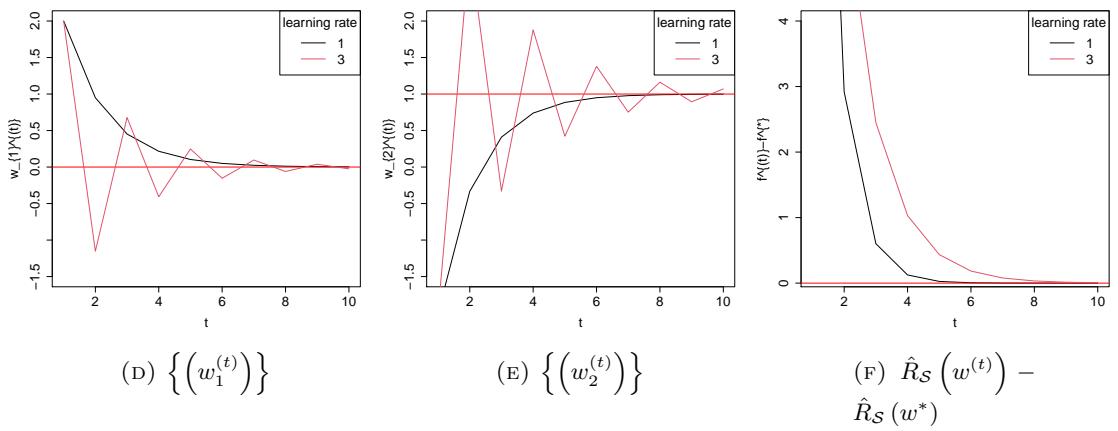
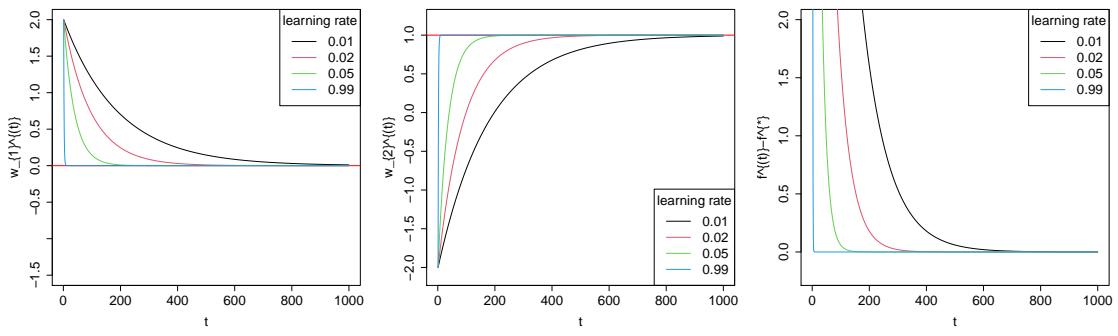
This is because $\hat{R}_{\mathcal{D}}(w)$ is differentiable in \mathbb{R}^2 so $\partial \hat{R}_{\mathcal{D}}(w) = \{\nabla \hat{R}_{\mathcal{D}}(w)\}$ and because

$$\nabla \hat{R}_{\mathcal{D}}(w) = \begin{pmatrix} \frac{d}{dw_1} \hat{R}_{\mathcal{S}}(w) \\ \frac{d}{dw_2} \hat{R}_{\mathcal{S}}(w) \end{pmatrix} = \dots = \begin{pmatrix} 2w_1^{(t)} + 2w_2^{(t)}\bar{x} - 2\bar{y} \\ 2w_1^{(t)}\bar{x} + 2w_2^{(t)}\bar{x^2} - 2y^\top x \end{pmatrix}$$

Consider data size $n = 100$, and parameters $\rho = 0.2$, $\sigma_y^2 = 1$ and $\sigma_x^2 = 1$. Then the real value (5.2) that I need to learn equals to $w^* = (0, 1)^\top$. Consider a GD seed $w_0 = (2, -2)$, and total number of iterations $T = 1000$.

Figures 5.1a, 5.1b, and 5.1c present trace plots of the chain $\{(w^{(t)})\}$ and error $\hat{R}_{\mathcal{S}}(w^{(t)}) - \hat{R}_{\mathcal{S}}(w^*)$ produced by running GD for $T = 1000$ total iterations and for different (each time) constant learning rates $\eta \in \{0.01, 0.02, 0.05, 0.99\}$. We observe that the larger learning rates under consideration were able to converge faster to the minimum w^* . This is because they perform larger steps and can learn faster -this is not a panacea.

Figures 5.1d, 5.1e, and 5.1f present trace plots of the chain $\{(w^{(t)})\}$, and of the error $\hat{R}_{\mathcal{S}}(w^{(t)}) - \hat{R}_{\mathcal{S}}(w^*)$ produced by running GD for $T = 1000$ total iterations and for learning rate $\eta = 1.0$ (previously considered) and a very big learning rate $\eta = 3.0$. We observe that the very big learning rate $\eta = 3.0$ presents slower convergence to the minimum w^* . This is because it creates unreasonably big steps in (3.2) that the produced chain overshoots the global minimum; see the cartoon in Figures 5.1g and 5.1h.



(G) Unnecessarily large learning rate

(H) Unnecessarily small learning rate

FIGURE 5.1

Lecture notes 5: Stochastic gradient descent

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim. To introduce the stochastic gradient descent (motivation, description, practical tricks, analysis in the convex scenario, and implementation).

Reading list & references:

- (1) Shalev-Shwartz, S., & Ben-David, S. (2014). Understanding machine learning: From theory to algorithms. Cambridge university press.
 - Ch. 14.3 Stochastic Gradient Descent (SGD), 14.5 Variants, 14.5 Learning with SGD
- (2) Bottou, L. (2012). Stochastic gradient descent tricks. In Neural networks: Tricks of the trade (pp. 421-436). Springer, Berlin, Heidelberg.

1. MOTIVATIONS FOR STOCHASTIC GRADIENT DESCENT

Problem 1. Consider a learning problem $(\mathcal{H}, \mathcal{Z}, \ell)$. Learning may involve the computation of the minimizer $w^* \in \mathcal{W}$ of the risk function (RF) $R_g(w) = \mathbb{E}_{z \sim g}(\ell(h_w, z))$ given an unknown data generating model $g(\cdot)$, using a known tractable loss $\ell(\cdot, \cdot)$, and hypothesis $h_w \in \mathcal{H}$; that is

$$(1.1) \quad w^* = \arg \min_w (R_g(w)) = \arg \min_w (\mathbb{E}_{z \sim g}(\ell(h_w, z)))$$

Note 2. Gradient descent (GD) cannot be directly utilized to address Problem 1 (i.e., minimize the Risk function) because g is unknown, and because (1.1) may involve a computationally intractable integration. Instead GD aims to minimize the ERF $\hat{R}(w) = \frac{1}{n} \sum_{i=1}^n \ell(h_w, z_i)$ which ideally can be considered as a proxy when data size n is big (big-data).

Note 3. The implementation of GD may be computationally impractical even in problems aiming to minimize ERF $\hat{R}_n(w)$ if we have big data ($n \approx \text{big}$). This is because GD requires the recursive computation of the exact gradient $\partial_w \hat{R}_n(w) = \left\{ \frac{1}{n} \sum_{i=1}^n v_i : v_i \in \partial_w \ell(h_w, z_i) \right\}$ using (required to scan through) all the data $\{z_i\}$ at each iteration. That may be too slow.

Note 4. Stochastic gradient descent (SGD) aims at solving (1.1), and overcoming the issues in Notes 2 & 3 by using an unbiased estimator of the actual gradient (or some sub-gradient) based on a sample (a single example or a set of examples) properly drawn from g .

2. STOCHASTIC GRADIENT DESCENT

Problem 5. For the sake of notation simplicity and generalization, we present Stochastic Gradient Descent (SGD) in the following minimization problem

$$(2.1) \quad w^* = \arg \min_w (f(w))$$

where here $f : \mathbb{R}^d \rightarrow \mathbb{R}$, and $w \in \mathcal{W} \subseteq \mathbb{R}^d$; $f(\cdot)$ is the unknown function to be minimized, e.g., $f(\cdot)$ can be the risk function $R_g(w) = \mathbb{E}_{z \sim g}(\ell(h_w, z))$.

Algorithm 6. *Stochastic Gradient Descent (SGD) with learning rate $\eta_t > 0$ for Problem 5*

For $t = 1, 2, 3, \dots$ iterate:

(1) compute

$$(2.2) \quad w^{(t+1)} = w^{(t)} - \eta_t v_t,$$

where v_t is a random vector such that $E(v_t | w^{(t)}) \in \partial f(w^{(t)})$

(2) terminate if a termination criterion is satisfied, e.g.

If $t \geq T_{\max}$ then STOP

Note 7. If f is differentiable at $w^{(t)}$, it is $\partial f(w^{(t)}) = \{\nabla f(w^{(t)})\}$. Hence v_t is such as $E(v_t | w^{(t)}) = \nabla f(w^{(t)})$ in Algorithm 6 step 1.

Note 8. Assume f is differentiable (for simplicity). To compare SGD with GD, we can re-write (2.2) in the SGD Algorithm 6 as

$$(2.3) \quad w^{(t+1)} = w^{(t)} - \eta_t [\nabla f(w^{(t)}) + \xi_t],$$

where

$$\xi_t := v_t - \nabla f(w^{(t)})$$

represents the (observed) noise introduced in (2.2) by using a random realization of the exact gradient.

Note 9. Given T SGD algorithm iterations, the output of SGD can be (but not exclusively)

(1) the average (after discarding the first few iterations of $w^{(t)}$ for stability reasons)

$$(2.4) \quad w_{\text{SGD}}^{(T)} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$$

(2) or the best value discovered

$$w_{\text{SGD}}^{(T)} = \arg \min_{\{w_t\}} (f(w^{(t)}))$$

(3) or the last value discovered

$$w_{\text{SGD}}^{(T)} = w^{(T)}$$

Note 10. SGD output converges to a local minimum, $w_{\text{SGD}}^{(T)} \rightarrow w_*$ (in some sense), under different sets of regularity conditions –Section 3 has a brief analysis. To achieve this, Conditions 11 on the learning rate are rather inevitable and should be satisfied.

Condition 11. Regarding the learning rate (or gain) $\{\eta_t\}$ should satisfy conditions

(1) $\eta_t \geq 0$,

- (2) $\sum_{t=1}^{\infty} \eta_t = \infty$
- (3) $\sum_{t=1}^{\infty} \eta_t^2 < \infty$

Note 12. The popular learning rates $\{\eta_t\}$ in Note 14 in Lect. notes 4 “Gradient descent” satisfy Condition 11 and hence can be used in SGD too. Once parameterized, η_t can be tuned based on pilot runs using a reasonably small fraction of the training data set.

Note 13. (Intuition on Condition 11). Assume that v_t is bounded. Condition 11(3) aims at reducing the effect of the randomness in v_t (introduced noise ξ_t) because it implies $\eta_t \searrow 0$ as $t \rightarrow \infty$; if this was not the case then

$$w^{(t+1)} - w^{(t)} = -\eta_t v_t \rightarrow 0$$

may not be satisfied and the chain $\{w^{(t)}\}$ may not converge. Condition 11(2) prevents η_t from reducing too fast and allows the generated chain $\{w^{(t)}\}$ to be able to converge. E.g., after t iterations

$$\begin{aligned} \|w^{(t)} - w^*\| &= \|w^{(t)} \pm w^{(0)} - w^*\| \geq \|w^{(0)} - w^*\| - \|w^{(t)} - w^{(0)}\| \\ &\geq \|w^{(0)} - w^*\| - \sum_{t=0}^{\infty} \|w^{(t+1)} - w^{(t)}\| = \|w^{(0)} - w^*\| - \sum_{t=0}^{T-1} \|\eta_t v_t\| \end{aligned}$$

However if it was $\sum_{t=1}^{\infty} \eta_t < \infty$ it would be $\sum_{t=0}^{\infty} \|\eta_t v_t\| < \infty$ and hence $w^{(t)}$ would never converge to w^* if the seed $w^{(0)}$ is far enough from w^* .

3. ANALYSIS OF SGD (ALGORITHM 6)

Note 14. Recall (Note 8) that the stochasticity of SGD comes from the stochastic sub-gradients $\{v_t\}$ in (2.2); hence the expectations below are under these random vectors’ distributions.

Proposition 15. *Let $f(\cdot)$ be a convex function. If we run SGD algorithm of f with learning rate $\eta_t > 0$ for T steps, the output $w_{\text{SGD}}^{(T)} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$ satisfies*

$$(3.1) \quad \mathbb{E} \left(f \left(w_{\text{SGD}}^{(T)} \right) \right) - f(w^*) \leq \frac{\|w^*\|^2}{2\eta T} + \frac{\eta}{2} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \|v_t\|^2$$

Proof. Let $v_{1:t} = (v_1, \dots, v_t)$. By Jensens’ inequality (or see 4.3 in Lect. notes 4)

$$(3.2) \quad \mathbb{E} \left(f \left(w_{\text{SGD}}^{(T)} \right) - f(w^*) \right) \leq \mathbb{E} \left(\frac{1}{T} \sum_{t=1}^T \left(f \left(w^{(t)} \right) - f(w^*) \right) \right) = \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left(f \left(w^{(t)} \right) - f(w^*) \right)$$

I will try to use Lemma 17 from Lect. notes 4, hence I need to show

$$(3.3) \quad \mathbb{E} \left(f \left(w^{(t)} \right) - f(w^*) \right) \leq \mathbb{E} \left(\langle w^{(t)} - w^*, v_t \rangle \right)$$

where the expectation is under $v_{1:T}$. It is

$$\begin{aligned} \mathbb{E}_{v_{1:T}} \left(\langle w^{(t)} - w^*, v_t \rangle \right) &= \mathbb{E}_{v_{1:t}} \left(\langle w^{(t)} - w^*, v_t \rangle \right) \\ &= \mathbb{E}_{v_{1:t-1}} \left(\mathbb{E}_{v_{1:t}} \left(\langle w^{(t)} - w^*, v_t \rangle | v_{1:t-1} \right) \right) \quad (\text{law of total expectation}) \end{aligned}$$

But $w^{(t)}$ is fully determined by $v_{1:t-1}$, (see (2.2)) so

$$\mathbb{E}_{v_{1:t-1}} \left(\mathbb{E}_{v_{1:t}} \left(\langle w^{(t)} - w^*, v_t \rangle | v_{1:t-1} \right) \right) = \mathbb{E}_{v_{1:t-1}} \left(\langle w^{(t)} - w^*, \mathbb{E}_{v_{1:t}} (v_t | v_{1:t-1}) \rangle \right)$$

As $w^{(t)}$ is fully determined by $v_{1:t-1}$ then $\mathbb{E}_{v_{1:t}} (v_t | v_{1:t-1}) = \mathbb{E}_{v_{1:t}} (v_t | w^{(t)}) \in \partial f(w^{(t)})$, hence $\mathbb{E}_{v_{1:t}} (v_t | v_{1:t-1})$ is a sub-gradient. By sub-gradient definition

$$\begin{aligned} \mathbb{E}_{v_{1:t-1}} \left(\langle w^{(t)} - w^*, \mathbb{E}_{v_{1:t}} (v_t | v_{1:t-1}) \rangle \right) &\geq \mathbb{E}_{v_{1:t-1}} \left(f(w^{(t)}) - f(w^*) \right) \\ (3.4) \quad &= \mathbb{E}_{v_{1:T}} \left(f(w^{(t)}) - f(w^*) \right) \end{aligned}$$

Hence combining (3.4), (3.3), and (3.2)

$$\mathbb{E} \left(f(w_{\text{SGD}}^{(T)}) - f(w^*) \right) \leq \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left(\langle w^{(t)} - w^*, v_t \rangle \right)$$

Then Lemma 17 in Lect. notes 4 “Gradient descent” implies

$$\mathbb{E} \left(f(w_{\text{SGD}}^{(T)}) - f(w^*) \right) \leq \mathbb{E} \left(\frac{1}{T} \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \frac{1}{T} \sum_{t=1}^T \|v_t\|^2 \right) = \frac{\mathbb{E} \|w^*\|^2}{2\eta T} + \frac{\eta}{2} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \|v_t\|^2$$

□

Note 16. The upper bound in (3.1) depends on $\mathbb{E} \|v_t\|^2$ which has to be bounded and hence the variation of v_t because

$$(3.5) \quad \mathbb{E} \|v_t\|^2 = \mathbb{E} (\|v_t - \mathbb{E}(v_t)\|^2) + \|\mathbb{E}(v_t)\|^2 = \text{tr}(\text{Var}(v_t)) + \|\mathbb{E}(v_t)\|^2$$

where $\mathbb{E} (\|v_t - \mathbb{E}(v_t)\|^2) = \text{tr}(\text{Var}(v_t))$ is the sum of all the variances at each dimension of v_t and $\|\mathbb{E}(v_t)\|^2$ is a finite constant as $\mathbb{E}(v_t) \in \partial_w f(w_t)$ is the unbiased estimator of the sub-gradient by construction.

Proposition 17. (Cont. Proposition 15) Let $f(\cdot)$ be a convex function over the set $\mathcal{W} = \{w : \|w\| \leq B\}$. Let $\mathbb{E} \|v_t\|^2 \leq \rho^2$. Assume we run SGD algorithm of $f(\cdot)$ with learning rate $\eta_t = \sqrt{\frac{B^2}{\rho^2 T}}$ for T steps, and output $w_{\text{SGD}}^{(T)} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$. Then

(1) upper bound on the sub-optimality is

$$(3.6) \quad \mathbb{E} \left(f(w_{\text{SGD}}^{(T)}) \right) - f(w^*) \leq \frac{B\rho}{\sqrt{T}}$$

(2) a given level of accuracy ε such that $\mathbb{E} \left(f(w_{\text{SGD}}^{(T)}) \right) - f(w^*) \leq \varepsilon$ can be achieved after T iterations

$$T \geq \frac{B^2 \rho^2}{\varepsilon^2}.$$

Proof. It follows from Proposition 15. Condition $\mathbb{E} \|v_t\|^2 \leq \rho^2$ can be achieved, for instance by assuming Lipschitz (See Lemma 21 from “Lect. notes 4: Gradient descent”). □

4. IMPLEMENTATION OF SGD IN THE LEARNING PROBLEM 1

Note 18. Note 19 implements SGD to the learning problem (Problem 1)

$$w^* = \arg \min_w (R_g(w)) = \arg \min_w (\mathbb{E}_{z \sim g} (\ell(h_w, z)))$$

Note 19. For a randomly drawn example $z \sim g(\cdot)$, the sub-gradient v of $\ell(w, z)$ at point w is an unbiased estimator of the sub-gradient of the risk $R_g(w)$ at point w . I.e. if $v \in \partial_w \ell(w, z)$ where $z \sim g(\cdot)$ then $\mathbb{E}_{z \sim g} (v|w) \in \partial_w R_g(w)$.

Proof. Let v be a sub-gradient of $\ell(w, z)$ at point w , then

$$(4.1) \quad \ell(u, z) - \ell(w, z) \geq \langle u - w, v \rangle$$

It is

$$\begin{aligned} R_g(u) - R_g(w) &= \mathbb{E}_{z \sim g} (\ell(u, z) - \ell(w, z) | w) \geq \mathbb{E}_{z \sim g} (\langle u - w, v \rangle | w) \\ &= \langle u - w, \mathbb{E}_{z \sim g} (v|w) \rangle \end{aligned}$$

Hence, by definition, v is such that $\mathbb{E}_{z \sim g} (v|w)$ is a sub-gradient of $R_g(w)$. \square

Note 20. Similarly, the average $\bar{v} = \frac{1}{m} \sum_{i=1}^m v_i$ of sub-gradients $v_i \in \partial_w \ell(w, \tilde{z}_i)$ at point w given a randomly drawn set of m examples $\{\tilde{z}_i \sim g(\cdot); i = 1, \dots, m\}$ is an unbiased estimator of the risk function sub-gradient; i.e. $\mathbb{E}_{z \sim g} (\bar{v}|w) \in \partial_w R_g(w)$.

Example 21. Assume a differentiable loss function $\ell(\cdot, z)$ at each z and set sub-gradient $v = \nabla_w \ell(w, z)$, then

$$\mathbb{E}_{z \sim g} (v|w) = \mathbb{E}_{z \sim g} (\nabla_w \ell(h_w, z) | w) = \nabla_w \mathbb{E}_{z \sim g} (\ell(h_w, z) | w) = \nabla_w R_g(w)$$

Note 22. Assume there is available a finite dataset $\mathcal{S}_n = \{z_i; i = 1, \dots, n\}$ of size n which consists of independent realizations z_i from the data generating distribution g ; $z_i \stackrel{\text{ind}}{\sim} g$. SGD (Algorithm 23) is an implementation of the SGD (Algorithm 6) in the learning Problem 1.

Algorithm 23. *Stochastic Gradient Descent with learning rate $\eta_t > 0$, and batch size m , for Problem 1.*

For $t = 1, 2, 3, \dots$ iterate:

- (1) randomly generate a set $\mathcal{J}^{(t)} \subseteq \{1, \dots, n\}^m$ of m indices from 1 to n with or without replacement.
- (2) compute

$$(4.2) \quad w^{(t+1)} = w^{(t)} - \eta_t v_t,$$

where $v_t = \frac{1}{|\mathcal{J}^{(t)}|} \sum_{j \in \mathcal{J}^{(t)}} \nu_{t,j}$ and $\nu_{t,j} \in \partial_w \ell(w^{(t)}, z_j)$ for any $j \in \mathcal{J}^{(t)}$.

- (3) terminate if a termination criterion is satisfied

Note 24. If it is possible to sample anytime fresh examples z_i directly from the data generation model g instead of just having access to only a given finite dataset of examples \mathcal{S}_n , then step 1 in Algorithm 23 can become

- (1) sample $\tilde{z}_j^{(t)} \sim g(\cdot)$ for $j = 1, \dots, m$.

Note 25. Algorithm 23 may be called online SGD when using sub-samples of size one ($m = 1$) and batch SGD when using sub-samples of larger sizes ($m > 1$).

Note 26. In theory, using larger batch size m has the benefit that reduces the variance of v_t at iteration t due to averaging effect, stabilizes the SGD algorithm, and reduces the error bound (3.1); see Remark 16.

Note 27. In practice, for a given fixed computational time, using smaller batch size m has the benefit that the algorithm iterates faster as each iteration processes less number of examples. E.g., consider the extreme cases GD vs online SGD utilized in a scenario of big-data (large training data set): if the dataset consists of several replications of the same values, GD (using all the data) has to process the same information multiple times, while the online SGD (using only one example at a time) would avoid this issue.

Note 28. In practice, for a given fixed computational time, it is possible for a SGD with smaller batch size m to present better generalization properties (wrt the theoretical assumptions) than those with larger m (GD is included). It is observed for the former to be often less prone to getting stuck in shallow local minima because of the additional amount of “noise” E.g., consider the extreme cases GD ($\mathcal{J}^{(t)} \equiv \{1, \dots, n\}$) vs online SGD ($m = 1$) in a scenario with non-convex risk function (e.g. our theoretical assumptions as violated): if the Risk function presents local minima, considering less examples randomly chosen each time may cause fluctuations in the gradient that allow the chain to accidentally jump/escape to an area with a lower minimum.

Proposition 29. (*Implementing Proposition 17*) Consider a convex-Lipschitz-bounded problem $(\mathcal{H}, \mathcal{Z}, \ell)$ with parameters ρ and B . Then for every $\epsilon > 0$, if we run SGD Algorithm 6 to minimize Problem 1 for $T \geq \frac{B^2}{\epsilon^2} \rho^2$ iterations then it produces output $w_{SGD}^{(T)}$ satisfying

$$E\left(R_g\left(w_{SGD}^{(T)}\right)\right) \leq \min_{\forall w \in \mathcal{H}} (R_g(w)) + \epsilon$$

and hence we can achieve PAC guarantees.

Proof. It is just re-writing the formula in part 2 of Proposition 17. Notice that here, condition $E\|v_t\|^2 \leq \rho^2$ is satisfied by the self-bounding property from Lipschitzness of the loss (see Lemma 21 from “Lect. notes 4: Gradient descent”); i.e. if ℓ is ρ -Lipschitz, then $\|v_t\| \leq \rho$ where $v_t \in \partial_w \ell(w^{(t)}, z_j)$ at $j \in \mathcal{J}^{(t)}$. \square

Example 30.¹ We continue Example 26 in Section 4 of Lect. notes 4. Recall, we considered a hypothesis space \mathcal{H} of linear functions $h : \mathbb{R}^2 \rightarrow \mathbb{R}$ with $h(w) = w_1 + w_2x$, $w = (w_1, w_2)^\top$,

¹https://github.com/georgios-stats/Machine_Learning_and_Neural_Networks_III_Epiphanys_2025/tree/main/Lecture_notes/code/example_SGD.R

and $\ell(w, z = (x, y)^\top) = (y_i - w_1 - w_2 x)^2$. Here we consider a big dataset $\mathcal{S}_n = \{z_1, \dots, z_n\}$ with $n = 10^6$ examples.

The batch SGD algorithm (Algorithm 23) with learning rate η_t and batch size $m = 10$ is
For $t = 1, 2, 3, \dots$ iterate:

- (1) Randomly generate a set $\mathcal{J}^{(t)}$ by drawing $m = 10$ numbers from $\{1, \dots, n = 10^6\}$
- (2) compute

$$w^{(t+1)} = w^{(t)} - \eta_t v_t,$$

where

$$v_t = \begin{pmatrix} 2w_1^{(t)} + 2w_2^{(t)} \frac{1}{m} \sum_{j \in \mathcal{J}^{(t)}} x_j - 2 \frac{1}{m} \sum_{j \in \mathcal{J}^{(t)}} y_j \\ 2w_1^{(t)} \bar{x} + 2w_2^{(t)} \frac{1}{m} \sum_{j \in \mathcal{J}^{(t)}} x_j^2 - 2 \sum_{j \in \mathcal{J}^{(t)}} y_j x_j \end{pmatrix}$$

- (3) if $t \geq T = 1000$ STOP

because

$$v_t = \frac{1}{m} \sum_{j \in \mathcal{J}^{(t)}} \nabla_w \ell(w^{(t)}, z_j) = \dots = \begin{pmatrix} 2w_1^{(t)} + 2w_2^{(t)} \frac{1}{m} \sum_{j \in \mathcal{J}^{(t)}} x_j - 2 \frac{1}{m} \sum_{j \in \mathcal{J}^{(t)}} y_j \\ 2w_1^{(t)} \bar{x} + 2w_2^{(t)} \frac{1}{m} \sum_{j \in \mathcal{J}^{(t)}} x_j^2 - 2 \sum_{j \in \mathcal{J}^{(t)}} y_j x_j \end{pmatrix}$$

In Figures A.1a & A.1b, we observe that increasing the batch size has improved the convergence however this is not a panacea. Also it had reduced the oscillations of chain $\{w^{(t)}\}$.

APPENDIX A. PLOTS

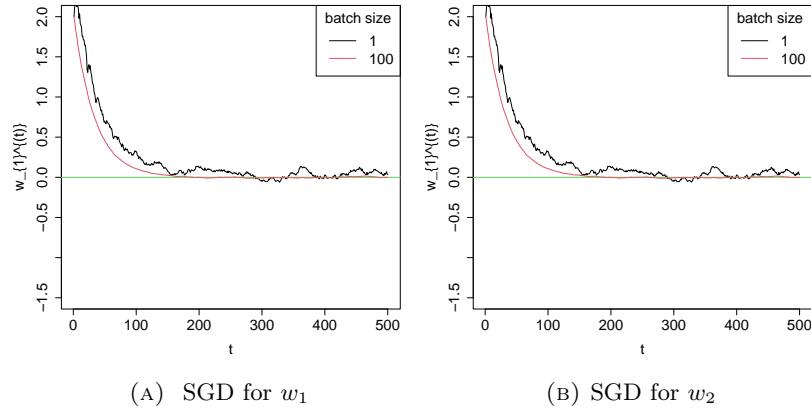


FIGURE A.1. Simulations of the Example 30

Lecture notes 6: Variations of stochastic gradient descent

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim. To introduce the stochastic gradient descent variations for restricted sets, adaptive step functions, and acceleration/variance reduction.

Reading list & references:

- (1) Johnson, Rie, and Tong Zhang. "Accelerating stochastic gradient descent using predictive variance reduction." *Advances in neural information processing systems* 26 (2013).
- (2) Duchi, John, Elad Hazan, and Yoram Singer. "Adaptive subgradient methods for online learning and stochastic optimization." *Journal of machine learning research* 12, no. 7 (2011).

1. STOCHASTIC GRADIENT DESCENT WITH PROJECTION

Problem 1. Consider the minimization problem

$$(1.1) \quad w^* = \arg \min_w (f(w))$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$, $w \in \mathcal{W}$, and assume $\mathcal{W} \subset \mathbb{R}^d$ is a restricted domain.

Note 2. In Problem 1 if for instance, $\mathcal{W} = \{w \in \mathbb{R}^d : \|w\| \leq B\}$ for $B > 0$, f is convex over \mathcal{W} , but f is non-convex in $\mathbb{R}^d / \mathcal{W}$, direct implementation of (standard) SGD (Algorithm 6, Lect. Notes 5) may produce a chain stepping out \mathcal{W} and hence an output $w_{\text{SGD}} \notin \mathcal{W}$.

Note 3. Stochastic gradient descent with projection (prjSGD) is a modified SGD suitable for applications with restricted sets and relies on the consideration of a projection step guaranteeing $w \in \mathcal{W}$

Algorithm 4. *Stochastic Gradient Descent with learning rate $\eta_t > 0$ and with projection in \mathcal{W} for Problem 1*

For $t = 1, 2, 3, \dots$ iterate:

- (1) compute

$$(1.2) \quad w^{(t+\frac{1}{2})} = w^{(t)} - \eta_t v_t,$$

where v_t is a random vector such that $E(v_t | w^{(t)}) \in \partial f(w^{(t)})$

- (2) compute

$$(1.3) \quad w^{(t+1)} = \arg \min_{w \in \mathcal{W}} (\|w - w^{(t+\frac{1}{2})}\|)$$

- (3) terminate if a termination criterion is satisfied

Note 5. The analysis of SGD with projection is given as a Exercise 13 from the Exercise sheet.

Example 6.¹ Consider a (rather naive) loss function $\ell(w, z = (x, y)) = -\cos(0.5(y - w_1 - w_2 x))$, a hypothesis class $\mathcal{H} = \{h_w(x) = w^\top x : w \in \mathcal{W}\}$ and $\mathcal{W} = \{w \in \mathbb{R}^2 : \|w\| \leq 1.5\}$, and assume that inputs x in dataset S are such that $x \in [-1, 1]$. Note that $-\cos(\cdot)$ is convex in $[-1.5, 1.5]$ and non-convex in \mathbb{R} . Consider learning rate $\eta_t = 50/t$ reducing to zero, and seed $w^{(0)} = (1.5, 1.5)$. An unconstrained SGD may produce a minimizer/solution outside \mathcal{H} because η_t is too large at the first few iterations. We can design the online SGD ($m = 1$) with projection to \mathcal{H} as

For $t = 1, 2, 3, \dots$ iterate:

- (1) randomly generate a set $\mathcal{J}^{(t)} = \{j^*\}$ by drawing one number from $\{1, \dots, n = 10^6\}$, and set $\tilde{\mathcal{S}}_1 = \{z_{j^*}\}$
- (2) compute

$$w^{(t+1/2)} = w^{(t)} - \frac{25}{t} \begin{pmatrix} \sin(0.5(y_{j^*} - w_1^{(t)} - w_2^{(t)} x_{j^*})) \\ \sin(0.5(y_{j^*} - w_1^{(t)} - w_2^{(t)} x_{j^*})) x_{j^*} \end{pmatrix}$$

$$w^{(t+1/2)} = \arg \min_{\|w\| \leq 1.5} (\|w - w^{(t+1/2)}\|)$$

- (3) if $t \geq T = 1000$ STOP

because

$$v_t = \nabla \ell(w^{(t)}, z_{j^*}) = \begin{pmatrix} \frac{\partial}{\partial w_1^{(t)}} \ell(w^{(t)}, z_{j^*}) \\ \frac{\partial}{\partial w_2^{(t)}} \ell(w^{(t)}, z_{j^*}) \end{pmatrix} = \begin{pmatrix} 0.5 \sin(0.5(y_{j^*} - w_1^{(t)} - w_2^{(t)} x_{j^*})) \\ 0.5 \sin(0.5(y_{j^*} - w_1^{(t)} - w_2^{(t)} x_{j^*})) x_{j^*} \end{pmatrix}$$

In Figures A.1b & A.1e, we observe that the SGD got trapped outside \mathcal{H} due to the unreasonably large learning rate at the beginning of the iterations, while the SGD with projection step managed to stay in \mathcal{H} and converge.

2. PRECONDITIONED SGD ABD THE ADAGRAD ALGORITHM

Problem 7. Consider the minimization problem

$$(2.1) \quad w^* = \arg \min_{w \in \mathcal{W}} (f(w))$$

where $f : \mathcal{W} \rightarrow \mathbb{R}$, and assume $\mathcal{W} \equiv \mathbb{R}^d$.

Note 8. All SGD (introduced earlier) are first order methods, in the sense they consider only the gradient, and hence they ignore the curvature of the space. Consequently they may present slow converge in cases the curvature changes, eg. among dimensions of w .

Note 9. Preconditioned Stochastic Gradient Descent (Algorithm 10) is a modified SGD aiming to expedite convergence of SGD by using a preconditioner P_t that accounts for the curvature (or geometry in general) of f .

¹https://github.com/georgios-stats/Machine_Learning_and_Neural_Networks_III_Epiphanys_2025/tree/main/Lecture_notes/code/example_projSGD.R

Algorithm 10. Preconditioned Stochastic Gradient Descent with learning rate $\eta_t > 0$, and preconditioner P_t for the solution of the minimization problem (3.1)

For $t = 1, 2, 3, \dots$ iterate:

(1) compute

$$(2.2) \quad w^{(t+1)} = w^{(t)} - \eta_t P_t v_t,$$

where v_t is a random vector such that $E(v_t|w^{(t)}) \in \partial f(w^{(t)})$, and P_t is a preconditioner

(2) terminate if a termination criterion is satisfied, e.g.

If $t \geq T$ then STOP

Note 11. A natural choice of P_t can be $P_t := [H_t + \epsilon I_d]^{-1}$, where H_t is the Hessian matrix $[H_t]_{i,j} = \frac{\partial^2}{\partial w_i \partial w_j} f(w)|_{w=w^{(t)}}$ (ie. the gradient of the gradient's elements), and ϵ is a tiny $\epsilon > 0$ to mitigate machine error when Hessian elements are close to zero.

Note 12. Using the inverse of the full Hessian as preconditioner P_t may be too expensive to perform matrix operations in (2.2) and operations can be too unstable/inaccurate due to the random error induced by the stochasticity of the gradient.

2.1. Adaptive Stochastic Gradient Decent (AdaGrad).

Note 13. AdaGrad (Algorithm 15) is a modified SGD algorithm that aims to adaptively adjust the learning rate of SGD, perform more informative gradient-based learning, and by improve convergence performance over standard SGD by dynamically incorporating knowledge of the geometry of function $f(\cdot)$ in earlier iterations.

Note 14. It is particularly suitable in cases where the data are sparse and sparse features w 's are more informative. It aims to perform larger updates (i.e. high learning rates) for those dimensions of w that are related to infrequent features (largest partial derivative) and smaller updates (i.e. low learning rates) for frequent ones (smaller partial derivative).

Algorithm 15. Adaptive Stochastic Gradient Decent (AdaGrad) is a preconditioned SGD (Algorithm 10) with preconditioner $P_t = [\text{diag}(G_t) + \epsilon I_d]^{-1/2}$ where notation $\text{diag}(A)$ denotes a $d \times d$ matrix whose diagonal is the d dimensional diagonal vector $(A_{1,1}, A_{2,2}, \dots, A_{d,d})$ of $d \times d$ matrix A and whose off-diagonal elements are zero, $G_t = \sum_{\tau=1}^t v_\tau v_\tau^\top$ is the sum of the outer products of the gradients $\{v_\tau; \tau \leq t\}$ up to the state t , and $\epsilon > 0$ is a tiny value (eg, 10^{-6}) set for computational stability in case the gradient becomes too close to zero. I.e.

$$(2.3) \quad w^{(t+1)} = w^{(t)} - \eta_t [\text{diag}(G_t) + \epsilon I_d]^{-1/2} v_t.$$

Note 16. AdaGrad individually adapts the learning rate of each dimension of w_t by scaling them inversely proportional to the square root of the sum of all the past squared values of the gradient

$\{v_\tau; \tau \leq t\}$. This is because

$$(2.4) \quad [G_t]_{j,j} = \sum_{\tau=1}^t (v_{\tau,j})^2$$

where j denotes the j -th dimension of w . Hence (2.3) and (2.3) imply that the j -th dimension of w is updated as

$$w_j^{(t+1)} = w_j^{(t)} - \eta_t \frac{1}{\sqrt{[G_t]_{j,j} + \epsilon}} v_{t,j}.$$

Note 17. The accumulation of positive terms in (2.4) makes the sum keep growing during training and causes the learning rate to shrink and become infinitesimally small. This offers an automatic way to choose a decreasing learning rate simplifying setting the learning rate; however it may result in a premature and excessive decrease in the effective learning rate. This can be mitigated by still considering in (2.3) a (user specified rate) $\eta_t \geq 0$ and tuning it properly via pilot runs.

Example 18. ² AdaGrad with $\eta_t = 1$, $\epsilon = 10^{-6}$, and batch size $m = 1$ is

- Set $G_0 = (0, 0)^\top$.
- For $t = 1, 2, 3, \dots$ iterate:
 - (1) randomly generate $j^{(t)}$ from $\{1, \dots, n = 10^6\}$
 - (2) compute

$$\begin{aligned} v_t &= \begin{pmatrix} 2w_1^{(t)} + 2w_2^{(t)}x_{j^{(t)}} - 2y_{j^{(t)}} \\ 2w_1^{(t)}\bar{x} + 2w_2^{(t)}x_{j^{(t)}}^2 - 2y_{j^{(t)}}x_{j^{(t)}} \end{pmatrix} \\ G_t &= G_{t-1} + \begin{pmatrix} v_{t,1}^2 \\ v_{t,2}^2 \end{pmatrix} \\ w^{(t+1)} &= \begin{pmatrix} w_1^{(t)} - \frac{\eta_t}{\sqrt{G_{t,1} + \epsilon}} v_{t,1} \\ w_2^{(t)} - \frac{\eta_t}{\sqrt{G_{t,2} + \epsilon}} v_{t,2} \end{pmatrix}, \end{aligned}$$

- (3) if $t \geq T = 1000$ STOP

because

$$v_t = \nabla \ell(w^{(t)}, z_{j^{(t)}}) = \begin{pmatrix} \frac{\partial}{\partial w_1^{(t)}} \ell(w^{(t)}, z_{j^{(t)}}) \\ \frac{\partial}{\partial w_2^{(t)}} \ell(w^{(t)}, z_{j^{(t)}}) \end{pmatrix} = \begin{pmatrix} 2w_1^{(t)} + 2w_2^{(t)}x_{j^{(t)}} - 2y_{j^{(t)}} \\ 2w_1^{(t)}\bar{x} + 2w_2^{(t)}x_{j^{(t)}}^2 - 2y_{j^{(t)}}x_{j^{(t)}} \end{pmatrix}$$

In Figures A.1g & A.1h, we see that AdaGrad with $\eta = 1$ works (I did not try to tune it), however to make vanilla SGD to work I have to tune $\eta = 0.03$ otherwise for $\eta = 1.0$ it did not work.

²https://github.com/georgios-stats/Machine_Learning_and_Neural_Networks_III_Epiphanys_2025/tree/main/Lecture_notes/code/example_AdaGrad.R

3. STOCHASTIC VARIANCE REDUCED GRADIENT (SVRG)

Note 19. (Recall) Consider minimization problem $w^* = \arg \min_w \{f(w)\}$. Consider SGD with learning rate $\eta > 0$, output $w_{\text{SGD}}^{(T)}$, and general recursion

$$w^{(t+1)} = w^{(t)} - \eta_t v_t$$

where $E(v_t|w^{(t)}) \in \partial f(w^{(t)})$. Recall that the upper bound of the error depends on the variance of the stochastic sub-gradient i.e. $E\|v_t - E(v_t)\|^2$.

$$(3.1) \quad E\left(f\left(w_{\text{SGD}}^{(T)}\right)\right) - f(w^*) \leq \frac{\|w^*\|^2}{2\eta T} + \frac{\eta}{2} \frac{1}{T} \sum_{t=1}^T E\|v_t\|^2$$

$$(3.2) \quad E\|v_t\|^2 = \underbrace{E\|v_t - E(v_t)\|^2}_{=\text{tr}(\text{Var}(v_t))} + \underbrace{\|E(v_t)\|^2}_{\text{const.}}$$

A way to improve the SGD may be by reducing / controlling the variances at each dimension of v_t .

Problem 20. Consider a learning problem $(\mathcal{H}, \mathcal{Z}, \ell)$ with risk function $R_g(w) = E_{z \sim g}(\ell(h_w, z))$ given an unknown data generating model $g(\cdot)$, using a known tractable loss $\ell(\cdot, \cdot)$, and hypothesis $h_w \in \mathcal{H}$. Consider a training data set $\mathcal{S}_n = \{z_i = (x_i, y_i); z_i \in \mathcal{Z}, i = 1, \dots, n\}$. We need to find minimizer $w^* = \arg \min_w (R_g(w))$.

3.1. Stochastic gradient descent using (mini-batches).

Note 21. A way to reduce variance in (3.2) (and hence achieve faster convergence rate) is to consider a larger batch size $m = |\mathcal{J}^{(t)}|$. Note that

$$E_{z \sim g}\|v_t - E(v_t)\|^2 = E_{z \sim g} \left\| \nabla_w \frac{1}{|\mathcal{J}^{(t)}|} \sum_{j \in \mathcal{J}^{(t)}} \ell(w, z_j) - \nabla_w R_g(w) \right\|^2 = \frac{E_{z_j \sim g} \|\nabla_w \ell(w, z_j) - \nabla_w R_g(w)\|^2}{|\mathcal{J}^{(t)}|}$$

where top part of the fraction is the gradient variance for SGD with $m = 1$. Using larger batch size $m = |\mathcal{J}^{(t)}|$ reduces the stochastic gradient variance and hence the convergence rate.

3.2. Stochastic Variance Reduced Gradient (SVRG).

Control variate methods.

Note 22. Control variate is a general way to perform variance reduction. Let $x \in \mathbb{R}^d$ be an unbiased estimator of θ . The goal is to compute $v = x + c(y - E(y))$ by finding a suitable control variable $y \in \mathbb{R}^d$ and constant $c \in \mathbb{R}$ such that $\text{tr}(\text{Var}_c(v)) < \text{tr}(\text{Var}(x))$ with purpose to use it as an unbiased estimator of θ instead of x $E_c(v) = E(x) = \theta$.

Example 23. Let random variables $x \in \mathbb{R}^d$, and $y \in \mathbb{R}^d$. Let $v = x + c(y - E(y))$ for some constant $c \in \mathbb{R}$. It is $E_c(v) = E(x)$. Also

$$E_c\|v - E_c(v)\|^2 = E\|x - E(x)\|^2 + c^2 E\|y - E(y)\|^2 + 2c E\langle x - E(x), y - E(y) \rangle$$

which is minimized for

$$0 = \frac{d}{dc} E_c\|v - E_c(v)\|^2 \Big|_{c=c^*} \implies c^* = -\frac{E\langle x - E(x), y - E(y) \rangle}{E\|y - E(y)\|^2}$$

with

$$(3.3) \quad E_{c^*} \|v - E_c(v)\|^2 = E \|x - E(x)\|^2 - \frac{E \langle x - E(x), y - E(y) \rangle}{E \|y - E(y)\|^2}$$

Variance reduced methods.

Definition 24. A SGD algorithm with learning rate $\eta > 0$ and general recursion

$$w^{(t+1)} = w^{(t)} - \eta_t v_t, \quad t = 1, 2, \dots$$

where $E(v_t | w^{(t)}) \in \partial f(w^{(t)})$ is called variance reduced SGD if it presents smaller stochastic gradient variance compared to standard SGD and has the property

$$E \|v_t - \nabla f(w^{(t)})\|^2 \rightarrow 0, \text{ as } t \rightarrow \infty$$

where $E(\cdot)$ is w.r.t. all random variables till iteration t .

Stochastic Variance Reduced Gradient (SVRG) .

Algorithm 25. Stochastic Variance Reduced Gradient with learning rate $\eta_t > 0$ for Problem 1.

For $t = 1, 2, 3, \dots$ iterate:

- (1) randomly draw an example $z_{j^{(t)}}$ indexed by $j^{(t)} \in \{1, \dots, n\}$.
- (2) compute

$$(3.4) \quad w^{(t+1)} = w^{(t)} - \eta_t v_t$$

where

$$(3.5) \quad v_t = \underbrace{\nabla \ell(w^{(t)}, z_{j^{(t)}})}_{=x} - \underbrace{\left(\nabla \ell(\tilde{w}, z_{j^{(t)}}) - \frac{1}{n} \sum_{i=1}^n \nabla \ell(\tilde{w}, z_i) \right)}_{=y}$$

- (3) if $\text{modulo}(t, \kappa) = 0$, set $\tilde{w} = w^{(t)}$
- (4) if a termination criterion is satisfied STOP

Note 26. Prepend, we do not know c in (3.3). The variance of gradient v_t^{SVRG} is

$$\begin{aligned} E_{z \sim g} \|v_t^{\text{SVRG}} - E(v_t^{\text{SVRG}})\|^2 &= E_{z \sim g} \left\| \overbrace{\nabla_w \ell(w^{(t)}, z_{j^{(t)}}) + c \left(\nabla_w \ell(\tilde{w}, z_{j^{(t)}}) - \frac{1}{n} \sum_{i=1}^n \nabla_w \ell(\tilde{w}, z_i) \right)}^{v_t^{\text{SVRG}} =} - \nabla_w R_g(w^{(t)}) \right\|^2 \\ &\leq E_{z \sim g} \left\| \nabla_w \ell(w^{(t)}, z_{j^{(t)}}) + c \nabla_w \ell(\tilde{w}, z_{j^{(t)}}) \right\|^2 \end{aligned}$$

As optimal c is often intractable, we sub-optimally set $c = -1$ and get

$$E_{z \sim g} \|v_t^{\text{SVRG}} - E(v_t^{\text{SVRG}})\|^2 \leq E_{z \sim g} \left\| \nabla_w \ell(w^{(t)}, z_{j^{(t)}}) - \nabla_w \ell(\tilde{w}, z_{j^{(t)}}) \right\|^2$$

If $\ell(\cdot, z_j)$ is β_j -smooth, then

$$\mathbb{E}_{z \sim g} \|v_t^{\text{SVRG}} - \mathbb{E}(v_t^{\text{SVRG}})\|^2 \leq \max_j (\{\beta_j\}) \|w^{(t)} - \tilde{w}\|^2$$

Hence to bound/control variance, point \tilde{w} should be chosen to be as close as possible to $w^{(t)}$. Perhaps by taking a snapshot $\tilde{w} \leftarrow w^{(t)}$ every κ -th iteration.

Note 27. The frequency of the snapshots κ is specified by the researcher. The smaller the κ , the more frequent snapshots, and the more correlated the baseline y will be with the objective x in (3.5) and hence the bigger the performance improvement (See (3.3)); however the iterations will be slower.

Note 28. Iterations of SVRG are computationally faster than those of full GD but SVRG can still match the theoretical convergence rate of GD.

Example 29.³ The SVRG with learning rate $\eta_t > 0$ and batch size $m = 1$ is

For $t = 1, 2, 3, \dots$ iterate:

- (1) randomly generate a set $\mathcal{J}^{(t)} = \{j^{(t)}\}$ by drawing one number $j^{(t)}$ from $\{1, \dots, n = 10^6\}$.
- (2) compute

$$(3.6) \quad w^{(t+1)} = \begin{bmatrix} w_1^{(t)} \\ w_2^{(t)} \end{bmatrix} - \eta_t \left(\begin{bmatrix} 2(w_1^{(t)} - \tilde{w}_1) + 2(w_2^{(t)} - \tilde{w}_2)x_{j^{(t)}} \\ 2(w_2^{(t)} - \tilde{w}_2)x_{j^{(t)}} + 2(w_2^{(t)} - \tilde{w}_2)(x_{j^{(t)}})^2 \end{bmatrix} + \nabla \hat{R}_{\mathcal{S}}(\tilde{w}) \right)$$

- (3) if modulo $(t, \kappa) = 0$, set $\tilde{w} = w^{(t)}$

- (4) compute

$$(3.7) \quad \frac{1}{n} \sum_{i=1}^n \ell(\tilde{w}, z_i) = \begin{pmatrix} 2\tilde{w}_1 + 2\tilde{w}_2\bar{x} - 2\bar{y} \\ 2\tilde{w}_1\bar{x} + 2\tilde{w}_2\bar{x}^2 - 2y^\top x \end{pmatrix} = \nabla \hat{R}_{\mathcal{S}}(\tilde{w})$$

- (5) if a termination criterion is satisfied STOP

Because (3.7) is actually the gradient of the Risk function at \tilde{w} and

$$\nabla \ell(w^{(t)}, z_{j^{(t)}}) - \nabla \ell(\tilde{w}, z_{j^{(t)}}) = \begin{pmatrix} 2(w_1^{(t)} - \tilde{w}_1) + 2(w_2^{(t)} - \tilde{w}_2)x_{j^{(t)}} \\ 2(w_2^{(t)} - \tilde{w}_2)x_{j^{(t)}} + 2(w_2^{(t)}(x_{j^{(t)}})^2 - \tilde{w}_2(x_{j^{(t)}})^2) \end{pmatrix}$$

In Figure A.1c we observe the frequency of the snapshots has improved the convergence; however this is not a panacea as seen in Figure A.1f.

³https://github.com/georgios-stats/Machine_Learning_and_Neural_Networks_III_Epiphanys_2025/tree/main/Lecture_notes/code/example_SVRG.R

APPENDIX A. PLOTS

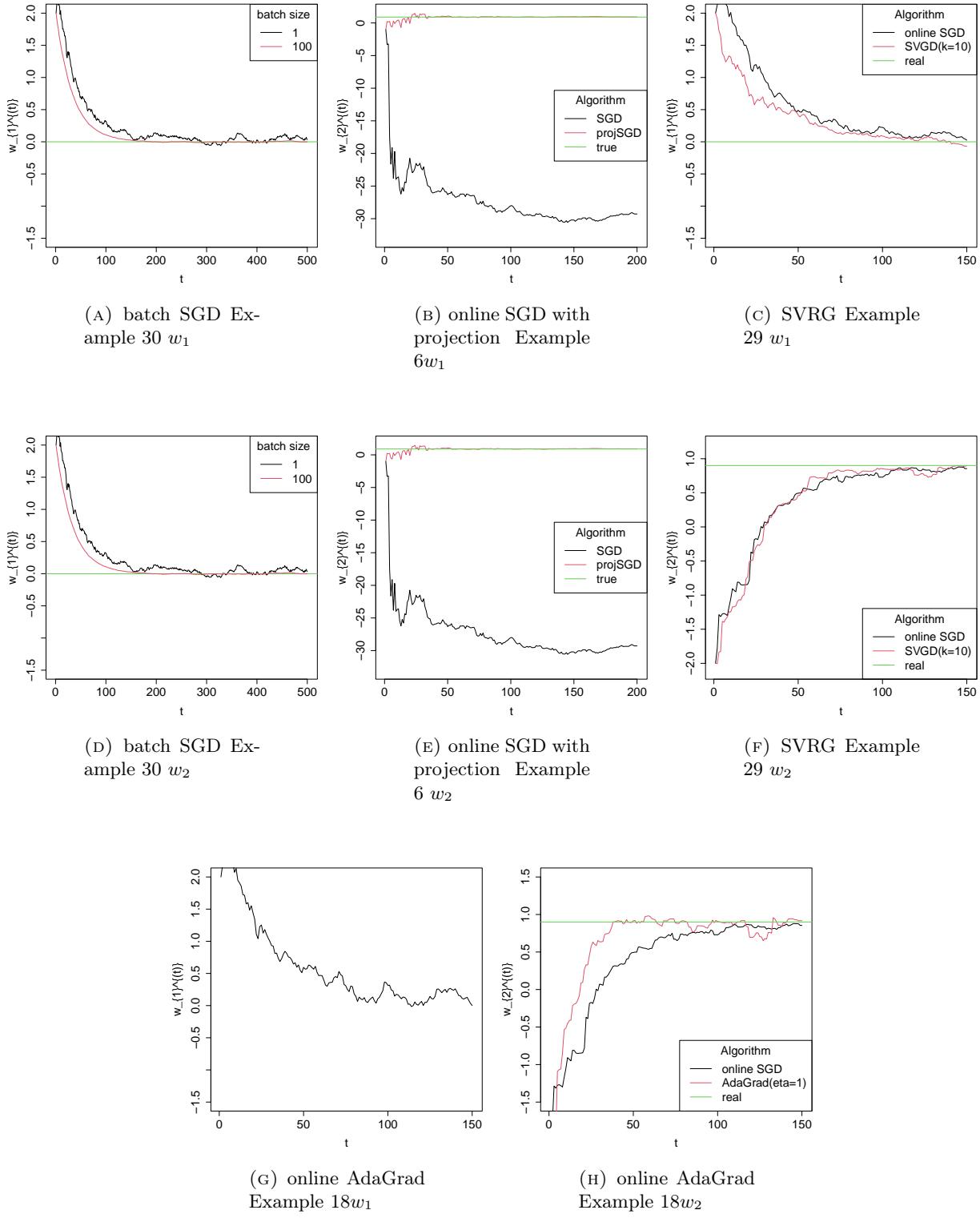


FIGURE A.1. Simulations of the Examples

Lecture notes 7: Bayesian Learning via Stochastic gradient and Stochastic gradient Langevin dynamics

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim. To introduce the Bayesian Learning, and Stochastic gradient Langevin dynamics (description, heuristics, and implementation). Stochastic gradient descent will be implemented in the Bayesian learning too.

Reading list & references:

- (1) Welling, M., & Teh, Y. W. (2011). Bayesian learning via stochastic gradient Langevin dynamics. In Proceedings of the 28th international conference on machine learning (ICML-11) (pp. 681-688).
- (2) Bottou, L. (2012). Stochastic gradient descent tricks. In Neural networks: Tricks of the trade (pp. 421-436). Springer, Berlin, Heidelberg.
- (3) Sato, I., & Nakagawa, H. (2014). Approximation analysis of stochastic gradient Langevin dynamics by using Fokker-Planck equation and Ito process. In International Conference on Machine Learning (pp. 982-990). PMLR.
- (4) Teh, Y. W., Thiery, A. H., & Vollmer, S. J. (2016). Consistency and fluctuations for stochastic gradient Langevin dynamics. Journal of Machine Learning Research, 17.
- (5) Vollmer, S. J., Zygalakis, K. C., & Teh, Y. W. (2016). Exploration of the (non-) asymptotic bias and variance of stochastic gradient Langevin dynamics. The Journal of Machine Learning Research, 17(1), 5504-5548. →further reading for theory
- (6) Nemeth, C., & Fearnhead, P. (2021). Stochastic gradient Markov chain Monte Carlo. Journal of the American Statistical Association, 116(533), 433-450.

1. BAYESIAN LEARNING AND MOTIVATIONS

Note 1. Bayesian methods are appealing in their ability to capture uncertainty in learned parameters and avoid overfitting. Arguably with large datasets there will be little overfitting. Alternatively, as we have access to larger datasets and more computational resources, we become interested in building more complex models (eg from a logistic regression to a deep neural network), so that there will always be a need to quantify the amount of parameter uncertainty.

Note 2. Consider a Bayesian statistical model with sampling distribution (statistical model) $f(z|w)$ labeled by an unknown parameter $w \in \mathcal{W} \subseteq \mathbb{R}^d$ that follows a prior distribution $f(w)$. Assume a dataset $\mathcal{S}_n = \{z_i; i = 1, \dots, n\}$ of size n containing independently drawn examples. Let $L_n(w) := f(z_{1:n}|w)$ denote the likelihood of the observables $\{z_i \in \mathcal{Z}\}_{i=1}^n$ given parameter w . The Bayesian

model is denoted as

$$(1.1) \quad \begin{cases} z_i|w & \stackrel{\text{ind}}{\sim} f(z_i|w), \ i = 1, \dots, n \\ w & \sim f(w) \end{cases}$$

Note 3. With regards to the Bayesian model in Note 2, we denote the likelihood of the observables $\{z_i \in \mathcal{Z}\}_{i=1}^n$ given the parameter w as

$$L_n(w) := f(z_{1:n}|w) = \prod_{i=1}^n f(z_i|w)$$

Note 4. Bayesian learning (inference) relies on the posterior distribution density

$$(1.2) \quad f(w|z_{1:n}) = \frac{L_n(w) f(w)}{\int L_n(w) f(w) dw}$$

which quantifies the researcher's belief (or uncertainty) about the unknown parameter w learned given examples $\{z_i \in \mathcal{Z}\}_{i=1}^n$. It is often intractable; hence there is often a need to numerically compute it. (Section 3)

Note 5. Point estimation of a function h of w is often performed via computation of the posterior expectation w given the examples in \mathcal{S}_n

$$(1.3) \quad E_f(h(w)|z_{1:n}) = \int h(w) f(w|z_{1:n}) dw$$

It is often intractable; hence there is often a need to numerically compute it. (Section 3)

Note 6. Point estimation of w can also be performed via maximum a-posteriori (MAP) estimator w^* of w that is the maximizer w^* of (1.2) i.e.

$$(1.4) \quad w^* = \arg \max_w (f(w|z_{1:n}))$$

$$(1.5) \quad = \arg \max_w \left(\underbrace{\log(L_n(w))}_{(\text{I})} + \underbrace{\log(f(w))}_{(\text{II})} \right)$$

Note that (I) may be interpreted as an empirical risk function, and (II) can be interpreted as a shrinkage term in terms of shrinkage methods (like LASSO, Ridge). It is often intractable; hence there is often a need to numerically compute it. (Section 2)

Note 7. To describe the learning algorithms Gradient Descent (GD), Stochastic Gradient Descent (SGD), and Stochastic Gradient Langevin Dynamics (SGLD), we consider that the examples (data) in (1.1) are independent realizations from the sampling distribution i.e. $z_i|w \sim f(\cdot|w)$ for $i = 1, \dots, n$ and that w is continuous (not discrete).

Note 8. In what follows, we first present the implementation of GD and SGD addressing MAP learning, and then we introduce the implementation of SGLD addressing posterior density and expectation learning.

2. MAXIMUM A POSTERIORI (MAP) LEARNING VIA GD AND SGD

Problem 9. Given the Bayesian model (1.1), and rearranging (1.4), MAP estimate w^* of w can be computed as

$$(2.1) \quad w^* = \arg \max_w (\log(L_n(w)) + f(w)) = \arg \max_w \left(\sum_{i=1}^n \log(f(z_i|w)) + \log(f(w)) \right)$$

Note 10. GD is particularly suitable to solve (2.1) when w has high dimensionality.

Algorithm 11. Gradient descent implementation in Problem 9 with learning rate $\eta_t \geq 0$.

For $t = 1, 2, 3, \dots$ iterate:

- (1) Compute

$$(2.2) \quad w^{(t+1)} = w^{(t)} + \eta_t \left(\sum_{i=1}^n \nabla_w \log(f(z_i|w^{(t)})) + \nabla_w \log(f(w^{(t)})) \right)$$

Note 12. SGD is particularly suitable to solve (2.1) when w has high dimensionality, and in big-data problems since the repetitive computation of the sum in (2.2) is prohibitively expensive. Yet consider the benefits of SGD against GD as discussed in (Notes 27 and 28 of Lect. Notes 5).

Algorithm 13. Stochastic Gradient Descent implementation in with learning rate $\eta_t \geq 0$ and batch size m

For $t = 1, 2, 3, \dots$ iterate:

- (1) generate a random set $\mathcal{J}^{(t)} \subseteq \{1, \dots, n\}^m$ of m indices from 1 to n with or without replacement
- (2) compute

$$(2.3) \quad w^{(t+1)} = w^{(t)} + \eta_t \left(\frac{n}{m} \sum_{j \in \mathcal{J}^{(t)}} \nabla_w \log(f(z_j|w^{(t)})) + \nabla_w \log(f(w^{(t)})) \right)$$

Note 14. Recursion (2.3) is justified in terms of SGD theory as

$$(2.4) \quad \mathbb{E}_{\mathcal{J}^{(t)} \sim \text{simple-random-sampling}} \left(\frac{n}{m} \sum_{j \in \mathcal{J}^{(t)}} \nabla_w \log(f(z_j|w^{(t)})) \right) = \sum_{i=1}^n \nabla_w \log(f(z_i|w^{(t)}))$$

see Exercise 15 from the Exercise sheet.

Note 15. The implementation of the SGD variants (Algorithms 23, 25, 10, 15 in Handout 3)) is straightforward based on Algorithm 13 and (2.4).

3. FULLY BAYESIAN LEARNING VIA SGLD

Problem 16. Fully Bayesian learning, computationally, is the problem of recovering the posterior distribution $f(w|z_{1:n})$ of w given $z_{1:n}$ that admits density (1.2). For a given Bayesian model (1.1),

the Bayesian estimator of $h := h(w)$ is often computed as the posterior expectation of $h(w)$ given the data \mathcal{S}_n

$$(1.3) \quad \text{E}_f(h(w)|z_{1:n}) = \int h(w) f(w|z_{1:n}) dw$$

Note 17. Monte Carlo integration aims at approximating (1.3) when intractable, by using Central Limit Theorem or Law of Large Numbers arguments as $\hat{h}_T \approx \text{E}_f(h(w)|z_{1:n})$ where

$$(3.1) \quad \hat{h}_T = \frac{1}{T} \sum_{t=1}^T h(w^{(t)})$$

where $\{w^{(t)}\}$ are T simulations drawn (approximately) from the posterior distribution (1.2). This theory is subject to conditions we skip.

Note 18. Stochastic gradient Langevin dynamics (SGLD) algorithm generates as output a random chain $\{w^{(t)}\}$ approximately distributed according to a distribution admitting density such as

$$(3.2) \quad f_\tau(w|z_{1:n}) \propto \exp\left(\frac{1}{\tau} (\log L_n(w) + \log f(w))\right)$$

$$(3.3) \quad \propto \exp\left(\frac{1}{\tau} \left(\sum_{i=1}^n \log(f(z_i|w)) + \log(f(w)) \right)\right)$$

for any $\tau > 0$ under regularity conditions. That allows to recover the whole posterior distribution (1.2) (hence account for uncertainty in $h = h(w)$) and approximate posterior expectations based on the Monte Carlo integration (Note 17).

Algorithm 19. *Stochastic Gradient Langevin Dynamics (SGLD) with learning rate $\eta_t > 0$, batch size m , and temperature $\tau > 0$ is*

-
- (1) Generate a random set $\mathcal{J}^{(t)} \subseteq \{1, \dots, n\}^m$ of m indices from 1 to n with or without replacement
 - (2) Compute

$$(3.4) \quad w^{(t+1)} = w^{(t)} + \eta_t \left(\frac{n}{m} \sum_{i \in \mathcal{J}^{(t)}} \nabla \log f(z_i|w^{(t)}) + \nabla \log f(w^{(t)}) \right) + \sqrt{2\eta_t \tau} \epsilon_t$$

where $\epsilon_t \stackrel{iid}{\sim} N(0, I)$.

- (3) Terminate if a termination criterion is satisfied; E.g., $t \geq T_{max}$ for a prespecified $T_{max} > 0$.

Ref [1]

How / why it works...

Note 20. (Mathematically speaking) The stochastic chain in (3.4) can be viewed as a discretization of the continuous-time Langevin diffusion described by the stochastic differential equation Ref [3]

$$(3.5) \quad dW(t) = -\nabla_w [-\log f(W(t)|z_{1:n})] dt + \sqrt{2\tau} dB(t), \quad t \geq 0$$

where $\{B(t)\}$ is a standard Brownian motion¹ in \mathbb{R}^d (i.e.). Under suitable assumptions on f , it can be shown that a Gibbs distribution with PDF such as

$$(3.6) \quad f^*(w|z_{1:n}) \propto \exp\left(-\frac{1}{\tau}[-\log f(w|z_{1:n})]\right)$$

is the unique invariant distribution of (3.5), and that the distributions of $W(t)$ converge rapidly to f^* as $t \rightarrow \infty$.

Note 21. (Heuristically speaking) SGLD relies on injecting the ‘right’ amount of noise (3.4) to a standard stochastic gradient optimization recursion (2.2), such that, as the stepsize η_t properly reduces, the produced chain $\{w^{(t+1)}\}$ converges to samples that could have been drawn from distribution with density (3.2). In the initial phase of running, the stochastic gradient noise will dominate the injected noise ϵ_t and the algorithm will imitate an efficient SGD Algorithm 11 -but this is until η_t or $\nabla \log(L_n(w))$ become small enough. In the later phase of running, the injected noise ϵ_t will dominate the stochastic gradient noise, so the SGLD will imitate a Langevin dynamics for the target distribution (1.2). The aim is for the algorithm to transition smoothly between the two phases. Whether the algorithm is in the stochastic optimization phase or Langevin dynamics phase depends on the variance of the injected noise versus that of the stochastic gradient.

Note 22. One can argue that, the output of SGLD is also an “almost” maximizer of the (minus) empirical risk with shrinkage for large enough t . A draw from the Gibbs distribution (3.6) is approximately a maximizer of (2.1). Also one can show that the SGLD recursion tracks the Langevin diffusion (3.5) in a suitable sense. Hence, both imply that the distributions of $W(t)$ will be close to the Gibbs distribution (3.6) for all sufficiently large t .

About the learning rate alg. parameter...

Note 23. To guarantee the algorithm to work it is important for the step sizes η_t to decrease to zero, so that the mixing rate of the algorithm will slow down with increasing number of iterations t . Then, we can keep the step size η_t constant once it has decreased below a critical level.

Condition 24. Regarding the learning rate (or gain) $\{\eta_t\}$ should satisfy the following conditions in order for SGLD Algorithm 19 to generate an output $\{w^{(t)}\}$ approximately distributed according to (3.2)

- (1) $\eta_t \geq 0$,
- (2) $\sum_{t=1}^{\infty} \eta_t = \infty$
- (3) $\sum_{t=1}^{\infty} \eta_t^2 < \infty$

Note 25. The popular learning rates $\{\eta_t\}$ in Note 12 in Lect. Notes 4 can be used. . Once parametrized, η_t can be tuned based on pilot runs using a reasonably small number of data.

¹A continuous-time stochastic process: (1) $B(0) = 0$; (2) $B(t)$ is almost surely continuous ; (3) $B(t)$ has independent increments ; (4) $B(t) - B(s) \sim N(0, t-s)$ for $0 \leq s \leq t$.

About the temperature alg. parameter...

Note 26. The temperature parameter $\tau > 0$ is user specified and aims at controlling (eg; inflating) the variance of the produced chain for instance with practical purpose to escape from local modes (otherwise energy barriers) in non-convex problems.

About the output of the alg...

Note 27. The first few iterations of SGLD (Algorithm 19) involve values generated at the beginning of the running algorithm while the chain have not yet converged to (or reached) an area of substantial posterior mass. Hence they are discarded from the output of the SGLD. These values are called burn-in.

Note 28. The output of SGLD (Algorithm 19) $\{w^{(t)}\}$ includes the generated values of w produced during the last few iterations of the running algorithm and after discarding the burn-in generated values.

Note 29. SGLD for $\tau = 1$ approximately simulates from the posterior (1.2).

About facilitating inference...

Note 30. Expectation (1.3), can be estimated as an arithmetic average

$$(3.7) \quad \widehat{h_T(w)} = \frac{1}{T} \sum_{t=1}^T h(w^{(t)})$$

as $\widehat{h_T(w)} \rightarrow E_f(h(w) | z_{1:n})$ based on LLN arguments.

Note 31. Another more efficient estimator for the expectation (1.3) is the weighted arithmetic average

$$(3.8) \quad \widehat{\tilde{h}_T(w)} = \sum_{t=T_0+1}^T \frac{\eta_t}{\sum_{t=T_0+1}^T \eta_{t'}} h(w^{(t)})$$

Note 32. Estimator (3.8) has a smaller standard error than estimator (3.7). As the step size η_t decreases, the mixing rate of the chain $\{w^{(t)}\}$ decreases as well and the simple sample average (3.7) overemphasizes in the tail end of the sequence where there is higher correlation among the samples resulting in higher variance in the estimator.

About the exploding gradient phenomenon...

Note 33. ‘Exploding gradients’ is the practical phenomenon in which large updates to weights during training can cause a numerical overflow or underflow due to the machine error of the computer.

Note 34. Practical solutions to ‘Exploding gradient’ involve, at each iteration t , checking the magnitude of the gradient v_t , (e.g., Euclidean norm $\|v_t\|$), and instantly changing it (e.g., gradient scaling or clipping) if it is gonna result an overflow.

Note 35. Gradient scaling involves normalizing the gradient vector such that vector norm (magnitude) equals a user-specified value. More formally, given a gradient v_t , gradient clipping can be used in a standard recursion

$$w^{(t+1)} = w^{(t)} + \eta_t v_t$$

as

$$w^{(t+1)} = w^{(t)} + \eta_t \text{clip}(v_t, c)$$

where

$$\text{clip}(v, c) = v \min\left(1, \frac{c}{\|v\|}\right)$$

and $c > 0$ is a clipping threshold implying that clipping will take place at iteration t if $\|v_t\| > c$.

Note 36. Gradient clipping involves forcing the gradient values (element-wise) to a specific minimum or maximum value if the gradient exceeded an expected range.

Note 37. Unreasonably frequent scaling or clipping may introduce significant bias and hence it should not be applied carelessly.

4. EXAMPLES ²

We continue the Example 26 in 4: Gradient descent.

Specifying the Bayesian model. Consider the Bayesian Normal linear regression model

$$\begin{cases} y_i|\beta, \sigma^2 \sim N(x_i^\top \beta, \sigma^2) & \text{sampling distribution } f(y_i|\beta, \sigma^2) \\ \beta|\sigma^2 \sim N(\mu, \sigma^2 V) & \text{prior } f(\beta|\sigma^2) \\ \sigma^2 \sim IG(\phi, \psi) & \text{prior } f(\sigma^2) \end{cases}$$

and $f(\beta, \sigma^2) = f(\beta|\sigma^2) f(\sigma^2)$, $\beta \in \mathbb{R}^d$, and $\sigma^2 \in \mathbb{R}_+$. Note given densities

$$\begin{aligned} N(x|\mu, \Sigma) &= \left(\frac{1}{2\pi}\right)^d \frac{1}{|\Sigma|} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1} (x - \mu)\right) \\ IG(x|a, b) &= \frac{b^a}{\Gamma(a)} x^{-a-1} \exp\left(-\frac{b}{x}\right) \mathbf{1}(x \geq 0) \end{aligned}$$

Performing computationally convenient transformations. Because SGD (Algorithm 13) and SGLD (Algorithm 19) can handle cases where $w \in \mathbb{R}^d$ in a straightforward manner than what they do when $w = (\beta, \sigma^2) \in \mathbb{R}^d \times \mathbb{R}_+$ which requires an additional projection step; we consider a transformation $w = (\beta, \gamma) \in \mathbb{R}^d \times \mathbb{R}$ with $\gamma = \log(\sigma^2)$. Hence, the Bayesian model becomes

$$\begin{cases} y_i|\beta, \sigma^2 \sim N(x_i^\top \beta, \exp(\gamma)) & \text{sampling distribution } f(y_i|\beta, \gamma) \\ \beta|\sigma^2 \sim N(\mu = 0, \exp(\gamma) V) & \text{prior } f(\beta|\gamma) \\ \gamma \sim f_\gamma(\gamma) & \text{prior } f(\gamma) \end{cases}$$

²Code is available from https://github.com/georgios-stats/Machine_Learning_and_Neural_Networks_III_Epiphanys_2025/tree/main/Lecture_notes/code/06.Stochastic_gradient_Langevine_dynamics

where $f_\gamma(\gamma)$ is computed according to the method of bijective transformation of random variables; i.e.

$$f_\gamma(\gamma) = \text{IG}(\exp(\gamma) | \phi, \psi) \left| \frac{d}{d\gamma} \exp(\gamma) \right| = \text{IG}(\exp(\gamma) | \phi, \psi) \exp(\gamma)$$

Computing the require quantities for SGLD. Then we can compute the required gradients in order to run the SGD, and SGLD with respect to $w = (\beta, \gamma)$ with $\gamma = \log(\sigma^2)$.

$$\begin{aligned} \log(f(z_i = (x_i, y_i) | w)) &= -\frac{1}{2} \log(2\pi) - \frac{1}{2}\gamma - \frac{1}{2} (y_j - x_i^\top \beta)^2 \exp(-\gamma) \\ \log(f(w = (\beta, \gamma))) &= \log(f(\beta | \gamma)) + \log(f(\gamma)) \\ \log(f(\beta | \gamma)) &= -\frac{d}{2} \log(2\pi) - \frac{d}{2}\gamma - \frac{1}{2} |V| - \frac{1}{2} \exp(-\gamma) (\beta - \mu)^\top V^{-1} (\beta - \mu) \\ \log(f(\gamma)) &= \psi \log(\phi) - \log(\Gamma(\phi)) - (\phi + 1)\gamma - \psi \exp(-\gamma) + \gamma \end{aligned}$$

Hence for the log sampling PDF we have

$$\begin{aligned} \nabla_w \log(f(z_i | w)) &= \left(\frac{d}{d\beta} \log(f(z_i | w)), \frac{d}{d\gamma} \log(f(z_i | w)) \right) \\ \frac{d}{d\beta} \log(f(z_i | w)) &= (y_i - x_i^\top \beta) x_i \exp(-\gamma) \\ \frac{d}{d\gamma} \log(f(z_i | w)) &= -\frac{1}{2} + \frac{1}{2} (y_j - x_i^\top \beta)^2 \exp(-\gamma) \end{aligned}$$

...so

$$(4.1) \quad \nabla_w \log(f(z_i | w)) = \begin{pmatrix} (y_i - x_i^\top \beta) x_i \exp(-\gamma) \\ -\frac{1}{2} + \frac{1}{2} (y_j - x_i^\top \beta)^2 \exp(-\gamma) \end{pmatrix}$$

Hence for the log a priori PDF we have

$$\begin{aligned} \nabla_w \log(f(w)) &= \left(\frac{d}{d\beta} \log(f(w)), \frac{d}{d\gamma} \log(f(w)) \right) \\ \frac{d}{d\beta} \log(f(w)) &= -\exp(-\gamma) V^{-1} (\beta - \mu) \\ \frac{d}{d\gamma} \log(f(w)) &= -\frac{d}{2} + \frac{1}{2} \exp(-\gamma) (\beta - \mu)^\top V^{-1} (\beta - \mu) - (\phi + 1) + \psi \exp(-\gamma) + 1 \end{aligned}$$

...so

$$(4.2) \quad \nabla_w \log(f(w)) = \begin{pmatrix} -\exp(-\gamma) V^{-1} (\beta - \mu) \\ \frac{d}{2} + \frac{1}{2} \exp(-\gamma) (\beta - \mu)^\top V^{-1} (\beta - \mu) - (\phi + 1) + \psi \exp(-\gamma) + 1 \end{pmatrix}$$

Implementation of SGLD. We consider $\mu = 0$, $\phi = 1$, $\psi = 1$, and $V = 100I_d$, for our simulations below.

To implement SGD (Algorithm 13) and SGLD (Algorithm 19), we just need to plug in the computed gradients (4.1) and (4.2) for $w = (\beta, \gamma) \in \mathbb{R}^d \times \mathbb{R}$ with $\gamma = \log(\sigma^2)$. After running SGD and SGLD with the computed gradients with respect to $w = (\beta, \gamma) \in \mathbb{R}^d \times \mathbb{R}$ with $\gamma = \log(\sigma^2)$, and obtaining chains $\{w^{(t)} = (\beta^{(t)}, \gamma^{(t)})\}_{t=1}^T$, we can just perform transformation $\{(\sigma^{(t)})^2 = \exp(\gamma^{(t)})\}_{t=1}^T$ if we are interested in learning (β, σ^2) .

In Figures 4.1, we ran the SGD for different batch sizes m and compared it against the exact MLE. We observe that SGD trace converges to the exact MLE. The oscillations are due to the stochastic gradient (ie, noise in the gradient).

In Figures 4.2, we ran the SGLD for different batch sizes m but the same temperature $\tau = 1$ and compared it against the exact posterior densities. We observe that the histograms of $\{w^{(t)}\}$ produced from SGLD are closer to the curves representing the exact posteriors when the batch size m is bigger. As we said this is not a panacea; if the landscape of the exact psterior density was multimodal (aka not convex but with had several maxima), then the SGLD using smaller batch sizes could have performed better, in the sense that the inflated noise from the stochastic gradient could accidentally make the generated chain to pass the low mass barrier and visit a different mode, unlike the one with larger batch-size and hence smaller variation.

In Figures 4.3, we ran the SGLD for different temperatures τ but the same batch sizes $m = 100$ and compared it against the exact posterior densities. We observe that increasing the temperature τ may increase the variation of the produced chain. We can use a large τ at the beginning of the run of the algorithm to perform an exploration of the space (this is particularly useful for non-convex/multimodal densities as it allows visiting different modes), and later on we can use a smaller temperature such as $\tau = 1$.

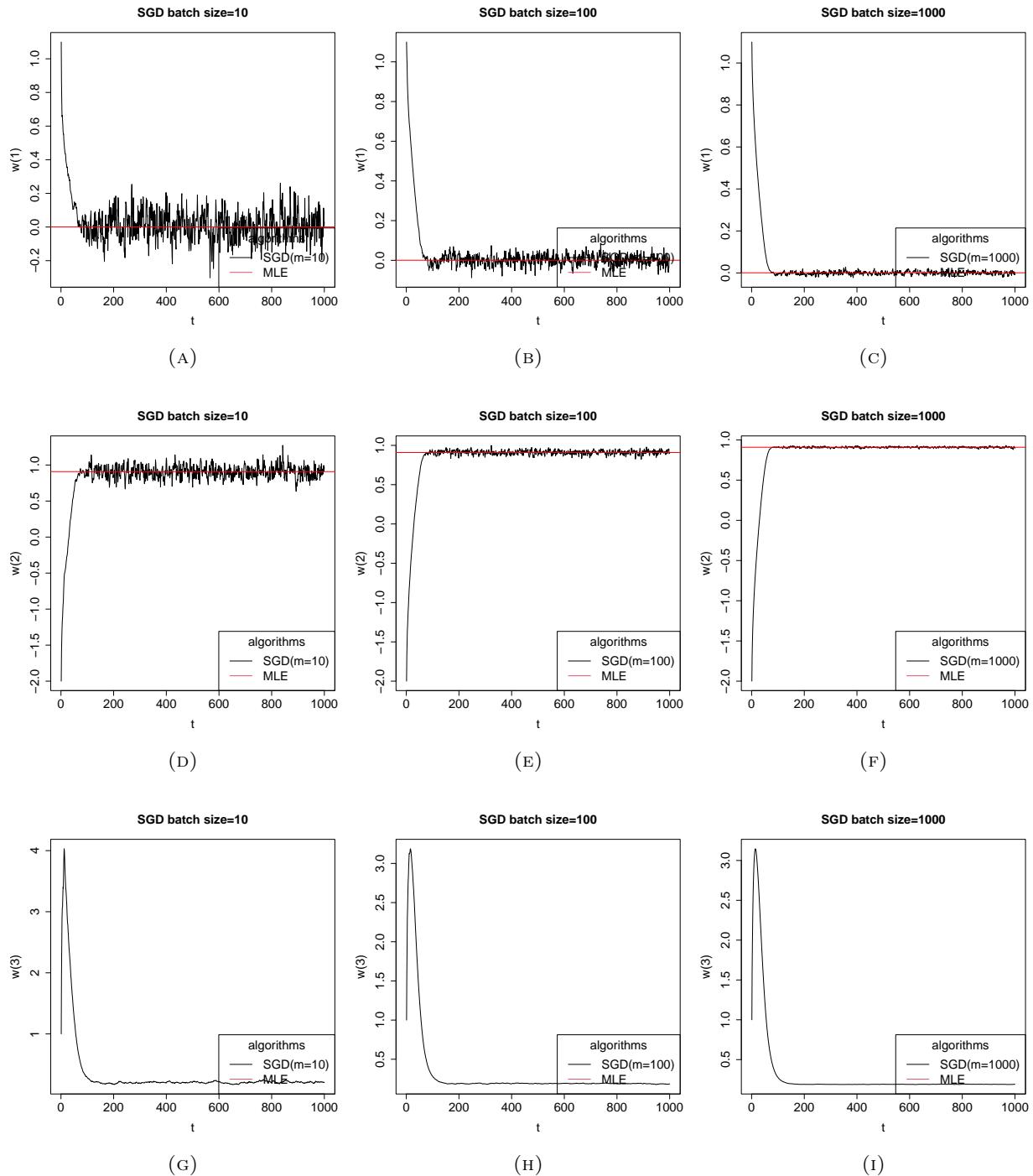


FIGURE 4.1. Bayesian learning via SGD (SGD vs (exact)MLE) Study on batch size m .

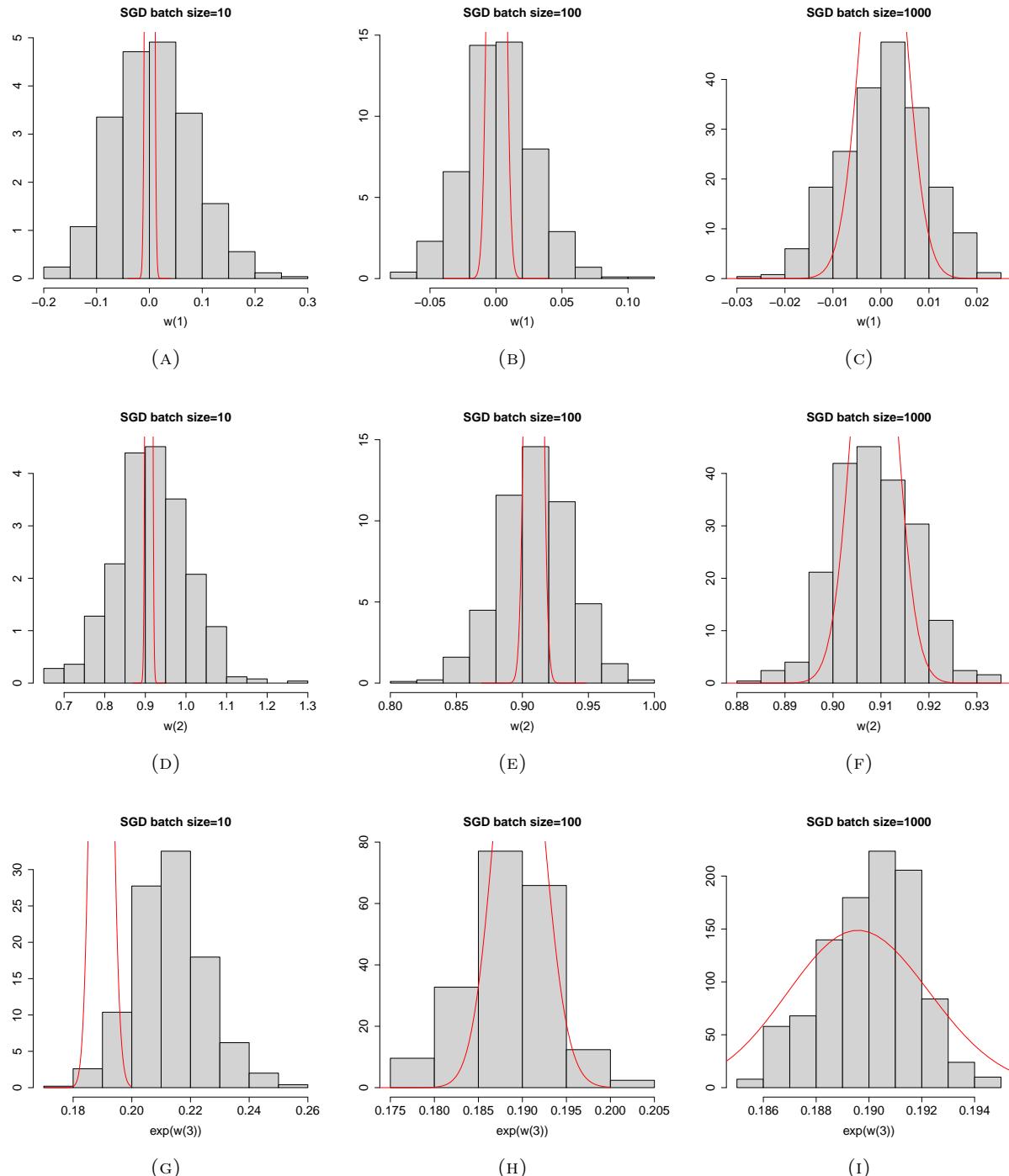


FIGURE 4.2. Bayesian learning: SGLD vs exact posterior (in red). Temperature $\tau = 1$. Study on batch size m

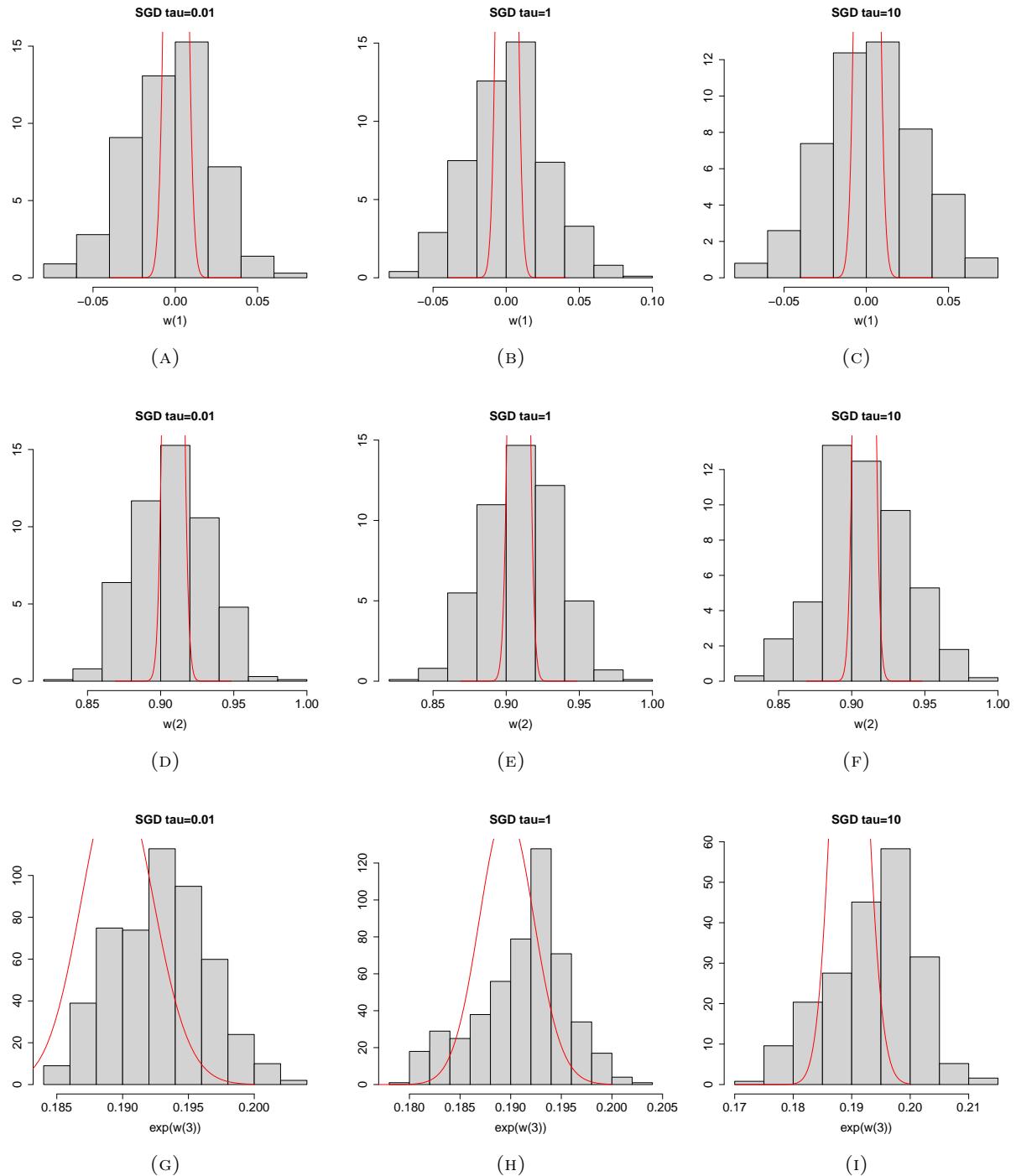


FIGURE 4.3. Bayesian learning: SGLD vs exact posterior (in red). Batch size = 100. Study on τ

Lecture notes 8: Support Vector Machines

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim. To introduce the Support Vector Machines as a procedure. Motivation, set-up, description, computation, and implementation. We focus on the classical treatment.

Reading list & references:

- (1) Shalev-Shwartz, S., & Ben-David, S. (2014). Understanding machine learning: From theory to algorithms. Cambridge university press.
 - Ch. 15 (pp. 167-170, 171-172, 176-177) Support Vector Machine
- (2) Bishop, C. M. (2006). Pattern recognition and machine learning (Vol. 4, No. 4, p. 738). New York: Springer.
 - Ch. 7.1 Sparse Kernel Machines/Maximum marginal classifiers
- (3) Vapnik, V. (2013). The nature of statistical learning theory. Springer science & business media.
- (4) Boyd, S. P., & Vandenberghe, L. (2004). Convex optimization. Cambridge university press.
 - Ch. 4, 5

1. INTRO AND MOTIVATION

Note 1. Support Vector Machines (SVM) is a ML procedure for learning linear predictors in high-dimensional feature spaces with regards the sample complexity challenges. Due to a duality property, SVM have sparse solutions, so that predictions for new inputs depend only on quantities evaluated at a subset of the whole training dataset.

Definition 2. Let $w \neq 0$. **Hyperplane** in space $\mathcal{X} \subseteq \mathbb{R}^d$ is called the sub-set

$$(1.1) \quad S = \left\{ x \in \mathbb{R}^d : \langle w, x \rangle + b = 0 \right\}.$$

Note 3. Hyperplane (1.1) separates \mathcal{X} in two half-spaces

$$S_+ = \left\{ x \in \mathbb{R}^d : \langle w, x \rangle + b > 0 \right\}$$

and

$$S_- = \left\{ x \in \mathbb{R}^d : \langle w, x \rangle + b < 0 \right\}$$

Definition 4. Halfspace (hypothesis space) is hypotheses class \mathcal{H} designed for binary classification problems, $\mathcal{X} \subseteq \mathbb{R}^d$ and $\mathcal{Y} = \{-1, +1\}$ defined as

$$\mathcal{H} = \left\{ x \mapsto \text{sign}(\langle w, x \rangle + b) : w \in \mathbb{R}^d, b \in \mathbb{R} \right\},$$

where b is called bias.

Definition 5. Each **halfspace hypothesis** $h \in \mathcal{H}$ has form

$$(1.2) \quad h_{w,b}(x) = \text{sign}(\langle w, x \rangle + b).$$

It takes an input in $\mathcal{X} \subseteq \mathbb{R}^d$ and returns an output in $\mathcal{Y} = \{-1, +1\}$. We may refer to it as halfspace (w, b) as well.

Note 6. Let $S = \{(x_i, y_i)\}_{i=1}^m$ be a training set of examples with $x_i \in \mathbb{R}^d$ the features and $y_i \in \{-1, +1\}$ the labels.

Note 7. Our goal is to train a halfspace hypothesis $h_{w,b}(x)$ in (1.2) against a training dataset S , with purpose to be able to classify a future feature x as $y = -1$ or $y = 1$.

Definition 8. The training set S is **linearly separable** if there exists a halfspace (w, b) such that for all $i = 1, \dots, n$

$$y_i = \text{sign}(\langle w, x_i \rangle + b)$$

or equivalently

$$(1.3) \quad y_i (\langle w, x_i \rangle + b) > 0$$

Note 9. Halfspaces (w^*, b^*) satisfying the linearly separable condition (1.3) are ERM hypothesis under the $0 - 1$ loss function $\ell^{0-1}((w, b), z) = 1_{(y_i \neq \text{sign}(\langle w, x_i \rangle + b))}$ and Empirical Risk $R_S^{0-1}(w, b) = \frac{1}{m} \sum_{i=1}^m \ell((w, b), z_i) \geq 0$, i.e.

$$(w^*, b^*) = \arg \min_{w,b} (R_S^{0-1}(w, b)) = \arg \min_{w,b} \left(\frac{1}{m} \sum_{i=1}^m \ell^{0-1}((w, b), z_i) \right)$$

as $R_S^{0-1}(w^*, b^*) = 0$.

Definition 10. Margin of a hyperplane with respect to a training set is defined to be the minimal distance between a point in the training set and the hyperplane.

Note 11. There are several different halfspaces (w, b) satisfying (1.3) for the same linearly separable training dataset S ; see Figure (1.1; left). In Figure (1.1; right) the margin γ is the distance from the hyperplane (solid line) to the closest points in either class (which touch the parallel dotted lines). Among the hyperplanes satisfying (1.3), a reasonable/desirable halfspace (w, b) is the one with the maximum margin γ ; see Figure (1.1; right). The rational is that if a hyperplane has a large margin, then it will still separate the training set even if we slightly perturb each instance.

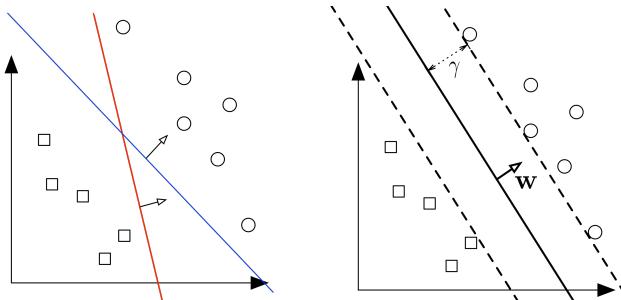


FIGURE 1.1

Note 12. Support Vector Machines (SVM) aims at learning the maximum margin separating hyperplane Figure (1.1; Right).

2. HARD SUPPORT VECTOR MACHINE

Assumption 13. Assume the training sample $S = \{(x_i, y_i)\}_{i=1}^m$ is linearly separable.

Definition 14. Hard Support Vector Machine (Hard-SVM) is the learning rule in which we return an ERM hyperplane that separates the training set with the largest possible margin given Assumption 13.

Problem 15. (Hard-SVM) Given a linearly separable training sample $S = \{(x_i, y_i)\}_{i=1}^m$ the Hard-SVM rule for the binary classification problem is the solution to the quadratic optimization problem:

$$(2.1) \quad \text{Solve } (\tilde{w}, \tilde{b}) = \arg \min_{(w,b)} \frac{1}{2} \|w\|_2^2$$

$$(2.2) \quad \text{subject to: } y_i(\langle w, x_i \rangle + b) \geq 1, \quad \forall i = 1, \dots, m$$

Scale

$$\hat{w} = \frac{\tilde{w}}{\|\tilde{w}\|}, \text{ and } \hat{b} = \frac{\tilde{b}}{\|\tilde{w}\|}$$

Note 16. Following we show why Problem 15 produces a Hard-SVM hyperplane stated in Note 14.

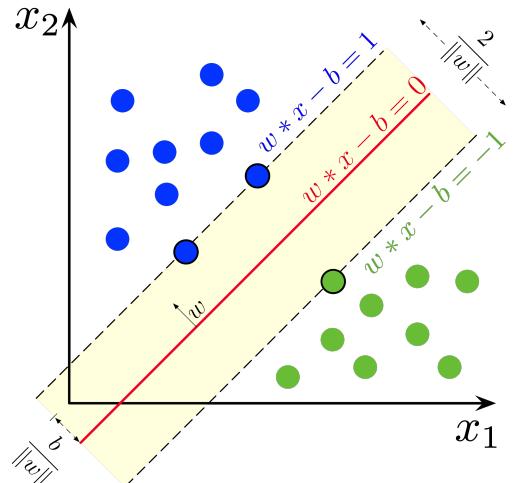
Lemma 17. The distance between a point x and the hyperplane defined by (w, b) with is $|\langle w, x \rangle + b| / \|w\|$.

Proof. We skip it. \square

Note 18. On the right, see the geometry of Problem 15.

Note 19. Hard-SVM selects two parallel hyperplanes that separate the two classes of data so that the distance between them is as large as possible. The predictive hyperplane (rule) is the hyperplane that lies halfway between them.

Note 20. Hard-SVM in Problem 15 searches for the hyperplane with minimum norm w among all those that separate the data and have distance greater or equal to 1.



Proof. (Sketch of the proof of Problem 15)

- (1) Let's formulate the Hard SVM problem according to the SVM definition under the linear separation assumption.

- (a) Based on Note 14, and Lemma 17, the closest point in the training set to the separating hyperplane has distance

$$(2.3) \quad \min_i (|\langle w, x_i \rangle + b| / \|w\|)$$

According to Definition 14, the Hard-SVM hypothesis can be formulated as having to find

$$(2.4) \quad (w^*, b^*) = \arg \max_{(w,b)} \left(\min_i \left(\frac{1}{\|w\|} |\langle w, x_i \rangle + b| \right) \right)$$

$$(2.5) \quad \text{subject to } y_i (\langle w, x_i \rangle + b) > 0, \forall i = 1, \dots, m$$

- (b) We observe that (2.4) & (2.5) is non-identifiable in the sense that $(\kappa w^*, \kappa b^*)$ for any $\kappa \neq 0$ are solutions. To make (2.3) identifiable (without loss of generality) we pick only (w^*, b^*) with $\|w\| = 1$. Additionally, because we consider only (w^*, b^*) which correctly separate the linearly separable training dataset into the two categories (i.e. (2.5)), it is $|\langle w, x_i \rangle + b| = y_i (\langle w, x_i \rangle + b)$. Hence, (2.4) & (2.5) is equivalent to computing

$$(2.6) \quad (w^*, b^*) = \arg \max_{(w,b): \|w\|=1} \left(\min_i (y_i (\langle w, x_i \rangle + b)) \right)$$

- (2) Next we show that the solution of (2.6) is equivalent to the solution of Problem 15; i.e. $(w^*, b^*) = (\hat{w}, \hat{b})$.

Let $\gamma^* := \min_i (|\langle w^*, x_i \rangle + b^*|)$ be the margin from (w^*, b^*) . Firstly, because

$$y_i (\langle w^*, x_i \rangle + b^*) \geq \gamma^* \Leftrightarrow y_i \left(\langle \frac{w^*}{\gamma^*}, x_i \rangle + \frac{b^*}{\gamma^*} \right) \geq 1$$

$\left(\frac{w^*}{\gamma^*}, \frac{b^*}{\gamma^*} \right)$ satisfies conditions (2.1) and (2.2). Secondly, I have $\|\tilde{w}\| \leq \left\| \frac{w^*}{\gamma^*} \right\| = \frac{1}{\gamma^*}$ because of (2.1) and because of $\|w^*\| = 1$. Hence, for all $i = 1, \dots, m$, it is

$$y_i \left(\langle \hat{w}, x_i \rangle + \hat{b} \right) = \frac{1}{\|\tilde{w}\|} y_i \left(\langle \tilde{w}, x_i \rangle + \tilde{b} \right) \geq \frac{1}{\|\tilde{w}\|} \geq \gamma^*$$

Hence (\hat{w}, \hat{b}) is the optimal solution of (2.6).

□

Definition 21. Homogeneous halfspaces in SVM is the case where the halfspaces pass from the origin; that is when the bias term in (2.2) is zero $b = 0$.

3. SOFT SUPPORT VECTOR MACHINE

Note 22. Hard-SVM assumes the strong Assumption 13 that the training set is linearly separable. This might not always be the case, and hence there is need to derive a procedure that weakens this assumption.

Note 23. Soft Support Vector Machine (Soft-SVM) aims to relax the strong assumption of Hard-SVM that the training set is linearly separable (2.5) with purpose to extend the scope of application.

Soft-SVM does not assume Assumption 13. Soft-SVM rule is given as the solution to the quadratic optimization problem 24.

Problem 24. (Soft-SVM) Given a training sample $S = \{(x_i, y_i)\}_{i=1}^m$ the Soft-SVM rule for the binary classification learning problem is solution to the quadratic optimization Problem

$$(3.1) \quad \text{Solve } (w^*, b^*, \xi^*) = \arg \min_{(w, b, \xi)} \left(\frac{1}{2} \|w\|_2^2 + C \frac{1}{m} \sum_{i=1}^m \xi_i \right)$$

$$(3.2) \quad \text{subject to: } y_i (\langle w, x_i \rangle + b) \geq 1 - \xi_i, \quad \forall i = 1, \dots, m$$

$$(3.3) \quad \xi_i \geq 0, \quad \forall i = 1, \dots, m$$

Note 25. To relax the linearly separable training set Assumption 13, Soft-SVM relies on replacing the “harder” constraint (2.2) with the “softer” one in (3.2) through the introduction of non-negative unknown quantities (slack variables) $\{\xi_i\}_{i=1}^m$ controlling how much the separability assumption (2.2) is violated. Because any point that is misclassified has $\xi_i > 1$, it follows that $\sum_{i=1}^m \xi_i$ is an upper bound on the number of misclassified points.

Note 26. Parameter $C > 0$ controls the trade-off between the slack variable penalty and the margin. It controls the trade-off between minimizing training errors and controlling model complexity.

Note 27. Soft-SVM learns all (w, b, ξ) via the minimization part in (3.1) where the trade off between the two terms is controlled via the user specified parameter C .

Note 28. Proposition 29 shows that the Soft-SVM is a binary classification learning problem under the hinge loss function and with a regularization term biasing toward low norm separators.

Proposition 29. *The solution of Problem 24 is equivalent to the Ridge regularized loss minimization problem*

$$(3.4) \quad (w^*, b^*) = \arg \min_{(w, b)} \left(R_S^{hinge} + \lambda \|w\|_2^2 \right)$$

with the hinge loss function

$$\ell^{hinge}((w, b), z) = \max(0, 1 - y(\langle w, x \rangle + b))$$

Empirical Risk Function

$$R_S^{hinge}((w, b)) = \frac{1}{m} \sum_{i=1}^m \max(0, 1 - y_i (\langle w, x_i \rangle + b))$$

regularization parameter $\lambda = \frac{1}{2C}$ and regularization term $\|\cdot\|_2^2$.

Proof. In Algorithm 24, we consider (3.1) as

$$(3.5) \quad \arg \min_{(w, b)} \left(\min_{\xi} \left(\lambda \|w\|_2^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \right) \right)$$

with $\lambda = \frac{1}{2C}$. Consider (w, b) fixed and focus on the inside minimization. From (3.2), it is $\xi_i \geq 1 - y_i (\langle w^*, x_i \rangle + b^*)$, and from (3.3), it is $\xi_i \geq 0$. If $y_i (\langle w, x_i \rangle + b) \geq 1$, the best assignment in 3.5 is $\xi_i = 0$ because it is $\xi_i \geq 0$ from (3.3) and I need to minimize (3.5) wrt ξ_i 's. If $y_i (\langle w, x_i \rangle + b) \leq 1$, the best assignment in (3.5) is $\xi_i = 1 - y_i (\langle w, x_i \rangle + b)$ because I need to minimize w.r.t ξ_i 's. Hence $\xi_i = \max(0, 1 - y_i (\langle w, x_i \rangle + b))$. \square

Example 30. Consider the Soft-SVM as a Ridge regularized loss minimization learning problem in (3.4) on a training dataset $S = \{(x_i, y_i)\}_{i=1}^m$ with $z_i = (x_i, y_i) \stackrel{\text{ind}}{\sim} g$ where g is a data generating process on $\mathcal{X} \times \{0, 1\}$ with $\mathcal{X} = \{x : \|x\| \leq \rho\}$. Let $\mathfrak{A}(S)$ be the solution of the Soft-SVM. Then:

- Because $\ell^{\text{hinge}}((w, b), z) = \max(0, 1 - y(\langle w, x \rangle + b))$ is $\|x\|$ -Lipschitz, it is

$$\mathbb{E}_{S \sim g}(R_g^{\text{hinge}}(\mathfrak{A}(S))) \leq R_g^{\text{hinge}}(u) + \lambda \|u\|^2 + \frac{2\rho^2}{\lambda m}$$

[Hint: Note 43 Handout 3: Learnability and stability in learning problems]

- If we assume a bounded learning problem, i.e. $\mathcal{H} = \{(w, b) : \|w, b\| \leq B\}$ for $B > 0$ then

$$\mathbb{E}_{S \sim g}(R_g^{\text{hinge}}(\mathfrak{A}(S))) \leq \min_{\|(w, b)\| \leq B} R_g^{\text{hinge}}((w, b)) + \lambda B^2 + \frac{2\rho^2}{\lambda m}$$

- If we choose $\lambda = \sqrt{\frac{2\rho^2}{B^2 m}}$ then

$$\mathbb{E}_{S \sim g}(R_g^{\text{hinge}}(\mathfrak{A}(S))) \leq \min_{\|(w, b)\| \leq B} R_g^{\text{hinge}}((w, b)) + \sqrt{\frac{8\rho^2 B^2}{m}}$$

- Assume one was interested in minimizing a $0 - 1$ Risk $R_g^{0-1}(\cdot)$ under the $0 - 1$ loss $\ell^{0-1}((w, b), z) := 1_{(y(w, x) \leq 0)}$. But $\ell^{0-1}(\cdot, z)$ is non-convex for all z . It is $\ell^{0-1}((w, b), z) \leq \ell^{\text{hinge}}((w, b), z)$ for all z . By using the convex hinge loss ℓ^{hinge} as a surrogate loss and implementing the Soft-SVM procedure, the error under $0 - 1$ Risk $R_g^{0-1}(\cdot)$ is

$$\mathbb{E}_{S \sim g}(R_g^{0-1}(\mathfrak{A}(S))) \leq \mathbb{E}_{S \sim g}(R_g^{\text{hinge}}(\mathfrak{A}(S))) \leq \min_{\|(w, b)\| \leq B} R_g^{\text{hinge}}((w, b)) + \sqrt{\frac{8\rho^2 B^2}{m}}$$

Example 31. Given Proposition 29, Soft-SVM in Problem 24 can be learned via any variation of SGD, eg online SGD (batch size $m = 1$) with recursion

$$\begin{aligned} \varpi^{(t+1)} &= \varpi^{(t)} - \eta_t v_t \\ \begin{bmatrix} w^{(t+1)} \\ b^{(t+1)} \end{bmatrix} &= \begin{bmatrix} w^{(t)} \\ b^{(t)} \end{bmatrix} - \eta_t \begin{bmatrix} v_{w,t} \\ v_{b,t} \end{bmatrix} \end{aligned}$$

where $v_{w,t} = \begin{cases} -y_i^{(t)} x_i^{(t)} & \text{if } y_i^{(t)} (\langle w^{(t)}, x_i^{(t)} \rangle + b^{(t)}) < 1 \\ 0 & \text{otherwise} \end{cases}$ and $v_{b,t} = \begin{cases} -y_i^{(t)} & \text{if } y_i^{(t)} (\langle w^{(t)}, x_i^{(t)} \rangle + b^{(t)}) < 1 \\ 0 & \text{otherwise} \end{cases}$.

4. DUALITY AND SPARSITY

Recall Lagrange multipliers and duality from previous years

4.1. Lagrangian duality (recall from previous years).

Notation 32. Let $f : \mathbb{R}^q \rightarrow \mathbb{R}$ be an objective function, let $\{g_i : \mathbb{R}^q \rightarrow \mathbb{R}\}_{i=1}^m$ inequality constraint convex functions, and let $\{h_j : \mathbb{R}^q \rightarrow \mathbb{R}\}_{j=1}^n$ equality constraint functions.

Definition 33. (Primal problem) Consider we have the following minimization problem that we will call Primal problem

$$(4.1) \quad \begin{aligned} p^* &= \min_x (f(x)) \\ \text{s.t. } &g_i(x) \leq 0, \quad \forall i = 1, \dots, m \\ &h_j(x) = 0, \quad \forall j = 1, \dots, n \end{aligned}$$

Definition 34. (Lagrangian function) To the Primal problem 33, we associate the Lagrangian function $L : \mathbb{R}^q \times \mathbb{R}^m \times \mathbb{R}^n$ with

$$(4.2) \quad L(x, \alpha, \beta) = f(x) + \sum_{i=1}^m \alpha_i g_i(x) + \sum_{j=1}^n \beta_j h_j(x)$$

Note 35. Note that the Primal problem is equivalent to

$$p^* = \min_x \left(\max_{\alpha \geq 0, \beta} (L(x, \alpha, \beta)) \right)$$

For instance, if I ignore the $\{h_j\}$ terms which is zero for a suitable solution for simplicity, I observe that

$$\max_{\alpha \geq 0, \beta} (L(x, \alpha, \beta)) = \begin{cases} f(x) & g_i(x) \leq 0, \forall i \\ \infty, & g_i(x) > 0, \exists i \end{cases}$$

Definition 36. (Lagrangian dual problem) The associated Lagrangian dual problem is

$$\begin{aligned} d^* &= \max_{\alpha, \beta} \left(\min_x (L(x, \alpha, \beta)) \right) \\ \text{s.t. } &\alpha_i \geq 0 \\ &\beta_j \in \mathbb{R} \end{aligned}$$

Definition 37. (Dual function) To the Dual problem, we associate a function $\tilde{L} : \mathbb{R}^m \times \mathbb{R}^n$ with

$$\tilde{L}(\alpha, \beta) := \min_x (L(x, \alpha, \beta))$$

called dual function.

Proposition 38. In general it is $p^* \geq d^*$; i.e.

$$\min_x \max_{\alpha \geq 0, \beta} (L(x, \alpha, \beta)) \geq \max_{\alpha \geq 0, \beta} \min_x (L(x, \alpha, \beta))$$

Definition 39. We call the general case $p^* \geq d^*$ weak duality.

Definition 40. When $p^* = d^*$ we say we have a strong duality.

Proposition 41. (*Strong duality via Slater condition*) If the primal problem (4.1) is convex, and satisfies the weak Slater's condition, i.e.

$$(\exists x_0 \in \mathcal{D}) : (g_i(x_0) < 0, \forall i = 1, \dots, n) \text{ and } (h_i(x_0) = 0, \forall i = 1, \dots, n)$$

then strong duality holds, that is: $p^* = d^*$. In other words

$$\min_x \max_{\alpha \geq 0, \beta} (L(x, \alpha, \beta)) = \max_{\alpha \geq 0, \beta} \min_x (L(x, \alpha, \beta))$$

Proposition 42. When $f, \{h_j\}, \{g_i\}$ are convex and the Slater condition holds, Karush–Kuhn–Tucker (KKT) conditions are necessary and sufficient conditions for x^* to be local minimum

$$(4.3) \quad \begin{aligned} 0 &= \nabla f(x^*) + \sum_{j=1}^n \beta_j \nabla h_j(x^*) + \sum_{i=1}^m \alpha_i \nabla g_i(x^*) && \text{Stationarity} \\ g_i(x^*) &\leq 0, \quad \forall i = 1, \dots, m && \text{Primal feasibility} \\ h_j(x^*) &= 0, \quad \forall j = 1, \dots, n \\ \alpha_i &\geq 0 \quad \forall i = 1, \dots, m && \text{Dual feasibility} \\ \alpha_i g_i(x^*) &= 0, \quad \forall i = 1, \dots, m && \text{Complementary slackness} \end{aligned}$$

4.2. Implementation in the Hard SVM.

Note 43. Consider the Primal problem in Hard-SVM Algorithm 15 (2.2) and (2.2).

Note 44. According to Definition 34 the Lagrangian function is

$$(4.4) \quad L(w, b, \alpha) = \frac{1}{2} \|w\|_2^2 - \sum_{i=1}^m \alpha_i (y_i (\langle w, x_i \rangle + b) - 1)$$

where $\{\alpha_i\}_{i=1}^m$ are non-negative Lagrange multipliers.

Insight with regards to Note 35

Note 45. (Insight with regards to Note 35) The Hard-SVM Problem 15 (2.2) and (2.2), can be rewritten in the form

$$(4.5) \quad \min_w \left(\frac{1}{2} \|w\|_2^2 + g(w) \right)$$

with

$$g(w) = \max_{\alpha \geq 0} \sum_{i=1}^m \alpha_i (1 - y_i (\langle w, x_i \rangle + b)) = \begin{cases} 0 & \text{if } y_i (\langle w, x_i \rangle + b) \geq 1 \\ \infty & \text{else} \end{cases}$$

Hence we get the weak duality

$$\begin{aligned}
 & \min_{w,b} \max_{\alpha \geq 0} \left(\frac{1}{2} \|w\|_2^2 + \sum_{i=1}^m \alpha_i (1 - y_i (\langle w, x_i \rangle + b)) \right) \\
 (4.6) \quad & \geq \max_{\alpha \geq 0} \min_{w,b} \underbrace{\left(\frac{1}{2} \|w\|_2^2 + \sum_{i=1}^m \alpha_i (1 - y_i (\langle w, x_i \rangle + b)) \right)}_{=L(w,b;\alpha)}
 \end{aligned}$$

Note 46. Strong duality (equality) holds because of the convexity. Here, keeping α fixed, $L(w, \alpha, b)$ is minimized when

$$(4.7) \quad 0 = \nabla_w L(w, \alpha, b)|_{(w^*, b^*)} \implies w^* = \sum_{i=1}^m \alpha_i y_i x_i$$

$$(4.8) \quad 0 = \nabla_b L(w, \alpha, b)|_{(w^*, b^*)} \implies 0 = \sum_{i=1}^m \alpha_i y_i$$

The dual Lagrangian function is

$$\begin{aligned}
 \tilde{L}(\alpha) &= \min_{w,b} (L(w, \alpha, b)) = \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i x_i \right\|_2^2 - \sum_{i=1}^m \alpha_i \left(y_i \left(\sum_{j=1}^m \alpha_j y_j x_j, x_i \right) + b \right) - 1 \\
 &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle x_j x_i \rangle
 \end{aligned}$$

Then dual problem is

$$\max_{\alpha \in \mathbb{R}^m : \alpha \geq 0} \left(\sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle x_j x_i \rangle \right)$$

Note 47. The Dual problem of the Hard-SVM (primal) problem 15 is

$$\begin{aligned}
 (4.9) \quad \alpha^* &= \arg \max_{\alpha \in \mathbb{R}^m : \alpha \geq 0} \left(\sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle x_j, x_i \rangle \right) \\
 \text{subject to } 0 &= \sum_{i=1}^m \alpha_i y_i
 \end{aligned}$$

Note 48. Once the optimal α^* is computed from (4.9), the optimal weights w^* can be computed (from (4.7)) as

$$(4.10) \quad w^* = \sum_{i=1}^m \alpha_i^* y_i x_i$$

Note 49. If $\alpha_i^* = 0$, the example (x_i, y_i) does not contribute to (4.10). If $\alpha_i^* \neq 0$, (x_i, y_i) contributes to (4.10). Moreover if $\alpha_i^* \neq 0$ the KKT condition (4.3)

$$\alpha_i^* (y_i (\langle w^*, x_i \rangle + b^*) - 1) = 0$$

implies that $y_i (\langle w^*, x_i \rangle + b^*) = 1$ meaning that x_i is on the boundary of the margin. Features x_i associated to $\alpha_i^* \neq 0$ which contribute to the evaluation of the optimal weights w^* and hence the evaluation of the separating predictive rule $h_{w,b}$ are called **supporting vectors** (see Figure 4.1).

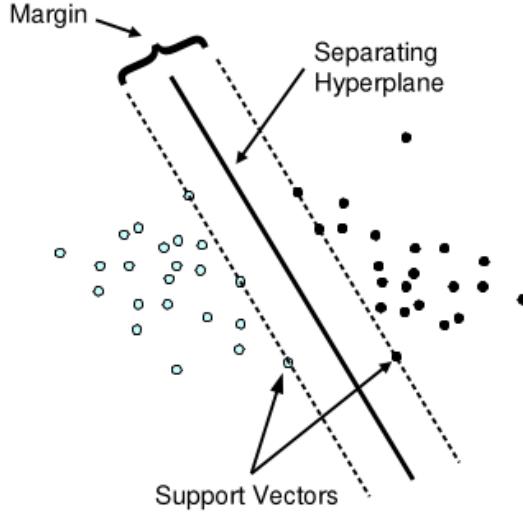


FIGURE 4.1

Note 50. Given optimal α^* (4.10) becomes

$$(4.11) \quad w^* = \sum_{i \in \mathcal{I}} \alpha_i^* y_i x_i$$

where $\mathcal{I} = \{i : y_i (\langle w^*, x_i \rangle + b^*) - 1 = 0\}$.

Note 51. The bias b parameter can be computed by KKT condition (4.3), it is

$$\alpha_i^* (y_i (\langle w^*, x_i \rangle + b^*) - 1) = 0$$

Consider $\mathcal{I} = \{i : y_i (\langle w^*, x_i \rangle + b^*) - 1 = 0\}$ then, for $i \in \mathcal{I}$ it is

$$y_i^2 (\langle w^*, x_i \rangle + b^*) - y_i = 0 \stackrel{y_i^2=1}{\implies} \langle w^*, x_i \rangle + b^* - y_i = 0$$

and summing up all $i \in \mathcal{I}$, I get

$$(4.12) \quad \sum_{i \in \mathcal{I}} (\langle w^*, x_i \rangle + b^* - y_i) = 0 \implies b^* = \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} (y_i - \langle w^*, x_i \rangle)$$

Note 52. Therefore the predictive halfspace hypothesis is given as

$$(4.13) \quad \begin{aligned} h_{w,b}(x) &= \text{sign}(\langle w^*, x \rangle + b^*) \\ &= \text{sign}\left(\sum_{i \in \mathcal{I}} \alpha_i^* y_i \langle x_i, x \rangle + b^*\right) \end{aligned}$$

with $\mathcal{I} = \{i : y_i (\langle w^*, x_i \rangle + b^*) - 1 = 0\}$.

Note 53. Compared to the Primal problem (Problem 15), the dual problem (*Note* 47) is computationally desirable in cases that the training data size m is smaller than the number of the dimensions d of the feature space; i.e. $d \gg m$. The solution of the Primal quadratic programming problem in $d + 1$ and complexity $O(d + 1)$, while the dual quadratic programming problem has m variables and complexity $O(m)$.

Note 54. Compared to the Primal problem, the dual problem (*Note* 47), can provide sparse solutions relying on a small number of support vectors (see Eq. 4.11, 4.12, and 4.13).

Note 55. (Snapshot of the next concept “The kernel trick”) The importance of the existence of the dual problem is that eg (4.9) and (4.13) involves the inner products between instances and does not require the direct access to specific elements of features within instance. For instance, if we consider the hypothesis (1.2) as an expansion of bases $h_{w,b}(x) = \text{sign}(\langle w, \phi(x) \rangle + b)$ with $\phi(x) = (\phi_1(x), \dots, \phi_d(x))$ with a large d such as $d \gg m$ then, based on (4.13), the separation rule is

$$h_{w,b}(x) = \text{sign} \left(\sum_{i=1}^m \alpha_i y_i k(x_i, x) + b \right)$$

and 4.9 becomes

$$\alpha^* = \arg \max_{\alpha \in \mathbb{R}^m: \alpha \geq 0} \left(\sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j k(x_i, x_j) \right)$$

with $k(x', x) = \langle \phi(x'), \phi(x) \rangle$. As we will see in the “kernel methods” lecture, in specific case, one can specify the “kernel” function $k(\cdot, \cdot)$ (having with specific desirable properties) avoiding the direct specification of a possibly high-dimensional dictionary of features $\{\phi_j(\cdot)\}$.

Note 56. In the Soft-SVM case, the solution is the same as in the Hard-SVM (4.10) and (4.13), the only difference is that in the quadratic programming problem (4.9) it is subject to $\alpha_j \in [0, C]$ where $C = 1/2\lambda$, it is called “cost” and controls the cost of having the constraint violation by adding the ξ'_i s. (See Exercise 16 from the Exercise sheet).

Note 57. Sparsity and duality in Soft SVM can be seen in Exercise 16 in the Exercise sheet.

Note 58. Relevance Vector Machine (a version of Bayesian SVM) can be seen in Exercise 17 in the Exercise sheet.

Lecture notes 9: The Kernel trick

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim. To introduce the ideas of learning machines by introducing data into high-dimensional feature spaces for accuracy gains; introduce the kernel trick, and kernel functions.

Reading list & references:

- (1) Shalev-Shwartz, S., & Ben-David, S. (2014). Understanding machine learning: From theory to algorithms. Cambridge university press.
 - Ch. 16.2 Support Vector Machine
- (2) Bishop, C. M. (2006). Pattern recognition and machine learning (Vol. 4, No. 4, p. 738). New York: Springer.
 - Ch. 6.1, 6.2 Kernel methods
- (3) Vapnik, V. (2013). The nature of statistical learning theory. Springer science & business media.

1. INTRO AND MOTIVATION

Note 1. Consider the Soft SVM with predictive rule $h(x) = \text{sign}(\eta(x))$ with separator $\eta(x) = \langle w, x \rangle + b$ where $w = (w_1, w_2)^\top \in \mathbb{R}^2$ and $x = (x_1, x_2)^\top \in \mathcal{X} \subseteq \mathbb{R}^2$. It can address learning problems where the data can (up to some degree of violation) be separated by a line (Figure 1.1a). In more challenging cases where the geometry is strongly non-linear this can totally fail (Figure 1.1b).

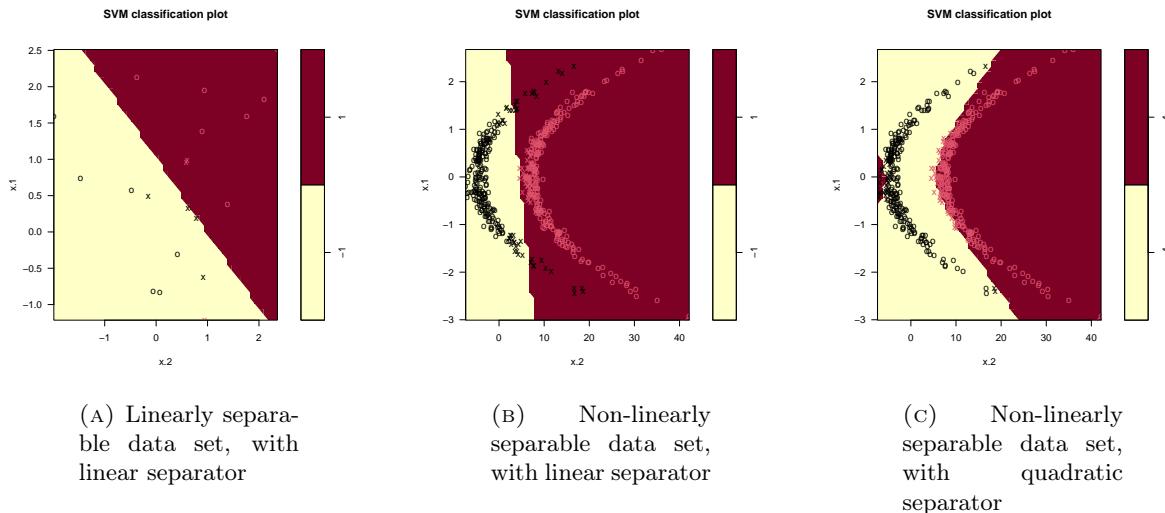


FIGURE 1.1. Soft SVM from the Computer lab

Note 2. The accuracy of predictive rule could be improved if I could take into account the curvature in the feature space \mathcal{X} by adding a quadratic term in the 2nd dimension as $h^\psi(x) = \text{sign}(\eta^\psi(x))$ with separator $\eta^\psi(x) = \langle w', \psi(x) \rangle + b'$ for some $\psi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$ with $\psi(x) = (x_1, x_2, x_2^2)$ and learning w' and b' . It works, (Figure 1.1c).

Note 3. In other words, in order to improve the expressiveness of a given hypothesis class $\mathcal{H} = \{x \mapsto f(\langle w, x \rangle)\}$ (for some function f) with purpose to learn a more accurate predictive rule, it is reasonable to consider an embedding $\psi(x)$ and work on the learning problem with $\mathcal{H} = \{x \mapsto f(\langle w, \psi(x) \rangle)\}$. Such an embedding $\psi(x)$ can possibly be a vector of basis functions such as polynomials, splines, etc...

Note 4. The above may drastically increase the dimensionality of the problem hence the computational cost and required training dataset size. This challenge is addressed by the Kernel trick which allows the design of cheap more expressive extensions of many well known algorithms.

2. IMPROVING EXPRESSIVE POWER VIA EMBEDDINGS IN FEATURE SPACES

Note 5. To make the class of hypotheses more expressive with purpose to improve accuracy, we can first map the original instance space $x \in \mathcal{X}$ into another feature space \mathcal{F} (possibly of a higher dimension) via an embedding $\psi : \mathcal{X} \rightarrow \mathcal{F}$ and then learn a hypothesis in that space.

Note 6. Consider a given ERM learning problem $(\mathcal{H}, \ell, \mathcal{Z})$ that involves hypothesis class \mathcal{H} , predictive rule $h \in \mathcal{H}$, loss function $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}_+$, training dataset $\mathcal{S} = \{z_i = (x_i, y_i)\}_{i=1}^m$ from data generating process $G(\cdot)$ defined over $\mathcal{Z} = \mathcal{Y} \times \mathcal{X}$.

Note 7. The introduction of embedding $\psi : \mathcal{X} \rightarrow \mathcal{F}$ induces

- (1) a probability distribution G^ψ over domain $\mathcal{F} \times \mathcal{Y}$ with $G^\psi(A) = G(\psi^{-1}(A))$ for every set $A \subseteq \mathcal{F} \times \mathcal{Y}$ given a data generating process $G(\cdot)$.
- (2) predictive rule $h^\psi(\cdot) := h \circ \psi(\cdot) = h(\psi(\cdot))$
- (3) risk function $R_{G^\psi}(h) := R_G(h \circ \psi)$, as

$$R_G(h \circ \psi) = \int \ell(h \circ \psi, z = (x, y)) dG(z) = \int \ell(h, z^\psi) dG^\psi(x, y) = R_{G^\psi}(h)$$

Note 8. The basic paradigm for improving expressive power of \mathcal{H} via embedding $\psi : \mathcal{X} \rightarrow \mathcal{F}$ involves:

- (1) Choose an embedding $\psi : \mathcal{X} \rightarrow \mathcal{F}$ with $\psi(x) := (\psi_1(x), \dots, \psi_d(x))^\top$ for some feature space \mathcal{F} .
- (2) Create the image sequence $\mathcal{S}^\psi = \left\{ z_i^\psi = (\psi(x_i), y_i) \right\}_{i=1}^m$ from the original training set \mathcal{S} .
- (3) Train a linear predictor h against \mathcal{S}^ψ .
- (4) Predict the label or the output of a new point x^{new} by $h^\psi(x^{\text{new}}) := h \circ \psi(x^{\text{new}}) = h(\psi(x^{\text{new}}))$

Note 9. For a specific learning task, the success of the learning paradigm in Note 6 depends on choosing an embedding ψ that provides a suitable deformation for the image of the data generating process (or the training dataset) to be as close as possible to what could be accurately addressed by the specific learning task.

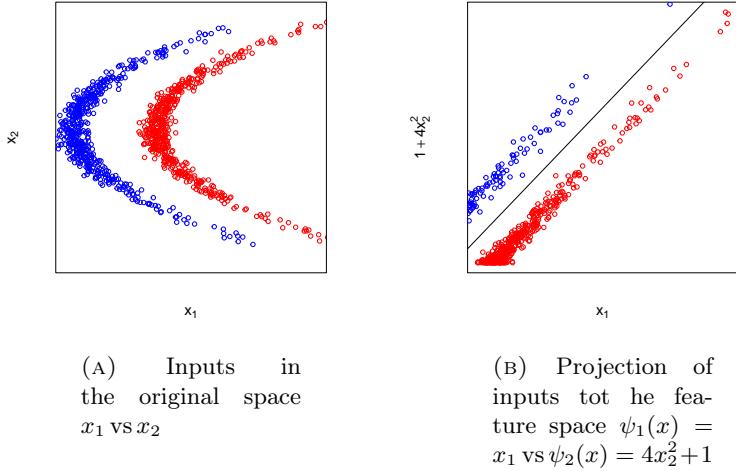


FIGURE 2.1. Projection of the inputs living in the original space to the feature space

Example 10. In SVM, ψ will make the image of the data distribution (close to being) linearly separable in the feature space \mathcal{F} , thus making the resulting learning algorithm a good learner for a given task (Figure 2.1). This requires prior knowledge of the problem (In Section 4, we see popular recipes for that).

Example 11. (Example) Consider the Soft-SVM in “Handout 8: Support Vector Machines”. Consider a given embedding $\psi : \mathcal{X} \rightarrow \mathcal{F}$ with $\psi(x) = (\psi_1(x), \psi_2(x) \dots)^\top$. The learning rule becomes

$$(2.1) \quad h^\psi(x) = h(\psi(x)) = \text{sign}(\langle w, \psi(x) \rangle + b).$$

In “Handout 8, Problem 24 becomes

Solve

$$\text{Solve: } (w^*, b^*, \xi^*) = \arg \min_{(w, b, \xi)} \left(\lambda \|w\|_2^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \right)$$

$$\text{subject to: } y_i (\langle w, \psi(x_i) \rangle + b) \geq 1 - \xi_i, \quad \forall i = 1, \dots, m$$

$$\xi_i \geq 0, \quad \forall i = 1, \dots, m$$

Note 12. Feature space \mathcal{F} is preferably a Hilbert space due to Lemma 15 that enables the Kernel trick via the representation Theorem 21. Eg, a Euclidean space such as \mathbb{R}^d for some d . That includes infinite dimensional spaces.

Note 13. Any function that maps the original instances \mathcal{X} into some Hilbert space \mathcal{F} can be used as feature mapping ψ .

Definition 14. A Hilbert space \mathcal{F} is a real or complex inner product (hence vector) space which is also complete^a metric space with respect to the distance function (norm) $\|\cdot\| : \mathcal{F} \rightarrow \mathbb{R}$ induced by the inner product $\langle \cdot, \cdot \rangle : \mathcal{F} \times \mathcal{F} \rightarrow \mathbb{C}$ as $\|\psi\| = \sqrt{\langle \psi, \psi \rangle}$ for all $\psi \in \mathcal{F}$.

Lemma 15. If \mathcal{X} is a linear subspace of a Hilbert space, then every $x \in \mathcal{X}$ can be written as $x = u + v$ where $u \in \mathcal{X}$ and $\langle u, v \rangle = 0$ for all $v \in \mathcal{X}$.

^aI.e., every Cauchy sequence converges to some element of \mathcal{F}

Note 16. Using a $\psi(x) = (\psi_1(x), \psi_2(x) \dots)^\top$ that is high dimensional (d is too large) may improve accuracy (expressiveness) of the learner (e.g. recall in polynomial regression increasing the polynomial degree). However this increases the computational effort/cost required to perform calculations to minimize the associated risk function in the high dimensional space, as well as we need more data. This is addressed via the Kernel trick.

3. THE KERNEL TRICK

Note 17. We discuss a duality in the learning Problem 18 that facilitates the implementation of the extension to a possibly high dimensional feature space (hence improving the expressiveness/accuracy) by using kernel functions (hence reducing dimensionality, computational cost, and required data size).

Problem 18. (RLM learning problem) Consider a prediction rule $h : \mathcal{X} \rightarrow \mathcal{Y}$ with $h(x) = \langle w, \psi(x) \rangle$ which is trained against a training sample $\{z_i = (x_i, y_i)\}_{i=1}^m$ with the following general optimization problem

$$(3.1) \quad \underset{w}{\text{minimize}}(f(\langle w, \psi(x_1) \rangle, \dots, \langle w, \psi(x_m) \rangle) + R(\|w\|)),$$

where $f : \mathbb{R}^m \rightarrow \mathbb{R}$ is an arbitrary function and $R : \mathbb{R}_+ \rightarrow \mathbb{R}$ is a monotonically non-decreasing function.

Example 19. (Cont. Example 11) Soft SVM problem has a solution equivalent to (see Proposition 29, Handout 8 Support Vector Machines)

$$(w^*, b^*) = \arg \min_{(w,b)} \left(\frac{1}{m} \sum_{i=1}^m \max(0, 1 - y_i (\langle w, \psi(x_i) \rangle + b)) + \lambda \|w\|_2^2 \right)$$

then in terms of (3.1) we get

$$f(\alpha_1, \dots, \alpha_m) = \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y_i \alpha_i\}, \text{ and } R(\beta) = \lambda \beta^2$$

Definition 20. Kernel function K is defined as $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{C}$ with $K(x, x') = \langle \psi(x), \psi(x') \rangle$ given an embedding $\psi(x)$ of some domain space \mathcal{X} into some Hilbert space \mathcal{F} . Kernel functions describe inner products in the feature space \mathcal{F} .

Theorem 21. (Representation theorem) Assume mapping $\psi : \mathcal{X} \rightarrow \mathcal{F}$ where \mathcal{F} is a Hilbert space. There exists a vector $\alpha \in \mathbb{R}^m$ such that $w = \sum_{i=1}^m \alpha_i \psi(x_i)$ is the optimal solution of (3.1) in Problem 18.

Proof. Let w^* be the optimal solution of (3.1). Because w^* is element of Hilbert space, it can be written as $w^* = \sum_{i=1}^m \alpha_i \psi(x_i) + u$ where $\langle u, \psi(x_i) \rangle = 0$ for all $i = 1, \dots, m$. Set $w := w^* - u$.

Because $\|w^*\|^2 = \|w\|^2 + \|u\|^2$ it is $\|w\| \leq \|w^*\|$ implying that

$$R(\|w\|) \leq R(\|w^*\|).$$

Because $\langle w, \psi(x_i) \rangle = \langle w^* - u, \psi(x_i) \rangle = \langle w^*, \psi(x_i) \rangle$ for all $i = 1, \dots, m$, it is

$$f(\langle w, \psi(x_1) \rangle, \dots, \langle w, \psi(x_m) \rangle) = f(\langle w^*, \psi(x_1) \rangle, \dots, \langle w^*, \psi(x_m) \rangle)$$

Then the objective function of (3.1) at w is less than or equal to that of the minimizer w^* which implies that $w = \sum_{i=1}^m \alpha_i \psi(x_i)$ is an optimal solution. \square

Note 22. Let $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a kernel function with $K(x, x') = \langle \psi(x), \psi(x') \rangle$. According to the representation Theorem 21, the Learning Problem 18, can be equivalently addressed by re-writing the learning predictive rule as

$$h_\alpha(x) = \sum_{i=1}^m \alpha_i K(x_i, x)$$

and learning $\{\alpha_i\}$ as the solutions of

$$(3.2) \quad \underset{\alpha}{\text{minimize}} \left(f \left(\sum_{i=1}^m \alpha_i K(x_i, x_1), \dots, \sum_{i=1}^m \alpha_i K(x_i, x_m) \right) + R \left(\sqrt{\sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j K(x_i, x_j)} \right) \right),$$

This is because

$$\langle w, \psi(x_j) \rangle = \langle \sum_{i=1}^m \alpha_i \psi(x_i), \psi(x_j) \rangle = \sum_{i=1}^m \alpha_i \langle \psi(x_i), \psi(x_j) \rangle = \sum_{i=1}^m \alpha_i K(x_i, x_j)$$

and

$$\|w\|^2 = \langle \sum_{i=1}^m \alpha_i \psi(x_i), \sum_{j=1}^m \alpha_j \psi(x_j) \rangle = \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j \langle \psi(x_i), \psi(x_j) \rangle = \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j K(x_i, x_j)$$

Note 23. In Learning Problem 18, direct access to elements $\psi(\cdot)$ in the feature space is not necessary because equivalently one can calculate or just specify the associated kernel function (that is inner products in the feature space).

Example 24. (Cont. Example 19) In Soft SVM the predictive rule becomes $h(x) = \text{sign} \left(\sum_j \alpha_j K(x_i, x_j) \right)$ and the learning task becomes

$$\underset{\alpha}{\text{minimize}} \left(\lambda \sum_i \sum_j \alpha_i \alpha_j K(x_i, x_j) + \frac{1}{m} \sum_i \max \left(0, 1 - y_i \sum_j \alpha_j K(x_i, x_j) \right) \right)$$

which can be addressed via SGD. This form of SVM is called Kernel SVM because we can just directly specify the kernel function $K(\cdot, \cdot)$ without the need to even think about feature mapping $\psi(\cdot)$ (which is eliminated and replaced by the kernel).

Example 25. (Polynomial Kernels) Let $x \in \mathcal{X} \subseteq \mathbb{R}^n$. We want to extend the linear mapping $x \mapsto \langle w, x \rangle$ to the k degree polynomial mapping $x \mapsto h(x)$. The multivariate polynomial can be written as $h(x) = \langle w, \psi(x) \rangle$, where $\psi : \mathbb{R}^n \rightarrow \mathbb{R}^d$ with $\psi(x)$ being a vector of elements $\psi_J(x) = \frac{\sqrt{k!}}{\sqrt{j_0!}} (\sqrt{c})^{j_0} \prod_{i=1}^n \frac{(k!)^{\frac{1}{2n}}}{(j_i!)^{\frac{1}{2}}} (x_i)^{j_i}$ for $J = (j_0, \dots, j_n)^\top$ and $j_0 + j_1 + \dots + j_n = k$. This learning problem can be equivalently be addressed with the k degree polynomial kernel

$$K(x, t) = (c + \langle x, t \rangle)^k; \text{ where } x, t \in \mathcal{X}$$

Solution: It is

$$\begin{aligned} K(x, t) &= (c + \langle x, t \rangle)^k = (\text{by setting } x_0 = t_0 = \sqrt{c}) = \left(\sum_{j=0}^n x_j t_j \right)^k = (\text{by Multinomial theorem}) = \\ &= \sum_{J:j_0+j_1+\dots+j_n=k} \frac{k!}{j_0! j_1! \cdots j_n!} \prod_{i=0}^n (x_i t_i)^{j_i} = \sum_{J:j_0+j_1+\dots+j_n=k} \underbrace{\left[\prod_{i=0}^n \frac{(k!)^{\frac{1}{2n}}}{(j_i!)^{\frac{1}{2}}} (x_i)^{j_i} \right]}_{\psi_J(x)} \underbrace{\left[\prod_{i=0}^n \frac{(k!)^{\frac{1}{2n}}}{(j_i!)^{\frac{1}{2}}} (t_i)^{j_i} \right]}_{\psi_J(t)} \\ &= \langle \psi(x), \psi(x') \rangle \quad \text{where } \psi(x) \text{ is as defined.} \end{aligned}$$

Solution. ► In Figure 1.1c, points can be distinguished by some ellipse, hence ψ could be defined as a vector of monomials up to order. Alternatively a degree 2 polynomial kernel could be used.

Example 26. (Radial basis kernel) Let the original input space be $x \in \mathcal{X} \subseteq \mathbb{R}$. Consider the Radial Basis Functions Kernel (or Gaussian kernel)

$$K(x, x') = \exp \left(-\frac{1}{2\sigma^2} \|x - x'\|_2^2 \right).$$

Show that it is a kernel indeed, by presenting it as an inner product in a feature space of infinite dimension, and state the bases of the mapping $\psi(\cdot)$.

Solution: It is

$$\begin{aligned} K(x, x') &= \exp \left(-\frac{1}{2\sigma^2} \|x - x'\|_2^2 \right) = \exp \left(\frac{1}{\sigma^2} xx' - \frac{1}{2} x^2 - \frac{1}{2} (x')^2 \right) \\ &= \exp \left(\frac{1}{\sigma^2} xx' \right) \exp \left(-\frac{1}{2\sigma^2} x^2 \right) \exp \left(-\frac{1}{2\sigma^2} (x')^2 \right) \\ &= \sum_{k=0}^{\infty} \frac{(xx'/\sigma^2)^k}{k!} \exp \left(-\frac{1}{2\sigma^2} x^2 \right) \exp \left(-\frac{1}{2\sigma^2} (x')^2 \right) \text{ (...by Taylor Expansion)} \\ &= \sum_{k=0}^{\infty} \left[\frac{x^k}{\sqrt{k!}\sigma^k} \exp \left(-\frac{1}{2\sigma^2} x^2 \right) \right] \left[\frac{(x')^k}{\sqrt{k!}\sigma^k} \exp \left(-\frac{1}{2\sigma^2} (x')^2 \right) \right] \end{aligned}$$

hence it is $K(x, x') = \langle \psi(x), \psi(x') \rangle$ with $\psi_k(x) = \frac{x^k}{\sqrt{k!}\sigma^k} \exp \left(-\frac{1}{2\sigma^2} x^2 \right)$.

4. CONSTRUCTION OF KERNELS

Note 27. The kernel formulated as an inner product in a feature space allows us to build interesting extensions of many well known algorithms by making use of the kernel trick and without the need to have direct access to the feature space (E.g. Example 24).

Note 28. Specifying a kernel function K is a way to express prior knowledge without the need to have direct access to the feature space. This is consequence of the Representation Theorem 21 that kernel is the inner product of feature mappings ψ which sufficiently replaces them in the learning problem, and the fact that ψ is a way to express and utilize prior knowledge about the problem at hand.

Note 29. Theorem 32 provides sufficient and necessary conditions to check whether the specified function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is indeed a kernel function; i.e. if K can be written as inner product $K(x, x') = \langle \psi(x), \psi(x') \rangle$ of feature functions $\psi(x)$.

Definition 30. Gram matrix is called the $m \times m$ matrix G s.t. $[G]_{i,j} = K(x_i, x_j)$.

Definition 31. A symmetric function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is positive semi-definite if its Gram matrix G , $[G]_{i,j} = K(x_i, x_j)$, is a positive semi-definite matrix.

Theorem 32. A symmetric function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ implements an inner product in some Hilbert space is a valid kernel function if and only if it is positive semi-definite i.e. its Gram matrix G , $[G]_{i,j} = K(x_i, x_j)$, is a positive semi-definite matrix.

Proof. Assume K is a valid kernel function (i.e. it implements an inner product in some Hilbert space) $K(x, x') = \langle \psi(x), \psi(x') \rangle$; let's consider $\psi : \mathcal{X} \rightarrow \mathbb{R}^d$ for simplicity. Let G be its Gram matrix with $G = \Psi^\top \Psi$ and $\psi(x_i)$ is the i -th column of Ψ . For any $\xi \in \mathbb{R}^d - \{0\}$

$$\begin{aligned}\xi^\top G \xi &= \sum_i \sum_j \xi_i K(x_i, x_j) \xi_j = \sum_i \sum_j \xi_i \langle \psi(x_i), \psi(x_j) \rangle \xi_j = \sum_i \sum_j \langle \xi_i \psi(x_i), \psi(x_j) \xi_j \rangle \\ &= \langle \sum_i \xi_i \psi(x_i), \sum_j \psi(x_j) \xi_j \rangle = \left\| \sum_i \xi_i \psi(x_i) \right\|_2^2 \geq 0\end{aligned}$$

Assume the symmetric function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is positive semi-definite. Let $\mathbb{R}^f = \{f : \mathcal{X} \rightarrow \mathbb{R}\}$. For $x \in \mathcal{X}$ let function ψ over \mathbb{R}^f with $\psi(x) = K(\cdot, x)$. This allows to define a vector space consisting of all the linear combinations of elements of the form $K(\cdot, x)$, having an inner product

$$\langle \sum_i \alpha_i K(\cdot, x_i), \sum_j \beta_j K(\cdot, x_j) \rangle = \sum_i \sum_j \alpha_i \beta_j \underbrace{\langle K(\cdot, x_i), K(\cdot, x_j) \rangle}_{=K(x_i, x_j)}.$$

This satisfies all the properties of inner product, s.t. it is symmetric, linearity, positive definite as $K(x, x') \geq 0$. Then there is some feature vector ψ such that $K(x, x') = \langle \psi(x), \psi(x') \rangle$. \square

Note 33. A powerful technique for constructing new kernels is to build them out of simpler kernels as building blocks. Below are some properties. Assume $K_1 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ and $K_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ are valid kernels, then the following are kernels too

- (1) $K(x, x') = K_1(x, x') + K_2(x, x')$
- (2) $K(x, x') = K_1(x, x') K_2(x, x')$
- (3) $K(x, x') = K_1(x_1, x'_1) + K_2(x_2, x'_2)$, where $x = (x_1, x_2)^\top$, $x' = (x'_1, x'_2)^\top$
- (4) $K(x, x') = K_1(x_1, x'_1) K_2(x_2, x'_2)$, where $x = (x_1, x_2)^\top$, $x' = (x'_1, x'_2)^\top$
- (5) $K(x, x') = f(x) K_1(x, x') f(x')$ for any function f
- (6) $K(x, x') = K_1(f(x), f(x'))$ for any function f

See the proof as a solution to an Example

Proof. We present the first two and the rest are proved similarly.

For (1). Let Gram matrix, G_j induced by kernel function K_j . For any $\xi \in \mathbb{R}^d - \{0\}$

$$\xi^\top G_3 \xi = \xi^\top (G_1 + G_2) \xi = \xi^\top G_1 \xi + \xi^\top G_2 \xi \geq 0$$

For (2). Assume that $K_j(x, x') = (\psi_j(x))^\top \psi_j(x)$. Then

$$\begin{aligned} K(x, x') &= K_1(x, x') K_2(x, x') = (\psi_1(x))^\top \psi_1(x') (\psi_2(x))^\top \psi_2(x') \\ &= (\psi_1(x))^\top \psi_2(x) (\psi_1(x'))^\top \psi_2(x') = ((\psi_1(x))^\top \psi_2(x))^\top (\psi_1(x'))^\top \psi_2(x') \end{aligned}$$

which can be represented as an inner product of feature vectors.

□

Note 34. The concept of a kernel formulated as an inner product in a feature space allows us to build interesting extensions of many well known algorithms by making use of the kernel trick. One example was the Kernel SVM. Some other popular cases are Gaussian process regression, and Kernel PCA.

Lecture notes 10: Multi-class classification

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim. To introduce the Support Vector Machines as a procedure. Motivation, set-up, description, computation, and implementation. We focus on the classical treatment.

Reading list & references:

- (1) Shalev-Shwartz, S., & Ben-David, S. (2014). Understanding machine learning: From theory to algorithms. Cambridge university press.
 - Ch. 17 (pp. 190-198)
- (2) Bishop, C. M. (2006). Pattern recognition and machine learning (Vol. 4, No. 4, p. 738). New York: Springer.
 - Ch. 7.1 Sparse Kernel Machines/Maximum marginal classifiers
- (3) Vapnik, V. (2013). The nature of statistical learning theory. Springer science & business media.
- (4) Boyd, S. P., & Vandenberghe, L. (2004). Convex optimization. Cambridge university press.
- (5) Strang, G. (2019). Linear algebra and learning from data. Wellesley-Cambridge Press.

1. INTRO AND MOTIVATION

Note 1. Multi-class classification is the machine learning problem of classifying instances into one of several possible target classes.

Example 2. Applications include: (i.) categorizing documents according to topic (input features is the set of documents and target is the set of possible topics); (ii.) determining which object appears in a given image (input is the set of images and target is the set of possible objects); etc...

Note 3. We wish to classify objects with given features $x \in \mathcal{X}$ into one of the $k \in \mathbb{N} - \{0\}$ categories $\{\mathcal{C}_1, \dots, \mathcal{C}_k\}$.

Note 4. To achieve that we aim to learn a predictive rule $h : \mathcal{X} \rightarrow \mathcal{Y}$ from a hypothesis class \mathcal{H} i.e. $h \in \mathcal{H}$. Here, to easy the notation we define a set of targets $\mathcal{Y} = \{1, \dots, k\}$ implying that the multi-class categorization involves $k \in \mathbb{N} - \{0\}$ possible target classes . If the object $z = (x, y)$ belongs to class \mathcal{C}_j i.e. $z \in \mathcal{C}_j$ then $y = j$.

Note 5. There is available a training data-set $\mathcal{S} = \{z_i = (x_i, y_i) ; i = 1, \dots, n\}$ of n examples where $x_i \in \mathcal{X}$ and $y_i \in \mathcal{Y}$ independently generated from a data-generating process g ; i.e. $\mathcal{S} \sim g$.

2. REDUCTION TO THE BINARY CLASSIFICATION APPROACHES

Note 6. The multi-class categorization problem can be addressed by a reduction to the binary classification. –We discuss two simple ways.

Note 7. Assume that there is available a binary classification learning algorithm $\mathfrak{A}_{\text{binary}}$, such as the soft-SVM or the Bernoulli regression.

2.1. One-versus-all approach.

Note 8. The one-versus-all approach involves training $k > 2$ binary classifiers, each of which discriminates between one class and the rest of the classes.

Note 9. It works as follows

- (1) For each class $j = 1, \dots, k$, construct a binary training data-sets

$$\mathcal{S}^{[j]} = \left\{ \left(x_1, t_1^{[j]} \right), \dots, \left(x_n, t_n^{[j]} \right) \right\}$$

where $t_i^{[j]} = \begin{cases} -1 & y_i \neq j \\ 1 & y_i = j \end{cases}$ from the original available training data-set $\mathcal{S} = \{z_i = (x_i, y_i); i = 1, \dots, n\}$.

- $\mathcal{S}^{[j]}$ is the set of examples labeled 1 if their label in \mathcal{S} was j , and -1 otherwise.
- (2) For each binary training data-set $\{\mathcal{S}^{[j]}\}$, consider a binary predictor rule $h^{[j]} : \mathcal{X} \rightarrow \{-1, +1\}$ with $h^{[j]} \in \mathcal{H}^{[j]}$ and train it against $\mathcal{S}^{[j]}$ such that

$$h^{[j]} = \mathfrak{A}_{\text{binary}}(\mathcal{S}^{[j]})$$

- (3) Given $\{h^{[j]}\}$, the one-versus-all multiclass hypothesis is defined as

$$h_{\text{ova}}(x) \in \arg \max_{j \in \mathcal{Y}} \left\{ h^{[j]}(x) \right\}$$

Note 10. In cases where there are ties, it is reasonable to decide the classification based on some reasonable weights. E.g. if algorithm $\mathfrak{A}_{\text{binary}}$ is a SVM with hypothesis of the form $h^{[j]}(x) = \text{sign}(\langle w^{[j]}, x \rangle)$, and I get $h^{[j]}(x) = 1 = h^{[j']}(x)$ with $\langle w^{[j]}, x \rangle > \langle w^{[j']}, x \rangle$ then $y = j$.

Note 11. In Note 9(2), we make the strong assumption that hoping that $h^{[j]}$ should equal 1 if and only if x belongs to class j . However, this is not always true. Figure 2.1 shows a case involving three classes $\mathcal{C}_?$ where this approach leads to regions $\mathcal{R}_?$ of input space that are ambiguously classified.

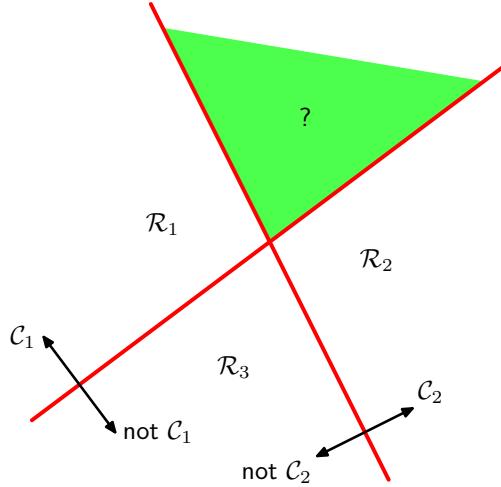


FIGURE 2.1. Two binary classification rules $h^{[1]}$ and $h^{[2]}$ designed to distinguish points in class \mathcal{C}_j from points not in class \mathcal{C}_j for $j \in \{1, 2\}$.

2.2. One-versus-one approach.

Note 12. The one-versus-one approach involves the introduction of $k(k - 1)/2$ binary classification rules, one for every possible pair of classes, and their comparison to each other, in a manner that each point is classified according to a majority vote amongst these binary classification rules.

Note 13. It works as follows

- (1) For each pair of classes (i, j) where $1 \leq i < j \leq k$, construct binary training data-sets $\mathcal{S}^{[i,j]}$ containing all examples from \mathcal{S} whose label is either i or j .

I.e.: For each binary training data-set $\mathcal{S}^{[i,j]}$ set the binary label as 1 if the class is i and as -1 if the class is j ; i.e. Starting with $\mathcal{S}^{[i,j]} = \emptyset$ for $t = 1, \dots, m$

- if $y_t = i$ add $(x_t, 1)$ to $\mathcal{S}^{[i,j]}$
- if $y_t = j$ add $(x_t, -1)$ to $\mathcal{S}^{[i,j]}$

- (2) For each pair (i, j) where $1 \leq i < j \leq k$, train a binary classification algorithm against $\mathcal{S}^{[i,j]}$ to get

$$h^{[i,j]} = \mathfrak{A}_{\text{binary}}(\mathcal{S}^{[i,j]})$$

- (3) Given $\{h^{[i,j]}\}$, the one-versus-one multiclass hypothesis is defined as the one with the highest number of “wins”; i.e.

$$h_{\text{ovo}}(x) \in \arg \max_{i \in \mathcal{Y}} \left\{ \sum_{j \in \mathcal{Y}} \text{sign}(j - \textcolor{red}{i}) h^{[i,j]}(x) \right\}$$

Note 14. The one-versus-one approach suffers from the problem of ambiguous regions as seen in Figure 2.2.

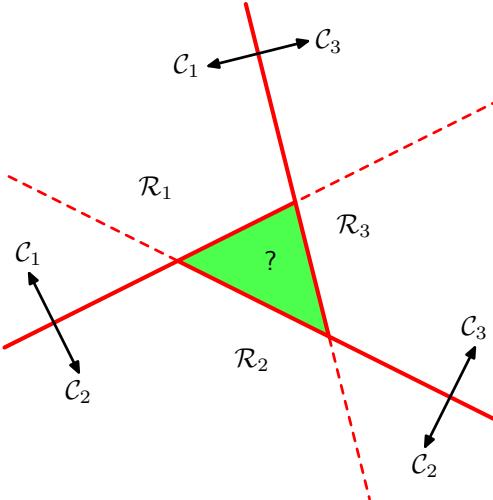


FIGURE 2.2. three binary classification rules each of which is used to separate a pair of classes \mathcal{C}_k and \mathcal{C}_j .

3. LINEAR MULTI-CLASS PREDICTORS

3.1. Setting up the predictive rule and class of hypotheses.

Note 15. We try to generalize the linear predictor for binary classifications taking the form

$$\begin{aligned} h_{\text{bin}}(x) &= \text{sign}(\langle w, x \rangle) \\ &= \arg \max_{y \in \{-1, +1\}} \langle w, yx \rangle \end{aligned}$$

into linear predictor for multi-class classifications.

Note 16. In the multi-class setting, we extend the binary SVM idea as follows.

- Define a multi-class predictor rule $h : \mathcal{X} \rightarrow \mathcal{Y}$ with formula

$$(3.1) \quad h_w(x) = \arg \max_{y \in \mathcal{Y}} \langle w, \Psi(x, y) \rangle$$

for $w \in \mathcal{W} \subset \mathbb{R}^d$.

- In (3.1), we need to specify a class-sensitive feature mapping $\Psi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$ whose function's $\Psi(x, y)$ role is to act as a score function that assesses how well the label y fits the instance x .
- The hypothesis class of multi-class predictors is

$$\mathcal{H}_{\Psi, \mathcal{W}} = \left\{ x \mapsto \arg \max_{y \in \mathcal{Y}} \langle w, \Psi(x, y) \rangle : w \in \mathcal{W} \right\}$$

Note 17. The equation of class-sensitive feature mapping Ψ depends on the application.

Example 18. An example of class-sensitive feature mapping Ψ such that

$$\Psi(x, y) = \begin{pmatrix} 0, \dots, 0, & \underbrace{x_1, \dots, x_n}_{n \text{ elements}}, & 0, \dots, 0, \\ (y-1)n \text{ elements} & & (k-y)n \text{ elements} \end{pmatrix}^\top$$

where $\Psi(x, y) \in \mathbb{R}^{nk}$ and $w \in \mathcal{W} \subset \mathbb{R}^{nk}$. In that way, the predictive rule simplifies to

$$\begin{aligned} h_w(x) &= \arg \max_{y \in \mathcal{Y}} \langle w, \Psi(x, y) \rangle \\ &= \arg \max_{y \in \mathcal{Y}} \langle w_y, x \rangle \end{aligned}$$

where $w = (w_1, \dots, w_k)^\top$ and each w_j has n elements.

3.2. About the learning problem.

Note 19. (0–1 loss function) The 0-1 loss function ℓ_{0-1} in the multi-class framework can be specified as

$$(3.2) \quad \ell_{0-1}(h_w, (x, y)) = 1_{\{y\}}(h_w(x)) = \begin{cases} 1, & h_w(x) \neq y \\ 0, & h_w(x) = y \end{cases}$$

as a straightforward extension of the one in binary classification.

Note 20. 0 – 1 loss (3.2) is non-convex as the one in the binary classification problem.

Note 21. (Hinge loss function) The hinge loss function ℓ_{hinge} in the multi-class framework can be specified as

$$(3.3) \quad \ell_{\text{hinge}}(h_w, (x, y)) = \max_{\xi \in \mathcal{Y}} (\ell_{0-1}(\xi, (x, y)) + \langle w, \Psi(x, \xi) - \Psi(x, y) \rangle)$$

extending on in binary classification. Example 24 provides a rational for its construction.

Example 22. The hinge loss function (3.3) is convex w.r.t. w

Solution: Straightforward as it is the maximum over linear functions of w .

Example 23. The hinge loss function (3.3) is ρ -Lipschitz w.r.t. w with $\rho = \max_{\xi \in \mathcal{Y}} (\Psi(x, \xi) - \Psi(x, y))$.

Solution: Let $g(w, \xi) = \ell_{0-1}(\xi, (x, y)) + \langle w, \Psi(x, \xi) - \Psi(x, y) \rangle$ and $f(w) = \ell_{\text{hinge}}(h_w, (x, y)) = \max_{\xi \in \mathcal{Y}} (g(w, \xi))$. It is

$$\begin{aligned} |f(w) - f(w')| &\leq \left| \max_{\xi \in \mathcal{Y}} (g(w, \xi)) - \max_{\xi \in \mathcal{Y}} (g(w', \xi')) \right| \leq \max_{\xi \in \mathcal{Y}} |g(w, \xi) - g(w', \xi)| \\ &\leq \max_{\xi \in \mathcal{Y}} |\langle w - w', \Psi(x, \xi) - \Psi(x, y) \rangle| \leq \max_{\xi \in \mathcal{Y}} (|\Psi(x, \xi) - \Psi(x, y)| |w - w'|) \\ &\leq \max_{\xi \in \mathcal{Y}} (|\Psi(x, \xi) - \Psi(x, y)|) |w - w'| \end{aligned}$$

Example 24. The hinge loss function (3.3) upper bounds the 0 – 1 loss function (3.2).

Solution: It is

$$\begin{aligned} \langle w, \Psi(x, y) \rangle &\leq \langle w, \Psi(x, h_w(x)) \rangle \quad \text{by construction} \\ \implies \ell_{0-1}(h_w, (x, y)) &\leq \ell_{0-1}(h_w(x), (x, y)) + \langle w, \Psi(x, h_w(x)) - \Psi(x, y) \rangle \\ &\leq \max_{\xi \in \mathcal{Y}} (\ell_{0-1}(\xi, (x, y)) + \langle w, \Psi(x, \xi) - \Psi(x, y) \rangle) \end{aligned}$$

Note 25. The hinge loss function ℓ_{hinge} can be used as a surrogate loss instead of the $0 - 1$ loss ℓ_{0-1} to address the non-convex learning problem with convex learning problem tools. This results from Note 22 and Example 24. Hence the computation of error bounds and the derivation of assumptions guarantee PAC learning algorithm is straightforward.

Example 26. The hinge loss function (3.3) reduces to the hinge loss $\ell_{\text{hinge}}(h_w(x), (x, y)) = \max(0, 1 - y \langle w, x \rangle)$ in binary classification.

Solution: In the binary scenario, I have $\mathcal{Y} = \{-1, +1\}$. If I set $\Psi(x, y) = \frac{1}{2}yx$ then (3.3) reduces to $\max(0, 1 - y \langle w, x \rangle)$.

Problem 27. The Multi-class SVM learning problem $(\mathcal{H}_{\Psi, W}, \mathcal{Z}, \ell_{\text{hinge}})$ using Ridge regularization with shrinkage parameter $\lambda > 0$ is

$$(3.4) \quad w^* = \arg \min_w \left(\underbrace{\frac{1}{m} \sum_{i=1}^m \max_{\xi \in \mathcal{Y}} (\ell_{0-1}(\xi, z_i) + \langle w, \Psi(x_i, \xi) - \Psi(x_i, y_i) \rangle)}_{\hat{R}_S(w) \text{ as empirical risk function}} + \lambda \|w\|_2^2 \right)$$

with predictive rule

$$h_{w^*}(x) = \arg \max_{y \in \mathcal{Y}} \langle w^*, \Psi(x, y) \rangle$$

3.3. About the learning algorithm.

Note 28. Problem 27 can be solved with linear programming algorithms, as well as any variation of the SGD algorithm. –We discuss the latter.

Example 29. The recursion of SGD applied in (3.4) is

$$w^{(t+1)} = w^{(t)} - \eta_t \left(\Psi(x^{(t)}, \hat{y}^{(t)}) - \Psi(x^{(t)}, y^{(t)}) \right)$$

with

$$\hat{y}^{(t)} \in \max_{\xi \in \mathcal{Y}} \left(\ell_{0-1}(\xi, z^{(t)}) + \langle w^{(t)}, \Psi(x^{(t)}, \xi) - \Psi(x^{(t)}, y^{(t)}) \rangle \right),$$

learning rate η_t and for a randomly chosen example $(x^{(t)}, y^{(t)}) \sim g$, at iteration t .

Solution: The recursion is

$$w^{(t+1)} = w^{(t)} - \eta_t v_t$$

where

$$v_t \in \partial_w \ell_{\text{hinge}}(h_w, z = (x, y))$$

I will use Proposition 11 from Lecture notes 2: Elements of convex learning problems. Because I can write

$$\ell_{\text{hinge}}(h_w, z) = \max_{\xi \in \mathcal{Y}} \{c(w; \xi)\}$$

with

$$c(w; \xi) = \ell_{0-1}(\xi, z_i) + \langle w, \Psi(x_i, \xi) - \Psi(x_i, y_i) \rangle$$

if \hat{y} is such as $c(w; \hat{y}) = \max_{\xi \in \mathcal{Y}} (c(w; \xi))$ then

$$\begin{aligned} v_t &\in \partial_w \ell_{\text{hinge}}(h_w, z = (x, y)) \\ &= \partial_w c(w; \hat{y}) \\ &= \partial_w \left(\ell_{0-1}(\hat{y}, z^{(t)}) + \langle w, \Psi(x^{(t)}, \hat{y}) - \Psi(x^{(t)}, y^{(t)}) \rangle \right) \\ &= \left\{ \nabla_w \langle w, \Psi(x^{(t)}, \hat{y}) - \Psi(x^{(t)}, y^{(t)}) \rangle \right\} \\ &= \left\{ \Psi(x^{(t)}, \hat{y}) - \Psi(x^{(t)}, y^{(t)}) \right\} \end{aligned}$$

Note 30. For the SGD implementation, the error bounds and the assumptions that guarantee a PAC learning algorithm can be derived similar to the binary case.

Lecture notes 11: Artificial neural networks

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim. To introduce the Artificial neural network as a model and procedure in classical and Bayesian framework. Motivation, set-up, description, computation, implementation, tricks. We focus on the Feedforward network.

Reading list & references:

- (1) Bishop, C. M., & Bishop, H. (2024). Deep learning: foundations and concepts. New York: Springer.^a
 - Ch. 5, 6, 8, 9
- (2) Bishop, C. M. (1995). Neural networks for pattern recognition. Oxford university press.
 - Ch. 4 The multi-layer perceptron
- (3) LeCun, Y., Bottou, L., Orr, G. B., & Müller, K. R. (2002). Efficient backprop. In Neural networks: Tricks of the trade (pp. 9–50). Berlin, Heidelberg: Springer Berlin Heidelberg.

^aAs alternative & more compact see Ch. 5 Neural Networks from “Bishop, C. M. (2006). Pattern recognition and machine learning (Vol. 4, No. 4, p. 738). New York: Springer.”

1. INTRO AND MOTIVATION

Note 1. Artificial Neural Networks (NN) are statistical models which have mostly been developed from the algorithmic perspective of machine learning. They were originally created as an attempt to model the act of thinking by modeling neurons in a brain. In ML, NN are used as global approximators.

Example 2. Consider a regression problem with predictive rule $h : \mathbb{R}^d \rightarrow \mathbb{R}^q$, suitable for cases where the examples (data) consist of input $x \in \mathbb{R}^d$, and output targets $y \in \mathbb{R}^q$. An example of a 2 layer artificial neural network formulating the predictive rule h is

$$h_k(x) = \sigma_2 \left(w_{2,k,0} + \sum_j w_{2,k,j} \sigma_1 \left(w_{1,j,0} + \sum_i w_{1,j,i} x_i \right) \right), \text{ for } k = 1, \dots, q$$

where $\{\sigma_j(\cdot)\}$ are known functions and $\{w_{\cdot,\cdot,\cdot}\}$ are unknown weights/coefficients. E.g., one may choose with $\sigma_2(\alpha) = \alpha$, $\sigma_1(\alpha) = 1/(1 + \exp(-\alpha))$.

Note 3. The original biological motivation for feed-forward NN stems from McCulloch & Pitts (1943) who published a seminal model of a NN as a binary thresholding device in discrete time, i.e.

$$n_j(t) = 1 \left(\sum_{\forall i \rightarrow j} w_{j,i} n_i(t-1) > \theta_i \right)$$

where the sum is over neuron i connected to neuron j ; $n_j(t)$ is the output of neuron i at time t and $0 < w_{j,i} < 1$ are attenuation weights. Thus the effect is to threshold a weighted sum of the inputs at value θ_i . Perhaps, such a mathematical model involving compositions of interconnected non-linear functions could be able to mimic human's learning mechanism and be implemented in a computing environment (with faster computational abilities) with purpose to discover patterns, make predictions, cluster, classify, etc...

Note 4. Mathematically, NN are rooted in the classical theorem by Kolmogorov stating (informally) that every continuous function $h(\cdot)$ on $[0, 1]^d$ can be written as

$$(1.1) \quad h(x) = \sum_{i=1}^{2d+1} F_i \left(\sum_{j=1}^d G_{i,j}(x_j) \right)$$

where $\{G_{i,j}\}$ and $\{F_i\}$ are continuous functions whose form depends on f . Perhaps, one may speculate that functions $\{G_{i,j}\}$ and $\{F_i\}$ can be approximated by sigmoids or threshold functions of the form $\sigma(w^\top x)$ allowing the number of tunable coefficients w to be high enough such that they can represent any function -hence the property of NN as global approximators.

no need
to
memorize
the
formula (1.1)

Example 5. Recall soft-SVM in (Example 11 in Handout 9: Kernel methods); the predictive rule h^ψ had a formula

$$(1.2) \quad h^\psi(x) = \text{sign}(b + \langle w, \psi(x) \rangle) = \text{sign} \left(b + \sum_j w_j \psi_j(x) \right),$$

with some embedding $\psi(x) = (\psi_1(x), \psi_2(x) \dots)^\top$ and unknown coefficients $\{w_j\}$. Assume there is no prior info about how to specify the formula of the “ideal” embedding ψ ; then the scientist could try to parameterize ψ such as

$$(1.3) \quad \psi_j(x) = \zeta \left(\omega_{j,0} + \sum_i \omega_{j,i} x_i \right),$$

for some chosen function $\zeta(\cdot)$ and unknown coefficients $\{\omega_{j,i}\}$, hoping that the unknown coefficients $\{\omega_{j,i}\}$ can be successfully tuned during the learning procedure against the dataset and hence result in an adaptive semi-parametric way to construct the embedding ψ from the training data without his/her interference. Based on Note 4, and assuming that a linear combination of $\{\zeta(\cdot)\}$, specified as sigmoids, could acceptably enough represent the fluctuations of $\{\psi_j(x)\}$ in (1.2), such an attempt to “adaptively construct the embedding ψ ” might be successful.

Example 6. (Cont Example 5) Having this in the mind, and while wanting to improve expressiveness, the scientist could become even more greedy, and consider the use of another embedding in (1.3) as

$$\psi_j(x) = \zeta \left(\omega_{j,0} + \sum_i \omega_{j,i} \phi_i(x) \right),$$

whose bases $\phi(x) = (\phi_1(x), \phi_2(x) \dots)^\top$ were parameterised as

$$\phi_i(x) = \xi \left(\varpi_{i,0} + \sum_s \varpi_{i,s} x_s \right)$$

for some chosen sigmoid $\xi(\cdot)$ and unknown coefficients $\{\varpi_{i,s}\}$. This recursive basis function composition creates a deep hierarchical structure which may improve approximation of the ideal embedding $\{\psi_j(x)\}$, and expressiveness of predictive rule h^ψ by sacrificing “interpretability” and “parsimony” of the “model”.

2. FEEDFORWARD NEURAL NETWORK (MATHEMATICAL SET-UP)

Note 7. An (Artificial) Neural Network (NN) is an interconnection of several sigmoids or threshold functions associated and arranged in layers, such that the outputs of one layer form the input of the next layer. Formally the structure of these interconnections can be depicted as a directed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ whose nodes \mathcal{V} correspond to vertices/neurons and edges \mathcal{E} correspond to links between them.

Note 8. A Feed-Forward Neural Network (FFNN) or else multi-layer perceptron is a special case of NN whose vertices can be numbered so that all connections go from a neuron (vertex) to one with a higher number. Hence, neurons (vertices) have one-way connections to other neurons such that the output of a lower numbered neuron feeds the input of the higher numbered neuron (in a forward manner). The neurons can be arranged in layers so that connections go from one layer to a later layer. FFNN can be depicted by a directed acyclic graph¹, $\mathcal{G} = (\mathcal{V}, \mathcal{E})$. (See Figure 2.1).

Note 9. Here we restrict ourselves to FFNN only.

Note 10. We assume that the network is organized in layers. The set of nodes is decomposed into a union of (nonempty) disjoint subsets

$$V = \cup_{t=0}^T V_t$$

such that every edge in \mathcal{E} connects a node from V_t to a node from V_{t+1} , for $t = 1, \dots, T$.

Note 11. The first layer V_0 is called **input layer**. If x has d dimensions, then the first layer V_0 contains d nodes.

Note 12. The last layer V_T is called **output layer**. If y has q dimensions, then the last layer V_T contains q nodes.

Note 13. The intermediate layers $\{V_1, \dots, V_{T-1}\}$ are called **hidden layers**.

Note 14. In a neural network, the nodes of the graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ correspond to neurons.

Note 15. The ***i*-th neuron of the *t*-th layer** is denoted as $v_{t,i}$.

Note 16. The output of neuron i in the input layer V_0 is simply x_i that is $o_{0,i}(x) = x_i$ for $i = \{1, \dots, d\}$.

¹(its vertices can be numbered so that all connections go from one vertex to another with a higher number)
Page 3 Created on 2025/05/03 at 06:58:39 by Georgios Karagiannis

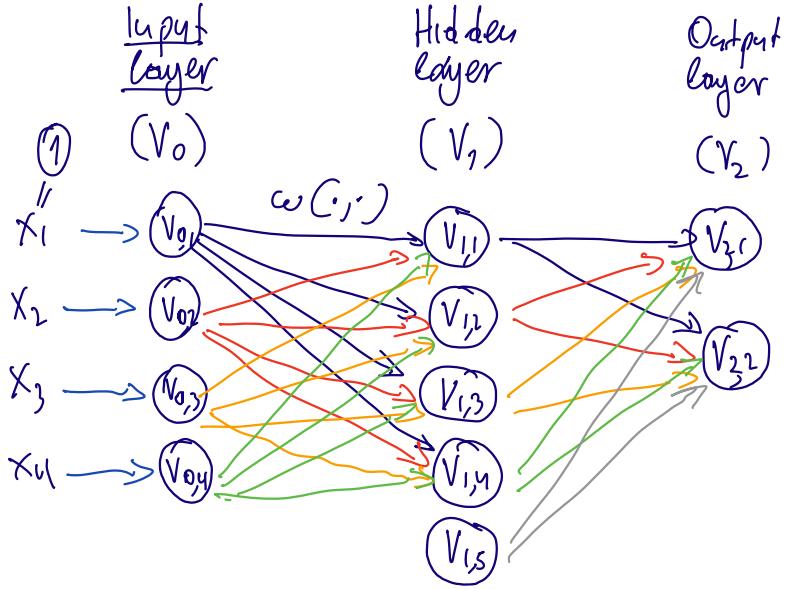


FIGURE 2.1. Feed forward neural network (1 hidden layer)

Note 17. Each edge in the graph $(v_{t,j}, v_{t+1,i})$ links the output of some neuron $v_{t,j}$ to the input of another neuron $v_{t+1,i}$; i.e. $(v_{t,j}, v_{t+1,i}) \in \mathcal{E}$.

Note 18. We define a **weight function** $w : \mathcal{E} \rightarrow \mathbb{R}$ over the edges \mathcal{E} .

Note 19. **Activation** of neuron i at hidden layer 1 is the weighted sum of the outputs $o_{0,i}(x) = x_i$ of the neurons in V_0 which are connected to $v_{1,i}$ where weighting is according to function w , that is

$$(2.1) \quad \alpha_{1,i}(x) = \sum_{\forall j: (v_{0,j}, v_{1,i}) \in \mathcal{E}} w((v_{0,j}, v_{1,i})) x_i$$

Note 20. Each single neuron $v_{t,i}$ is modeled as a simple scalar function, $\sigma_t(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$, called **activation function** at layer t .

Note 21. **The output of neuron** $v_{t,i}$ when the network is fed with the input x is denoted as $o_{t,i}(x) := \sigma_{t-1}(\alpha_{t-1,i}(x))$.

Note 22. Activation of neuron i at layer t is the weighted sum of the outputs $o_{t-1,j}(x)$ of the neurons in V_{t-1} which are connected to $v_{t,i}$ where weighting is according to function w , that is

$$(2.2) \quad \alpha_{t,i}(x) = \sum_{\forall j: (v_{t-1,j}, v_{t,i}) \in \mathcal{E}} w((v_{t-1,j}, v_{t,i})) o_{t-1,j}(x)$$

Note 23. The input to $v_{t+1,i}$ is activation $\alpha_{t+1,i}(x)$ namely a weighted sum of the outputs $o_{t,j}(x)$ of the neurons in V_t which are connected to $v_{t+1,i}$, where weighting is according to w . The output of $v_{t+1,i}$ is the application of the activation function $\sigma_{t+1}(\cdot)$ on its input $\alpha_{t+1,i}(x)$.

Algorithm 24. *The feed-forward NN formula in a layer by layer manner is performed (defined) according to the following recursion.*

At $t = 0, \text{ for } i = 1, \dots, |V_0|$

$$o_{0,i}(x) := x_i$$

At $t = 0, \dots, T - 1, \text{ for } i = 1, \dots, |V_{t+1}|$

$$\begin{aligned}\alpha_{t+1,i}(x) &= \sum_{\forall j: (v_{t,j}, v_{t+1,i}) \in \mathcal{E}} w((v_{t,j}, v_{t+1,i})) o_{t,j}(x) \\ o_{t+1,i}(x) &= \sigma_{t+1}(\alpha_{t+1,i}(x))\end{aligned}$$

Note 25. **Depth** of the NN is the number of the layers excluding the input layer; i.e. T .

Note 26. **Size** of the network is the number $|V|$.

Note 27. **Width** of the NN is the number $\max_{\forall t}(|V_t|)$.

Note 28. The **architecture** of the neural network is defined by the triplet $(\mathcal{V}, \mathcal{E}, \sigma_t)$.

Note 29. The neural network can be fully specified by the quadruplet $(\mathcal{V}, \mathcal{E}, \sigma_t, w)$.

Example 30. Figure 2.1 denotes a NN with depth 2, size 11, width 5. The neuron with no incoming edges has $o_{1,5} = \sigma(0)$.

Notation 31. To easy the notation, we denote the weights as $w_{(t+1),i,j} := w((v_{t,j}, v_{t+1,i}))$. Using this notation, $w_{(t+1),i,j} = 0$ is equivalent in (2.2) to $(v_{t,j}, v_{t+1,i}) \notin \mathcal{E}$ and means that the link $v_{t,j} \rightarrow v_{t+1,i}$ is not in the network.

Note 32. Often a **constant neuron** $v_{t,0}$ (at each layer t and $i = 0$) which outputs 1; i.e. $o_{0,0}(x) = 1$ and $o_{t,0}(x) = 1$. The corresponding weight $w_{(t),k,0}$ is called **bias**. This resembles to the constant term in the linear regression.

Example 33. (Cont. Example 2) The 2 layer neural network

$$(2.3) \quad h_k(x) = \sigma_{(2)} \left(w_{(2),k,0} + \sum_{\forall j} w_{(2),k,j} \sigma_{(1)} \left(w_{(1),j,0} + \sum_{\forall i} w_{(1),j,i} x_i \right) \right)$$

can be written according to the recursion in Algorithm 24 as

- Input layer

$$o_{(0),i}(x) = \begin{cases} 1 & i = 0 \\ x_i & i = 1, \dots, d \end{cases}$$

- Hidden layer

$$\begin{aligned}\alpha_{(1),j}(x) &= w_{(1),j,0} + \sum_j w_{(1),j,i} x_i \\ o_{(1),j}(x) &= \sigma_{(1)}(\alpha_{(1),j}(x)) \\ &= \sigma_{(1)}\left(w_{(1),j,0} + \sum_j w_{(1),j,i} x_i\right)\end{aligned}$$

- Output layer

$$\begin{aligned}\alpha_{(2),k}(x) &= w_{(2),k,0} + \sum_j w_{(2),k,j} o_{(1),k}(x) \\ o_{(2),k}(x) &= \sigma_{(2)}(\alpha_{(2),k}(x)) \\ &= \sigma_{(2)}\left(w_{(2),k,0} + \sum_j w_{(2),k,j} o_{(2),k}(x)\right)\end{aligned}$$

If $h_k(x)$ returns values in \mathbb{R} , we can choose $\sigma_{(2)}$ as the identity function i.e., $\sigma_{(2)}(\alpha) = \alpha$. Note that (2.3) is also presented more compact

$$(2.4) \quad h_k(x) = \sigma_{(2)}\left(\sum_j w_{(2),k,j} \sigma_{(1)}\left(\sum_i w_{(1),j,i} x_i\right)\right)$$

by considering the constant neuron associated with first x which is set equal to 1, e.g $x_1 = 1$, similar to the linear regression models.

Note 34. Activation functions σ_t are non-increasing functions often sigmoids or threshold functions. Their choice is problem-dependent. some examples

- Identity function: $\sigma(\alpha) = \alpha$ cannot be used in hidden layers
- Threshold sigmoid: $\sigma(\alpha) = 1(\alpha > 0)$
- Logistic sigmoid: $\sigma(\alpha) = \frac{1}{1+\exp(-\alpha)}$
- Rectified linear unit: $\text{ReLU}(\alpha) = \max(\alpha, 0)$

Example 35. Examples on the choice of the activation function at the output layer T :

- In the univariate regression problem with prediction rule $h(\cdot) \in \mathbb{R}$, and examples $z_i = (x_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$, we can choose $\sigma_T(\alpha) = \alpha$ to get $h(\alpha) = \sigma_T(\alpha) = \alpha$.
- In the binary logistic regression problem with prediction rule $h(\cdot) \in [0, 1]$ and examples $z_i = (x_i, y_i) \in \mathbb{R}^d \times \{0, 1\}$, we can choose $\sigma_T(\alpha) = \frac{1}{1+\exp(-\alpha)}$ to get $h(\alpha) = \sigma_T(\alpha) = \frac{1}{1+\exp(-\alpha)} = \frac{\exp(\alpha)}{1+\exp(\alpha)}$.

3. LEARNING NEURAL NETWORKS

Note 36. Assume we are interested in a prediction rule $h_{\mathcal{V}, \mathcal{E}, \sigma, w} : \mathbb{R}^{|V_0|} \rightarrow \mathbb{R}^{|V_T|}$ which is modeled as a feed-forward Neural Network with $(\mathcal{V}, \mathcal{E}, \sigma, w)$; that is

$$h_{\mathcal{V}, \mathcal{E}, \sigma, w}(x) = o_T(x)$$

where $o_T = (o_{T,1}, \dots, o_{T,|V_T|})^\top$ is according to the Algorithm 24.

Note 37. Learning the architecture $(\mathcal{V}, \mathcal{E}, \sigma)$ of a neural network is a model selection task (recall the variable selection in linear regression).

Note 38. We assume that the architecture $(\mathcal{V}, \mathcal{E}, \sigma)$ of the neural network $(\mathcal{V}, \mathcal{E}, \sigma, w)$ is fixed (given/known), and that there is interest in learning the weight function $w : \mathcal{E} \rightarrow \mathbb{R}$ or equivalently in vector form the vector of weights $\{w_{(t+1),i,j}\}$ where $w_{(t+1),i,j} := w((v_{t,j}, v_{t+1,i}))$.

Note 39. The class of hypotheses is

$$\mathcal{H}_{\mathcal{V}, \mathcal{E}, \sigma} = \{h_{\mathcal{V}, \mathcal{E}, \sigma, w} : \text{for all } w : \mathcal{E} \rightarrow \mathbb{R}\}$$

for given $(\mathcal{V}, \mathcal{E}, \sigma)$.

Notation 40. To simplify notation we will use h_w instead of $h_{\mathcal{V}, \mathcal{E}, \sigma, w}$ as $\mathcal{V}, \mathcal{E}, \sigma$ is fixed here.

Note 41. Assume that there is available a training set of examples (data-set) $\mathcal{S} = \{z_i = (x_i, y_i) ; i = 1, \dots, n\}$ with $x_i \in \mathcal{X} = \mathbb{R}^{|V_0|}$ and $y_i \in \mathcal{Y}$.

Note 42. To learn the unknown w , we need to specify a loss function $\ell(w, z)$ at some value of weight vector $w \in \mathbb{R}^{|\mathcal{E}|}$ and at some example $z = (x, y)$.

Definition 43. Error function is a performance measure that can be defined as

$$(3.1) \quad \text{EF}(w | \{z_i^*\}) = \sum_{i=1}^n \ell(w, z_i^*)$$

where $\mathcal{S}^* = \{z_i^* = (x_i^*, y_i^*) ; i = 1, \dots, n^*\}$ is a set of examples which does not necessarily need to be the training set \mathcal{S} .

Example 44. (Regression problem) Assume we wish to predict the mapping $x \xrightarrow{h(\cdot)} y$, where $x \in \mathbb{R}^d$ and $y \in \mathbb{R}$. Consider a predictive rule $h_w : \mathbb{R}^d \rightarrow \mathbb{R}$. Assume a training data-set $\{z_i = (x_i, y_i)\}$.

- The output activation function can be the identity function $\sigma_T(a) = a$. This is because $h_w(x) = o_T(x) = \sigma_T(\alpha_T(x))$ by Algorithm 24 and Note 36. Since h_w returns in \mathbb{R} and the output activation $\alpha_T(x)$ returns in \mathbb{R} , then mapping $\sigma_T(a) = a$ suffices.
- A reasonable loss can be the Euclidean distance

$$\ell(w, z = (x, y)) = \frac{1}{2} (h_w(x) - y)^2$$

- The same loss results as the minus log-likelihood of a specified statistical model

$$(3.2) \quad y_i | x, w \stackrel{\text{ind}}{\sim} N(\mu_i, \beta^{-1}), \text{ where } \mu_i = h_w(x_i)$$

for some fixed $\beta > 0$, i.e.

$$\begin{aligned}\ell(w, z = (x, y)) &= -\log(N(y|h_w(x), \beta^{-1})) \\ &= -\left(-\frac{1}{2} \log\left(\frac{1}{2\pi}\right) - \frac{1}{2} \log(\beta^{-1}) - \frac{1}{2} \frac{(h_w(x) - y)^2}{\beta^{-1}}\right) \\ &= \frac{1}{2}\beta(h_w(x) - y)^2 + \text{const}...\end{aligned}$$

- The Error function is

$$\text{EF}(w|z) = \sum_{i=1}^n \ell(w, z_i = (x_i, y_i)) = \frac{1}{2}\beta \sum_{i=1}^n (h_w(x_i) - y_i)^2 + \text{const}...$$

Example 45. (Multi-output regression problem) Assume we wish to predict the mapping $x \xrightarrow{h(\cdot)} y$, where $x \in \mathbb{R}^d$ and $y \in \mathbb{R}^q$. Consider a predictive rule $h_w : \mathbb{R}^d \rightarrow \mathbb{R}^q$. Assume a training data-set $\{z_i = (x_i, y_i)\}$.

- The output activation function is the identity function $\sigma_T(a) = a$. This is because $h_{w,k}(x) = o_{T,k}(x) = \sigma_T(o_{T,k}(x))$ for $k = 1, \dots, q$ by Algorithm 24 and Note 36. Since h_w returns in \mathbb{R}^q and the output activation is $\sigma_T(x)$ in \mathbb{R} , then mapping $\sigma_T(a) = a$ suffices.
- A reasonable loss can be an ellipsoidal norm

$$\ell(w, z = (x, y)) = (h_w(x) - y)^\top P(h_w(x) - y), \text{ for some } P > 0$$

- The same loss results as the minus log-likelihood of a specified statistical model

$$(3.3) \quad y_i|x_i, w \stackrel{\text{ind}}{\sim} N(\mu, P^{-1}), \text{ where } \mu_i = h_w(x_i)$$

for some fixed precision matrix $P > 0$, i.e.

$$\begin{aligned}\ell(w, z = (x, y)) &= -\log(N(y|h_w(x), P^{-1})) \\ &= -\left(-\frac{q}{2} \log\left(\frac{1}{2\pi}\right) + \frac{1}{2} \log(|P|) - \frac{1}{2} (h_w(x) - y)^\top P (h_w(x) - y)\right) \\ &= \frac{1}{2} (h_w(x) - y)^\top P (h_w(x) - y) + \text{const}...\end{aligned}$$

- The Error function is

$$\begin{aligned}\text{EF}(w|z) &= \sum_{i=1}^n \ell(w, z_i = (x_i, y_i)) \\ &= \frac{1}{2} \sum_{i=1}^n (h_{w,k}(x_i) - y_{k,i})^\top P (h_{w,k}(x) - y_{k,i}) + \text{const}...\end{aligned}$$

where $y_{k,i}$ is the k -th dimension of the i -th example in the dataset.

Example 46. (Binary classification problem) Assume we wish to classify objects with features $x \in \mathbb{R}^d$ in 2 categories. Consider a predictive rule $h_w : \mathbb{R}^d \rightarrow (0, 1)$ as a classification probability i.e., $h_w(x) = \Pr(x \text{ belongs to class 1})$. Assume a training data-set $\{z_i = (x_i, y_i)\}$ with $y_i \in \{0, \dots, q\}$ labeling the class.

- A suitable output activation function can be the logistic sigmoid

$$\sigma_T(a) = \frac{1}{1 + \exp(-a)}$$

This is because $h_w(x) = o_T(x) = \sigma_T(\alpha_T(x))$ by Algorithm 24 and Note 36. Since h_w returns in $(0, 1)$ and the output activation is $\alpha_T(x)$ in \mathbb{R} , then the aforesaid mapping suffices.

- A reasonable loss (among others) can be the cross entropy

$$\ell(w, z = (x, y)) = -y \log(h_w(x)) - (1 - y) \log(1 - h_w(x))$$

- The same loss results as the minus log-likelihood of a specified statistical model

$$(3.4) \quad y_i|x_i, w \stackrel{\text{ind}}{\sim} \text{Bernoulli}(p_i), \text{ where } p_i = h_w(x_i)$$

with mass function

$$f(y|x, w) = h_w(x)^y (1 - h_w(x))^{1-y}$$

i.e.

$$\begin{aligned} \ell(w, z = (x, y)) &= -\log(f(y|x, w)) \\ &= -y \log(h_w(x)) - (1 - y) \log(1 - h_w(x)) \end{aligned}$$

- The Error function given a set of examples $\{z_i^* = (x_i^*, y_i^*)\}$ is

$$\begin{aligned} \text{EF}(w|z^*) &= \sum_{i=1}^n \ell(w, z_i^* = (x_i^*, y_i^*)) \\ &= - \sum_{i=1}^n y_i^* \log(h_w(x_i^*)) - (1 - y_i^*) \log(1 - h_w(x_i^*)) \end{aligned}$$

Example 47. (Multi-class classification problem) Assume we wish to classify objects with features $x \in \mathbb{R}^d$ in q categories. Consider a predictive rule $h_w : \mathbb{R}^d \rightarrow \mathcal{P}$, with $\mathcal{P} = \left\{ \varpi \in (0, 1)^q : \sum_{j=1}^q \varpi_j = 1 \right\}$ and $h_w = (h_{w,1}, \dots, h_{w,q})^\top$, as a classification probability i.e., $h_{w,k}(x) = \Pr(x \text{ belongs to class } k)$. Assume a training data-set $\{z_i = (x_i, y_i)\}$ where y_i , with $y_{i,k} \in \{0, 1\}$ and $\sum_{k=1}^q y_{i,k} = 1$, labeling the class to which the i -the example belongs. Then it is

- A suitable output activation function can be the softmax function

$$(3.5) \quad \sigma_T(a_k) = \frac{\exp(a_k)}{\sum_{k'=1}^q \exp(a_{k'})}, \text{ for } k = 1, \dots, q$$

This is because $h_{w,k}(x) = o_{T,k}(x) = \sigma_{T,k}(\alpha_{T,k}(x))$ by Algorithm 24 and Note (36). Since h_w is essentially a probability vector and the output activation is $\sigma_{T,k}(x)$ in \mathbb{R} the aforesaid mapping suffices. Unfortunately, (3.5) is invariant to additive transformations $a \leftarrow a + c$ i.e., $\sigma_T(a_k) = \sigma_T(a_k + \text{const})$.

- Another suitable output activation function can be

$$(3.6) \quad \tilde{\sigma}_T(a_k) = \frac{\exp(a_k)}{1 + \sum_{k'=1}^{q-1} \exp(a_{k'})}, \text{ for } k = 1, \dots, q-1$$

This is because $h_{w,k}(x) = o_{T,k}(x) = \sigma_{T,k}(\alpha_{T,k}(x))$ by Algorithm 24 and Note 36. Since h_w is essentially a probability vector and the output activation is $\alpha_{T,k}(x)$ in \mathbb{R} the aforesaid mapping suffices as $h_{w,k}(x) = o_{T,k}(x)$ for $k = 1, \dots, q-1$ and $h_{w,q}(x) = 1 - \sum_{k=1}^{q-1} h_{w,k}(x)$. The advantage of (3.6) compared to (3.5) is that the former uses one less output neurons (less unknown parameters to learn). Also (3.6) is not invariant to additive transformations $a \leftarrow a + c$ unlike (3.5)i.e, $\tilde{\sigma}_T(a_k) \neq \tilde{\sigma}_T(a_k + \text{const})$.

- A reasonable loss can be the cross entropy

$$\ell(w, z = (x, y)) = - \sum_{k=1}^q y_k \log(h_{w,k}(x))$$

- The same loss results as the minus log-likelihood of a specified statistical model

$$y_i|x, w \sim \text{Multinomial}(p_i), \text{ where } p_i = h_w(x_i)$$

with mass function

$$f(y_i|x, w) = \prod_{k=1}^q h_{w,k}(x_i)^{y_{i,k}}$$

i.e.

$$\begin{aligned} \ell(w, z = (x, y)) &= -\log(f(y|x, w)) \\ &= - \sum_{k=1}^q y_k \log(h_{w,k}(x)) \end{aligned}$$

which is the cross-entropy.

- The Error function given a set of examples $\{z_i^* = (x_i^*, y_i^*)\}$ is

$$\begin{aligned} \text{EF}(w|z^*) &= \sum_{i=1}^n \ell(w, z_i^* = (x_i^*, y_i^*)) \\ &= - \sum_{i=1}^n \sum_{k=1}^q y_{k,i}^* \log(h_{w,k}(x_i^*)) \end{aligned}$$

where $y_{k,i}$ is the k -th dimension of the i -th example in the dataset.

4. CLASSICAL LEARNING OF NEURAL NETWORK

Note 48. Our purpose is to find optimal $h_w \in \mathcal{H}_{\mathcal{V}, \mathcal{E}, \sigma}$ under loss $\ell(\cdot, \cdot)$ and against training data-set $\mathcal{S} = \{z_i = (x_i, y_i); i = 1, \dots, n\}$.

4.1. Standard learning problem.

Note 49. Essentially, this is an optimization problem, where the objective is to minimize either the risk function $R_g(w)$ or the empirical risk function $\hat{R}_S(w)$.

Problem 50. Compute optimal $w^* \in \mathbb{R}^{|\mathcal{E}|}$ by minimizing the risk function $R_g(w)$

$$(4.1) \quad w^* = \arg \min_w (R_g(w)) = \arg \min_w (\mathbb{E}_{z \sim g} (\ell(w, z)))$$

Problem 51. Compute optimal $w^* \in \mathbb{R}^{|\mathcal{E}|}$ by minimizing the empirical risk function $\hat{R}_S(w)$

$$(4.2) \quad w^* = \arg \min_w \left(\hat{R}_S(w) \right) = \arg \min_w \left(\frac{1}{n} \sum_{i=1}^n \ell(w, z_i) \right)$$

4.2. Regularized loss optimization.

Note 52. Often neural network models are over-parameterized, in the sense that the dimensionality of $w \in \mathbb{R}^{|\mathcal{E}|}$ is too large, with negative consequences in its predictability. This can be addressed by learning the architecture NN, precisely by learning which of the edges of the FFNN significantly contribute to the NN model and should be kept, and which do not contribute and may be removed (or be inactive).

Note 53. To address Note 52, one can resort to shrinkage methods e.g., LASSO, Ridge, which indirectly allow edge/weight selection/elimination by shrinking the values of the weights $\{w_{(t),i,j}\}$ towards zero, and setting some of them as $w_{(t),i,j} = 0$ if their absolute value is small enough. This is based on the observation in (2.2) i.e., $w_{(t),i,j} = 0$ is equivalent to $(v_{t,j}, v_{t+1,i}) \notin \mathcal{E}$ which implies that the link $v_{t,j} \rightarrow v_{t+1,i}$ is not active (essentially not in the neural network).

Problem 54. Compute $h_w \in \mathcal{H}_{\mathcal{V}, \mathcal{E}, \sigma}$ (essentially compute $w \in \mathbb{R}^{|\mathcal{E}|}$) under loss $\ell(\cdot, \cdot)$ and shrinkage term (or weight decay) $J(w; \lambda)$, and against training data-set $\mathcal{S} = \{z_i = (x_i, y_i); i = 1, \dots, n\}$. Compute $w^* \in \mathbb{R}^{|\mathcal{E}|}$, according to the minimization:

$$(4.3) \quad w^* = \arg \min_w (R_g(w) + J(w; \lambda))$$

$$(4.4) \quad = \arg \min_w (\mathbb{E}_{z \sim g} (\ell(w, z) + J(w; \lambda)))$$

and set $w_{t,i,j}^* = 0$ if $w_{t,i,j}^*$ is less than a threshold user specific value $\xi > 0$ i.e. $|w_{t,i,j}^*| < \xi$.

Note 55. Popular shrinkage terms $J(w; \lambda)$ are

- Ridge: $J(w; \lambda) = \lambda \|w\|_2^2$
- LASSO: $J(w; \lambda) = \lambda \|w\|_1$
- Elastic net: $J(w; \lambda = (\lambda_1, \lambda_2)) = \lambda_1 \|w\|_1 + \lambda_2 \|w\|_2^2$

5. STOCHASTIC GRADIENT DESCENT FOR CLASSICAL TRAINING OF NN

Note 56. Training a neural network model is usually a high-dimensional problem (essentially w has high dimensionality to provide a better approximation) and a big-data problem (essentially we need a large number of training examples to learn a large number of weights). For this reason, Stochastic Gradient Descent (and its variations) is a suitable computational learning tool.

Note 57. To address Problem 50, the recursion of the SGD with batch size m is

$$w^{(t+1)} = w^{(t)} - \eta_t \frac{1}{m} \sum_{j=1}^m \partial_w \ell(w^{(t)}, z_j^{(t)})$$

Note 58. To address the Problem 54, the recursion of the SGD with batch size m is

$$w^{(t+1)} = w^{(t)} - \eta_t \left[\frac{1}{m} \sum_{j=1}^m \partial_w \ell \left(w^{(t)}, z_j^{(t)} \right) + \partial_w J(w; \lambda) \right]$$

for some positive λ which is user specified, or chosen via cross validation.

Note 59. The learning problem associated to the Neural network model is (almost always) non-convex due to the non-convex loss with respect to the w 's. Upon implementing SGD in the learning problem of Neural Network, the theoretical results in Section 4 (Handout 2) and Section 3 (Handout 3) will not be effective due to the violation of the assumptions.

Note 60. Practical guidelines for the use of SGD in the learning problem of NN:

Ref [3]

- (1) Utilize a learning rate η_t that changes over the iterations. The choice of the sequence η_t is more significant. In practice, it is tuned by a trial and error manner: given a validation data-set \mathcal{S}^* you may perform cross validation based on Error Function (3.1).
- (2) Re-run the SGD procedure several times and by using different settings (learning rate η_t , batch size m) and different seed $w^{(0)}$ (randomly chosen) each time. Possibly, by luck, at some trial, we will initialize the SGD process with a random seed $w^{(0)}$ producing a trace leading to a good local minimum w^* .
- (3) The output w_{SGD}^* returned by SGD is the best discovered w tested by using a performance measure (Error Function) using a validation set $\mathcal{S}^* = \{z_i^* = (x_i^*, y_i^*); i = 1, \dots, n^*\}$; Eg.

$$w_{\text{SGD}}^* = \arg \min_{w^{(t)}} \left(\text{EF} \left(w^{(t)} | \{z_i^*\} \right) \right).$$

5.1. Error backpropagation.

Note 61. The error back-propagation procedure is an efficient algorithm for the computation of the gradient $\nabla_w \ell(w, z)$ of the loss function $\ell(w, z)$ at some value of w and some example $z = (x, y)$ as required for the implementation of stochastic gradient based algorithm to train the FFNN.

Note 62. Error backpropagation assumes that $\ell(w, z)$ is differentiable at w for each value of z ; $\nabla_w \ell(w, z)$ exists –however extensions exist.

Notation 63. Let $V_t = \{v_{t,1}, \dots, v_{t,k_t}\}$ be the t -th layer of a NN and $k_t = |V_t|$ the number of neuron in layer $t = 1, \dots, T$.

Algorithm 64. (*Error back-propagation*)

Requires: The FFNN $(\mathcal{V}, \mathcal{E}, \sigma, w)$ with the values of the a weight vector $w \in \mathbb{R}^{|\mathcal{E}|}$, and example value $z = (x, y)$

Returns: $\nabla_w \ell(w, z) = \left(\frac{\partial}{\partial w_{t,i,j}} \ell(w, z); \forall t, i, j \right)$

Initialize:

$$\text{Set } w_{t+1,j,i} = \begin{cases} w(v_{t,i}, v_{t+1,j}) & \text{if } (v_{t,i}, v_{t+1,j}) \in \mathcal{E} \\ 0 & \text{if } (v_{t,i}, v_{t+1,j}) \notin \mathcal{E} \end{cases}$$

Forward pass:

(1) For $i = 1, \dots, d$

Set:

$$o_{0,i} = x_i$$

(2) For $t = 1, \dots, T$; For $i = 1, \dots, k_t$

Compute:

$$\alpha_{t,i} = \sum_{j=1}^{k_{t-1}} w_{t,i,j} o_{t-1,j}$$

Compute:

$$o_{t,i} = \sigma_t(\alpha_{t,i})$$

Backward pass:

(1) For $t = T$; For $i = 1, \dots, k_T$

Compute:

$$(5.1) \quad \delta_{T,i} = \frac{\partial \ell_T}{\partial \alpha_{T,i}} (\alpha_T)$$

where ℓ_T is the loss as a function of the vector of activations $\alpha_T = (\alpha_{T,1}, \dots, \alpha_{T,k_T})^\top$.

(2) For $t = T - 1, \dots, 1$; For $i = 1, \dots, k_t$

Compute:

$$(5.2) \quad \delta_{t,i} = \frac{d}{d \alpha_{t,i}} \sigma_t(\alpha_{t,i}) \sum_{j=1}^{k_{t+1}} w_{t+1,j,i} \delta_{t+1,j}$$

Output:

For $t = 1, \dots, T$; $i = 1, \dots, k_t$; $j = 1, \dots, k_{t-1}$

If edge $(v_{t-1,j}, v_{t,i}) \in \mathcal{E}$, set:

$$(5.3) \quad \frac{\partial \ell}{\partial w_{t,i,j}} (w, z) = \delta_{t,i} o_{t-1,j}$$

Note 65. In (5.1) and (5.2), $\{\delta_{t,i}\}$ are called errors.

Example 66. Consider the multi-output regression problem, assume a predictive rule $h_w : \mathbb{R}^d \rightarrow \mathbb{R}^q$ with $h_w = (h_{w,1}, \dots, h_{w,q})$ and

$$h_k(x) = \sigma_2 \left(\sum_{j=1}^c w_{2,k,j} \sigma_1 \left(\sum_{i=1}^d w_{1,j,i} x_i \right) \right)$$

with activation functions $\sigma_2(a) = a$, and $\sigma_1(a) = \tanh(a) = \frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)}$, and loss $\ell(w, z) = \frac{1}{2} \sum_{k=1}^q (h_k(x) - y_k)^2$ for example $z = (x, y)$. Perform the error back propagation steps to compute the elements of $\nabla_w \ell(w, z)$ at some w and z .

Solution.

Forward pass:

Set: $o_{0,i} = x_i$ for $i = 1, \dots, d$

Compute:

at $t = 1$: for $j = 1, \dots, c$

comp: $\alpha_{1,j} = \sum_{i=1}^d w_{1,j,i} x_i$

comp: $o_{1,j} = \tanh(\alpha_{1,j})$

at $t = 2$: for $k = 1, \dots, q$

comp: $\alpha_{2,k} = \sum_{j=1}^c w_{2,k,j} o_{1,j}$

comp: $o_{2,k} = \alpha_{2,k}$

get: $h_k = o_{2,k}$

Note that $\frac{d}{d\xi} \sigma_1(\xi) = 1 - (\sigma_1(\xi))^2$ and that $\frac{d}{d\xi} \sigma_2(\xi) = 1$.

Backward pass:

at $t = 2$: for $k = 1, \dots, q$

comp:

$$\delta_{2,k} = \frac{\partial}{\partial \alpha_{2,k}} \ell_T = \sum_{j=1}^q \frac{\partial \ell_T}{\partial o_{2,j}} (o_{2,j}) \frac{\partial o_{2,j}}{\partial \alpha_{2,k}} (\alpha_{2,k}) = \frac{\partial \ell_T}{\partial o_{2,k}} (o_{2,k}) \frac{\partial o_{2,k}}{\partial \alpha_{2,k}} (\alpha_{2,k}) = h_k - y_k$$

at $t = 1$: for $j = 1, \dots, c$

comp:

$$\begin{aligned} \delta_{1,j} &= \frac{d}{d\xi} \sigma_1(\xi) \Big|_{\xi=\alpha_{1,j}} \sum_{k=1}^q w_{2,k,j} \delta_{2,k} = \left(1 - \left(\underbrace{\sigma_1(\alpha_{1,j})}_{=o_{1,j}} \right)^2 \right) \sum_{k=1}^q w_{2,k,j} \delta_{2,k} \\ &= \left(1 - (o_{1,j})^2 \right) \sum_{k=1}^q w_{2,k,j} \delta_{2,k} \end{aligned}$$

Output:

$$\frac{\partial \ell(w, z)}{\partial w_{1,j,i}} = \delta_{1,j} x_i \text{ and } \frac{\partial \ell(w, z)}{\partial w_{2,k,j}} = \delta_{2,k} o_{1,j}$$

Note 67. We show that the Algorithm 64 is valid. The Forward pass is valid due to Algorithm 24. Now, we work on the Backward pass.

- Let $\alpha_t = (\alpha_{t,1}, \dots, \alpha_{t,k_t})^\top$ be the vector of the activations of the t -th layer of a NN.
- Let $\ell_t(\cdot)$ be the loss function $\ell(w, z)$ as a function of the vector of the activations α_t at t -th layer.

- It is

$$(5.4) \quad \frac{\partial \ell}{\partial w_{t,i,j}}(w, z) = \frac{\partial \ell_t}{\partial \alpha_{t,i}}(\alpha_t) \frac{\partial \alpha_{t,i}}{\partial w_{t,i,j}} = \delta_{t,i} \frac{\partial \alpha_{t,i}}{\partial w_{t,i,j}}$$

where

$$\delta_{t,i} = \frac{\partial \ell_t}{\partial \alpha_{t,i}}(\alpha_t)$$

for $i = 1, \dots, k_t$, $j = 1, \dots, k_{t-1}$ and $t = 1, \dots, T$.

- Regarding the second part of (5.4), it is

$$\frac{\partial \alpha_{t,i}}{\partial w_{t,i,j}} = \frac{\partial}{\partial w_{t,i,j}} \sum_{s=1}^{k_{t-1}} w_{t,i,s} o_{t-1,s} = o_{t-1,i}$$

for $t = 1, \dots, T$ and $i = 1, \dots, k_t$.

- I will try to find a way to compute δ_t given δ_{t+1} recursively from V_T to V_0 .

- For $t = T$, and $i = 1, \dots, k_T$, it is

$$\delta_{T,i} = \frac{\partial \ell_T}{\partial \alpha_{T,i}}(\alpha_T)$$

- For $t < T$, and $i = 1, \dots, k_t$, it is

$$\delta_{t,i} = \frac{\partial \ell_t}{\partial \alpha_{t,i}}(\alpha_t) = \sum_{j=1}^{k_{t+1}} \frac{\partial \ell_{t+1}}{\partial \alpha_{t+1,j}}(\alpha_{t+1}) \frac{\partial \alpha_{t+1,j}}{\partial \alpha_{t,i}} = \sum_{j=1}^{k_{t+1}} \delta_{t+1,j} \frac{\partial \alpha_{t+1,j}}{\partial \alpha_{t,i}}$$

and

$$\begin{aligned} \frac{\partial \alpha_{t+1,j}}{\partial \alpha_{t,i}} &= \frac{\partial}{\partial \alpha_{t,i}} \sum_{s=1}^{k_t} w_{t+1,j,s} o_{t,s} = \frac{\partial}{\partial \alpha_{t,i}} \sum_{s=1}^{k_t} w_{t+1,j,s} \sigma_{\textcolor{red}{t}}(\alpha_{\textcolor{red}{t},s}) \\ &= w_{t+1,j,i} \frac{d}{d \alpha_{\textcolor{red}{t},i}} \sigma_{\textcolor{red}{t}}(\alpha_{\textcolor{red}{t},i}) \end{aligned}$$

hence

$$\delta_{t,i} = \frac{d}{d \alpha_{\textcolor{red}{t},i}} \sigma_{\textcolor{red}{t}}(\alpha_{\textcolor{red}{t},i}) \sum_{j=1}^{k_{t+1}} \delta_{t+1,j} w_{t+1,j,i}$$

Note 68. Error back-propagation idea can also be used to efficiently compute the gradient $\nabla_x \ell(w, z = (x, y))$ of the loss function $\ell(w, z = (x, y))$ with respect to the inputs x . The idea is the same, use chain rule to find a recursive procedure from V_T to V_0 .

5.2. Preconditioning and computation of the Hessian.

Note 69. Recall (Handout 3, Algorithm 10), that SGD may be improved by using a suitable preconditioner $P_t > 0$ as

$$w^{(t+1)} = w^{(t)} - \eta_t P_t \nabla_w \ell(w^{(t)}, z^{(t)})$$

such a preconditioner can be the $P_t := [H_t + \epsilon I_d]^{-1}$ where H_t is the Hessian of $\ell(w^{(t)}, z^{(t)})$ and $\epsilon > 0$.

Note 70. The computation of the Hessian H_t can be done by using error propagation ideas for

$$[H_t]_{i,j} = \frac{\partial^2}{\partial w_i \partial w_j} f(w) \Big|_{w=w^{(t)}}$$

We will not go further to exact computations.

Note 71. The exact computation of the Hessian H_t of the loss $\ell(\cdot, z)$ in NN setting can be computationally expensive and hence approximations are often used (with hope to work well). One can implement the general purpose AdaGrad (Section 2.1, Handout 3). In what follows we present other alternatives tailored to the NN model.

Note 72. Consider the regression problem with predictive rule $h : \mathbb{R}^d \rightarrow \mathbb{R}$ with $h_w(x) = o_T(x) = \alpha_T(x)$, and loss function $\ell(w, z = (x, y)) = \frac{1}{2}(h_w(x) - y)^2$. Then the Hessian of ℓ is

$$\begin{aligned} (5.5) \quad H &= \frac{d(\nabla_w \ell)}{dw} = \frac{d}{dw} \left(\nabla_w h_w(x) (h_w(x) - y)^\top \right) \\ &= \left(\frac{d}{dw} \nabla_w h_w(x) \right) (h_w(x) - y)^\top + \nabla_w h_w(x) (\nabla_w h_w(x))^\top \end{aligned}$$

We can expect that $h_w(x) \approx y$ provided that the network is well trained due to the global approximation ability of the FFNN. Yet, we can expect that $h_w = E(y)$ provided that we train the network under the quadratic loss and because $\arg \min_h E_{y \sim g}(h - y)^2 = E_{y \sim g}(y)$. Hence a reasonable approximation can be

$$\begin{aligned} H &\approx \nabla_w h_w(x) (\nabla_w h_w(x))^\top \\ &\stackrel{\sigma_T(a)=a}{=} \nabla_w \alpha_T(x) (\nabla_w \alpha_T(x))^\top \end{aligned}$$

Consequently, the approximation of the Hessian of the Error function $EF(w | \{z_i\}) = \sum_{i=1}^n \ell(w, z_i = (x_i, y_i))$ is

$$(5.6) \quad H_n = \sum_{i=1}^n \frac{d(\nabla_w EF)}{dw} \approx \sum_{i=1}^n \nabla_w h_w(x_i) (\nabla_w h_w(x_i))^\top$$

Note 73. (Cont. Note 72) To efficiently compute the inverse H_n^{-1} of (5.6) I can utilize Woodbury identity

$$(5.7) \quad (M + vv^\top)^{-1} = M^{-1} - \frac{(M^{-1}v)(v^\top M^{-1})}{1 + v^\top M^{-1}v}$$

No need to
memorize
Woodbury
identity

offering a way to avoid the computational expensive task of directly inverting the high dimensional H_n . Given $v_i = \nabla_w h_w(x)$, it is

$$(5.8) \quad \begin{aligned} (H_n)^{-1} &= \left(\sum_{i=1}^n v_i v_i^\top \right)^{-1} = \left(\sum_{i=1}^{n-1} v_i v_i^\top + v_n v_n^\top \right)^{-1} = \left(H_{n-1} + v_n v_n^\top \right)^{-1} \\ &= H_{n-1}^{-1} - \frac{(H_{n-1}^{-1} v_n)(v_n^\top H_{n-1}^{-1})}{1 + v_n^\top H_{n-1}^{-1} v_n} \end{aligned}$$

In practice I start with $H_0 = \epsilon I$ with $\epsilon > 0$ small, and iterate (5.8).

Note 74. Likewise and by modifying (5.5), one can compute the corresponding approximations for the classification problems or problems with different loss functions.

6. BAYESIAN ARTIFICIAL NEURAL NETWORKS

Note 75. Consider a feed-forward Neural Network with $(\mathcal{V}, \mathcal{E}, \sigma, w)$. Assume the architecture $(\mathcal{V}, \mathcal{E}, \sigma)$ of a neural network is fixed/known, and interest lies in learning the weights $\{w_{t,i,j}\}$.

Problem 76. Assume there is available a training set of examples (data-set) $\mathcal{S} = \{z_i = (x_i, y_i); i = 1, \dots, n\}$ with $x_i \in \mathcal{X} = \mathbb{R}^{|V_0|}$ and $y_i \in \mathcal{Y}$. The Bayesian NN model is then

$$\begin{cases} y_i | x_i, w & \stackrel{\text{ind}}{\sim} f(y_i | x_i, w), i = 1, \dots, n \text{ (sampling distribution)} \\ w & \sim f(w) \text{ (prior distribution)} \end{cases}$$

where the sampling distribution is specified either according to the experimental design (how $\{z_i\}$) are collected or based on subjective judgments, while the prior distribution is specified based on subjective manner.

Note 77. In most of the real applications, it is difficult (almost impossible) to specify the sampling distribution based on the experimental design or judgments, and almost impossible to specify the prior of the weights based on subjective judgments.

Note 78. One could specify the sampling distribution based on the loss function $\ell(h_w(x), z = (x, y))$ as

$$(6.1) \quad f(y_i | x_i, w) \propto \exp(-\ell(h_w(x_i), (x_i, y_i))),$$

based on the argument that sampling distribution in the posterior distribution contraindicates how far the theoretical model $h_w(x_i)$ (as casted in a NN) is from the corresponding observation y_i via the likelihood.

Note 79. The prior density of the weights may be specified based on some shrinkage term $J(w; \lambda)$ (Note 55) as

$$(6.2) \quad f(w) \propto \exp(-J(w; \lambda))$$

based on the argument that I often model the predictive rule $h_w(\cdot)$ with an over-parametrized NN (with many weights, more than needed) and many of them should be shrunk to zero. Also as we

see in Note 88, careless training of NN tends to produce over-fitted NN with large weights (in abs values).

Note 80. Regarding prior it is often, $w \sim N(0, I\lambda^{-1})$ for some small λ controlling prior uncertainty about weights, i.e. 50% chances for the weight to be above or below zero.

Note 81. The corresponding posterior (according to the Bayes theorem) is

$$f(w|\{z_i\}) = \frac{\prod_{i=1}^n f(z_i|w) f(w)}{\int \prod_{i=1}^n f(z_i|w') f(w') dw'} \propto \prod_{i=1}^n f(z_i|w) f(w)$$

and given (6.1) and (6.2), in log scale I get

$$\log(f(w|\{z_i\})) = - \sum_{i=1}^n \ell(h_w(x_i), (x_i, y_i)) - J(w; \lambda) + \text{const}$$

that resembles the EF with a shrinkage term in classical learning.

Example 82. (Regression problem) Sampling distribution can be specified as in (6.1) based on (Euclidean) loss

$$\ell(w, z = (x, y)) = \frac{\beta}{2} (h_w(x) - y)^2$$

which (based on (3.4)) can build a density

$$f(y_i|x_i, w) \propto \exp(-\ell(h_w(x_i), (x_i, y_i))) = -\frac{\beta}{2} (h_w(x_i) - y_i)^2,$$

This implies a sampling distribution

$$(6.3) \quad y_i|x_i, w \sim N(\mu_i, \beta^{-1}), \text{ where } \mu_i = h_w(x_i)$$

for some fixed $\beta > 0$ based on (3.2). Prior can be specified with density (6.2) where

$$J(w; \lambda) = \frac{\lambda}{2} \|w\|_2^2$$

which is based on the Ridge shrinkage term. This implies a prior

$$(6.4) \quad w \sim N(0, I\lambda^{-1})$$

The resulted posterior density is

$$(6.5) \quad f(w|\{z_i\}) \propto \exp\left(-\frac{\beta}{2} \sum_{i=1}^n (h_w(x_i) - y_i)^2 - \frac{\lambda}{2} \|w\|_2^2\right)$$

Example 83. (Binary classification problem) Sampling distribution can be specified as in (6.1) based on (cross-entropy) loss

$$\ell(w, (x, y)) = y \log(h_w(x)) + (1 - y)(1 - \log(h_w(x)))$$

which (based on (3.4)) can build a density

$$\begin{aligned} f(y_i|x_i, w) &\propto \exp(-\ell(h_w(x_i), (x_i, y_i))) \\ &= -y_i \log(h_w(x_i)) - (1 - y_i)(1 - \log(h_w(x_i))), \end{aligned}$$

This implies a sampling distribution

$$y_i|x_i, w \sim \text{Bernoulli}(p_i), \text{ where } p_i = h_w(x_i)$$

To encourage sparsity, prior can be specified as in (6.2) by using the LASSO shrinkage term

$$J(w; \lambda) = \lambda \|w\|_1;$$

The resulted posterior density is

$$f(w|\{z_i\}) \propto \exp \left(- \sum_{i=1}^n [y_i \log(h_w(x_i)) + (1-y_i)(1-\log(h_w(x_i)))] - \lambda \|w\|_1 \right)$$

Note 84. Sampling from the posterior can be performed via SGLD due to the high-dimensionality in the weights w and the big size of the training data set. Recall the recursion of the SGLD with batch size m is

$$(6.6) \quad w^{(t+1)} = w^{(t)} + \eta_t \left(\frac{n}{m} \sum_{j \in \mathcal{J}^{(t)}} \nabla_w \log \left(f(z_j^{(t)}|w^{(t)}) \right) + \nabla_w \log \left(f(w^{(t)}) \right) \right) + \sqrt{2\eta_t \tau} \epsilon_t,$$

for $\epsilon_t \sim N(0, 1)$. (See Algorithm 19 in Lecture notes 7: Bayesian Learning via Stochastic gradient and Stochastic gradient Langevin dynamics).

Note 85. Given sample values $\{w^{(t)}\}_{t=1}^T$ produced from SGLD recursions (6.6), learning of any function $h_w(x)$ of the w 's can be performed via

- (1) Monte Carlo estimation namely averaging the samples values $\{w^{(t)}\}_{t=1}^T$ as

$$\widehat{h}_w(x) = \frac{1}{T} \sum_{t=1}^T h_{w^{(t)}}(x), \quad (\text{Monte Carlo estimator})$$

- (2) Maximum-A-posteriori (MAP) estimation namely find the best \hat{w}^* among the samples values $\{w^{(t)}\}_{t=1}^T$ that maximizes the posterior i.e

$$\hat{w}^* = \arg \max_{w \in \{w^{(t)}\}_{t=1}^T} (\log f(\{z_i\}|w) + \log f(w))$$

and compute

$$\widehat{h}_w(x) = h_{\hat{w}^*}(x), \quad (\text{MAP estimator})$$

7. THEORETICAL ASPECTS

Note 86. We will not go this direction here. For the interested student, I suggest for starters, Ch 5 from “Ripley, B. D. (2007). Pattern recognition and neural networks. Cambridge university press.”; and for advanced Ch 30 from “Devroye, L., Györfi, L., & Lugosi, G. (2013). A probabilistic theory of pattern recognition (Vol. 31). Springer Science & Business Media.”

8. COMMENTS, GUIDELINES, AND DISCUSSIONS

8.1. Over-fitting issues.

Note 87. Neural Networks can be highly “over-parametrized”. E.g. they may consist of a large number of layers each of them having a large number of neurons with purpose to represent each feature/characteristic of the pattern to be recognized –this increases the number of unknown weights to be learned. Careless training may produce NN models with bad generalization predictive properties.

Note 88. Training of NN models corresponds to an iterative reduction of the Error Function defined on the training data-set. On the other hand it has been observed that the Error Function defined on the validation data-set (independent to the training data set) often shows a decrease at first (while NN is being tuned), followed by an increase (while the NN starts to overfit). The reason is because the NN training is a high-dimensional and big-data problem; hence during the NN training against the EF defined on the training data, the NN memorises too much info from the specific training dataset, overfits to it, and presents bad generalization performance. This overfitting is often associated to the production of weight values larger than needed for the NN learning rule to have a good generalization performance. It is desirable to avoid this over-fitting with purpose to obtain a NN based learning rule with good generalization performance.

Note 89. **Early stopping** is a way of limiting the effective network complexity by halting training before a minimum of the training error has been reached. During training of NN against the MM defined on the training set, we monitor the NN performance against the EF defined on the validation data-set, and stop the training procedure just before the EF defined on the validation data set starts increasing (aka before overfitting) to prevent the model memorizing too much information about the training set.

Note 90. **Regularization** (Section 4.2) can be used to address the above issue by using a shrinkage term preventing the weights to grow too much away from zero.

8.2. Non-identifiabilities.

Note 91. NN Learning Problems (E.g, Problems 50, 51, 54, and 76) are often non-identifiable with respect to weights $\{w_{t,i,j}\}$.

Example 92. Consider the FFNN in Figure 8.2 with 1 hidden layer of M neurons, activation function $\sigma_1(a) = \tanh(a)$ and full connectivity in both layers. If we flip the sign of all of the w 's feeding into a particular hidden unit (for a given input pattern) then the sign of the activation of the hidden unit will be reversed because $\tanh(-a) = -\tanh(a)$. This transformation can be compensated by changing the sign of all of the w 's leading out of that hidden unit. By changing the signs of a particular group of w 's, the input-output mapping function represented by the FFNN is unchanged, and so there are two different weight vectors resulting the same mapping function. I can do M such flips (there are M hidden neurons) leading to 2^M equivalent parameter settings. Similarly, we can permute the labels of the neurons in hidden layer without changing the loss function; there are $M!$ such permutations. Hence the equivalent parameter settings in total are $2^M M!$. –This is not harmful in training or predictive rul computation as we just need to find one such parameter setting.

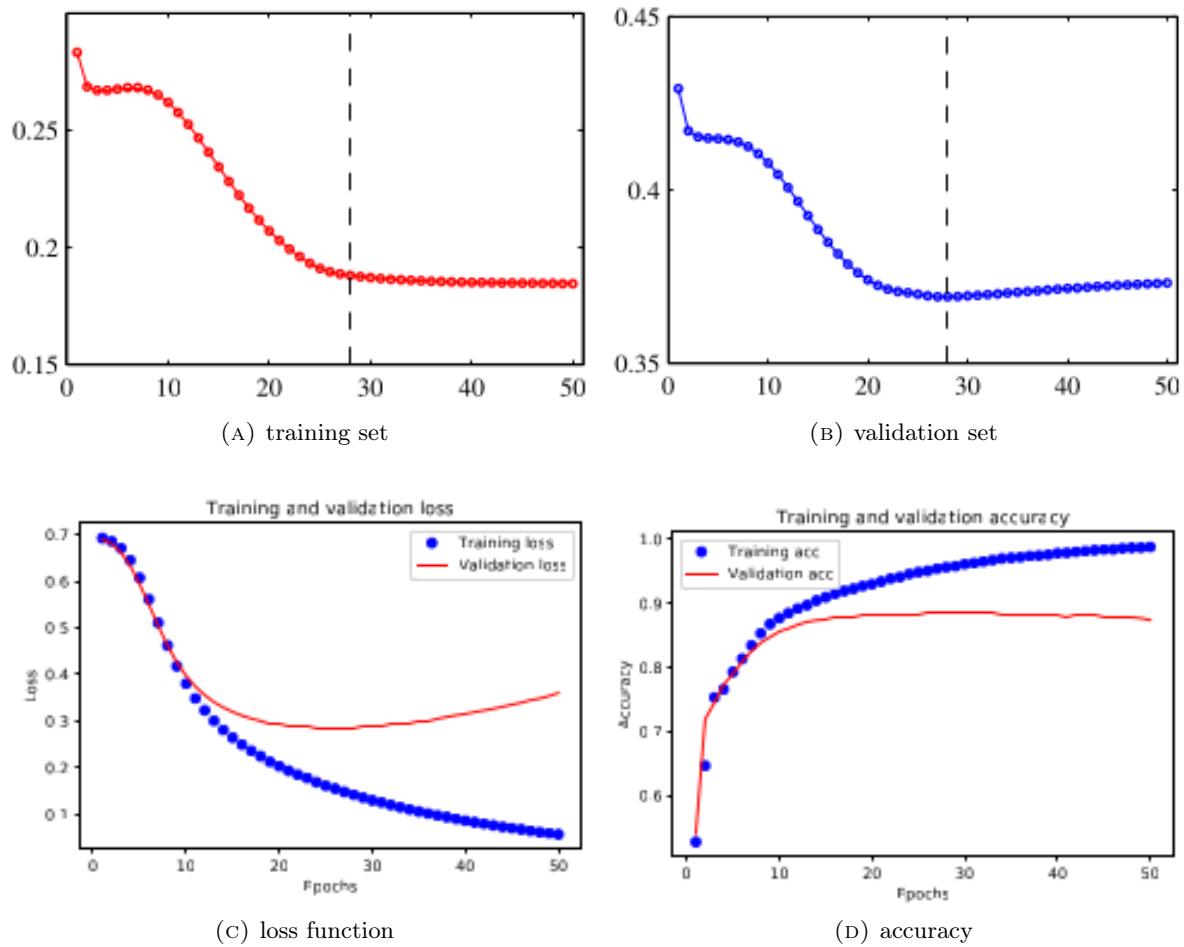


FIGURE 8.1. Behavior of the Error Function wrt the iterations, against a training set and against a validation set.

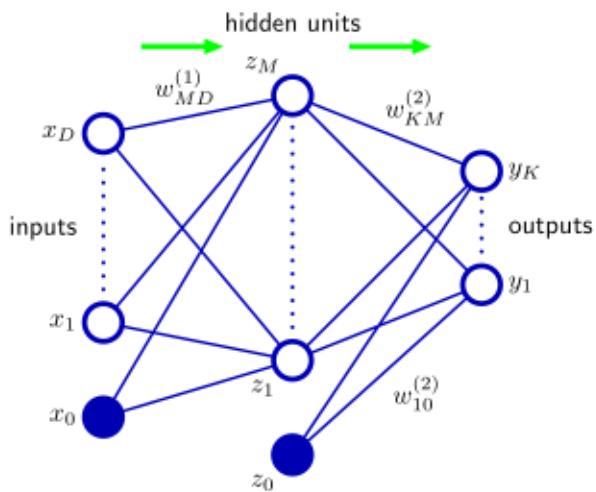


FIGURE 8.2. A FFNN

8.3. Non-convexity.

Note 93. Learning problems with NN are non-convex, hence the loss function to be minimized has many local minima. The consequence is that the SGD/SGLD learning algorithms may be trapped in a local optima and never reach any of the global ones. The number of local minima is exaggerated due to the symmetries discussed in Note 91. Due to the large number of data (big-data) used to train the NN (in real life) such local minima become more shallow while the global minima more picky. Due to this and the stochastic nature of SGD/SGLD algorithms, SGD/SGLD may be able (by chance) to escape from such local minima/maxima and reach the global ones.

Note 94. To address local optima issue, run SGD/SGLD learning algorithms multiple times by initializing them with different seeds each time.

8.4. Vanishing and exploding gradient.

Note 95. In the Error back-propagation Algorithm 64, for the computation of $\left\{ \frac{\partial \ell}{\partial w_{t,i,j}}(w, z) \right\}$ crucial quantities are the δ -derivatives the propagation of which follows the recursive rule (5.2), i.e.

$$\delta_{t,i} = \sum_{j=1}^{k_{t+1}} w_{t+1,j,i} \delta_{t+1,j} \sigma'_t(\alpha_{t,i}); t = T-1, \dots, 1; i = 1, \dots, k_t$$

δ_t -derivatives depend on the respective δ_{t+1} -derivatives, weights, and derivatives of the sigmoid in a multiplicative manner. Keeping on

$$\delta_{t,i} = \left[\sum_{j=1}^{k_{t+1}} \left(\sum_{s=1}^{k_{t+2}} \delta_{t+1,s} w_{t+2,s,j} \right) \sigma'_{t+1}(\alpha_{t+1,j}) w_{t+1,j,i} \right] \sigma'_t(\alpha_{t,i})$$

we see that as the backpropagation algorithm flows backwards, the derivatives of the sigmoids as well as the weights are multiplied and the number of the involved products grows. This may cause the two problems: vanishing or exploding gradient.

Note 96. This is more powerful in lower layers.

Note 97. This exaggerates in deeper structures of NN.

8.4.1. Vanishing gradient.

Note 98. As $\sigma'_t(\cdot)$ are often smaller than 1 if the corresponding weights are not large enough, the gradients of the loss wrt the weights of the lower layers can take vanishingly small values and make training via SGD too slow.

Note 99. A remedy to mitigate vanishing gradient is to use rectified linear unit (ReLU) activation functions

$$\sigma_t(\alpha) = \max(0, \alpha) + \xi \min(0, \alpha)$$

where $\xi \in \mathbb{R}$ is a user specified parameter. For instance: if $\xi = 0$, we observe that $\sigma'_t(\alpha) = 1$ when $\alpha > 0$ hence it does not suffer from saturation. However the training may freeze when $\alpha < 0$ since $\sigma'_t(\alpha) = 0$, and hence we often consider some $|\xi| \approx$ small depending on the application.

8.4.2. Exploding gradient.

Note 100. If the values of the weights get very large values, this can lead to an explosion of the gradient estimates. This can affect the learning by pushing the estimates of the weights to the wrong region in the parameter space when SGD is used.

Note 101. Remedies of this are clipping the gradient in SGD (Note 35 in Handout 7: Bayesian Learning via Stochastic gradient and Stochastic gradient Langevin dynamics), the use of shrinkage terms in the learning problem (Section 4.2), or the use of Bayesian framework with appropriate priors (Section 6).

8.5. Skip-layers.

Note 102. The general definition of feed forward Neural Network allows “skip-layer” connections that is the directed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ in NN architecture includes edges which may not necessarily connect neurons between consecutive layers, namely $(v_{t,j}, v_{t+k,i}) \in \mathcal{E}$ for some $k > 1$.

Example 103. For instance, a fully connected FFNN with 1 hidden layer is

$$(8.1) \quad h_k(x) = \sigma_{(2)} \left(\underbrace{\sum_{\forall i} w_{(1),k,i}^{\text{skip}} x_i}_{\text{skip layer terms}} + \sum_{\forall j} w_{(2),k,j} \sigma_{(1)} \left(\sum_{\forall i} w_{(1),j,i} x_i \right) \right)$$

Compared to (2.4) in Example 33 and (2.4), (8.1) links the input to the output.

Note 104. In skip-layer FFNN cases, the error back propagation (Section 5.1) has to be adjusted properly.

Note 105. FFNN are mainly used without skip-layers because theory (Section 7) states that FFNN are global approximators even without skip layers, as well as because of computational inconvenience and modeling parsimony.

APPENDIX A. ABOUT GRAPHS

Definition 106. A directed graph is an ordered pair $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ comprising

- a set of nodes \mathcal{V} (or vertices), where nodes are abstract objects, and
- a set of edges $\mathcal{E} = \{(v, u) | v \in \mathcal{E}, u \in \mathcal{E}, v \neq u\}$ (or directed edges, directed links, arrows) which are ordered pairs of vertices (that is, an edge is associated with two distinct vertices).

Definition 107. An edge-weighted graph or a network is a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ equipped with a weight function $w : \mathcal{E} \rightarrow \mathbb{R}$ that assigns a number (the weight) to each edge $e \in \mathcal{E}$.

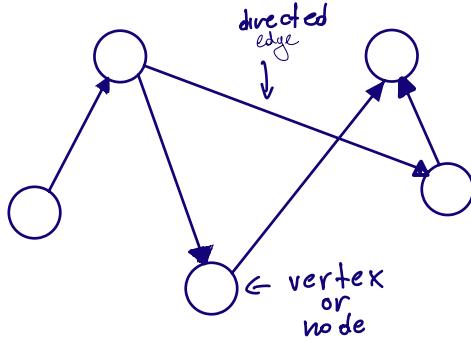


FIGURE A.1. A directed graph

APPENDIX B. ABOUT PARTIAL DERIVATIVES

Note 108. Consider a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$.

Definition 109. The partial derivative of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at the point $a = (a_1, \dots, a_n) \in U \subseteq \mathbb{R}^n$ with respect to the i -th variable is denoted as

$$\frac{\partial f}{\partial x_i}(a) \text{ or } \left. \frac{\partial f}{\partial x_i}(x) \right|_{x=a}$$

and defined as

$$\begin{aligned} \left. \frac{\partial f}{\partial x_i}(x) \right|_{x=a} &= \lim_{h \rightarrow 0} \frac{f(a_1, \dots, a_{i-1}, a_i + h, a_{i+1}, \dots, a_n) - f(a_1, \dots, a_{i-1}, a_i, a_{i+1}, \dots, a_n)}{h} \\ &= \lim_{h \rightarrow 0} \frac{f(a + he_i) - f(a)}{h} \end{aligned}$$

where e_i is a $0 - 1$ vector with only one ace in the i -th location.

Remark 110. Essentially, Definition 109 says that the partial derivative of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at the point $a = (a_1, \dots, a_n) \in U \subseteq \mathbb{R}^n$ with respect to the i -th variable is

$$\frac{\partial f}{\partial x_i}(a) = \left. \frac{dg}{dh}(h) \right|_{h=0}$$

the derivative of function $g(h) := f(a_1, \dots, a_{i-1}, a_i + h, a_{i+1}, \dots, a_n)$ at value 0.

Example 111. Consider a function f with $f(x_1, x_2) = x_1^2 + x_1x_2^3$. Compute its partial derivatives at $a = (2, 3)^\top$; i.e. $\frac{\partial f}{\partial x_1}(a)$ and $\frac{\partial f}{\partial x_2}(a)$.

Solution. It is

$$\begin{aligned}\frac{\partial f}{\partial x_1}(a) &= \frac{\partial f}{\partial x_1}(x)\Big|_{x=a} = \frac{d}{dx_1}(x_1^2 + x_1x_2^3)\Big|_{x=a} = 2x_1 + x_2^3|_{x=a} \\ &= 2a_1 + a_2^3 = 4 + 27 = 31 \\ \frac{\partial f}{\partial x_2}(a) &= \frac{\partial f}{\partial x_2}(x)\Big|_{x=a} = \frac{d}{dx_2}(x_1^2 + x_1x_2^3)\Big|_{x=a} = 2x_1x_2^2|_{x=a} \\ &= 2a_1a_2^2 = 4 + 27 = 72\end{aligned}$$

Lecture notes 12: Gaussian process regression

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim. To introduce the Gaussian process regression as a kernel method.

Reading list & references:

- (1) Bishop, C. M. (2006). Pattern recognition and machine learning (Vol. 4, No. 4, p. 738). New York: Springer.
 - Ch. 6.4 Gaussian process
- (2) Rasmussen, C. E., & Williams, C. K. (2006). Gaussian processes for machine learning (Vol. 1, p. 159). Cambridge, MA: MIT press.
 - Chapter 2, Regression (supplementary)
- (3) Roustant, O., Ginsbourger, D., & Deville, Y. (2012). DiceKriging, DiceOptim: Two R packages for the analysis of computer experiments by kriging-based metamodeling and optimization. *Journal of statistical software*, 51, 1-55.
 - Supplementary material related to the the implementation of GP in R computing environment.

1. INTRO AND MOTIVATION

Note 1. As motivation for the Gaussian process regression, we “Kernelize” the standard Bayesian normal linear regression in the machine learning framework.

Note 2. Consider the predictive rule $h(x) = \eta(x)$, and cast it in a linear form $\eta(x) = (\psi(x))^\top w$ where $\psi(x) = (\psi_1(x), \dots, \psi_d(x))$ is a vector of basis functions mapping the input space \mathcal{X} into a feature space \mathcal{F} . Assume there is available a set of observables $\{z_i = (x_i, y_i)\}_{i=1}^n$. We associate the learning problem with the Bayesian linear regression model

(1.1)

$$\begin{cases} y_i | \psi(x_i), w, \sigma^2 \stackrel{\text{ind}}{\sim} N(\eta(x_i), \sigma^2), \quad i = 1, \dots, n \\ \eta(\cdot) = (\psi(\cdot))^\top w \\ w \sim N(\mu_0, V_0) \end{cases} \quad \text{equiv.} \quad \begin{cases} y | \eta, \sigma^2 \sim N(\eta, I\sigma^2) & (\text{sampl. distr.}) \\ \eta = \Psi w & (\text{linear model restr.}) \\ w \sim N(\mu_0, V_0) & (\text{prior}) \end{cases}$$

where $[\Psi]_{i,j} = \psi_j(x_i)$.

Note 3. The marginal likelihood is

$$(1.2) \quad f(y) = N\left(y | \Psi^\top \mu_0, \Psi V_0 \Psi^\top + I\sigma^2\right)$$

where $N(y|\mu, \Sigma)$ denotes the pdf of the Normal distribution with mean μ , and covariance matrix Σ .

Note 4. ¹The predictive distribution of a new outcome y_* at a new input x_* given the observables $\{z_i = (x_i, y_i)\}_{i=1}^n$ is

$$f(y_*|x_*, \{(x_i, y_i)\}) = N(\mu_*(x_*), \sigma_*^2(x_*))$$

with

$$(1.3) \quad \mu_*(x_*) = \psi(x_*)^\top \mu_0 + \frac{1}{\sigma^2} \underbrace{\psi(x_*)^\top V \Psi}_{K(x_*, X)} \left(\underbrace{\Psi^\top V \Psi}_{K(X, X)} + \sigma^2 \right)^{-1} (\Psi^\top \mu_0 - y)$$

$$(1.4) \quad \sigma_*^2(x_*) = \left(\underbrace{\psi(x_*)^\top V \psi(x_*)}_{=K(x_*, x_*)} + \sigma^2 \right) - \underbrace{\psi(x_*)^\top V \Psi}_{=K(x_*, X)} \left(\underbrace{\Psi^\top V \Psi}_{=K(X, X)} + \sigma^2 \right)^{-1} \underbrace{(\psi(x_*)^\top V \Psi)^\top}_{=K(X, x_*)}$$

according to Proposition 44.

Note 5. In the prior part of (1.1), let's assume $\mu_0 = 0$ (arguably) denoting complete ignorance whether $\eta(\cdot)$ is positive or negative. By applying Kernel trick in (1.3) and (1.4), the feature space always enters in the form inner products. In fact we can define a kernel $K(x, x') = \langle L\psi(x), L\psi(x') \rangle = \psi(x)^\top V \psi(x')$ where L is such that $V = L^\top L$, in terms of Section 4 in Handout 7: Kernel methods. We can denote $K(x_*, X) = \psi(x_*)^\top V \Psi$, and $K(x_*, x_*) = \psi(x_*)^\top V \psi(x_*)$.

2. THE GAUSSIAN PROCESS REGRESSION MODEL

Definition 6. Gaussian process (GP) is a collection of random variables $\{f(x); x \in \mathcal{X}\}$, indexed by label x , where any finite collection of those variables has a multivariate normal distribution. It is fully specified by its mean and covariance functions. It is denoted as

$$f(\cdot) \sim GP(\mu(\cdot), C(\cdot, \cdot))$$

with mean

$$\mu(x) := E(f(x)), x \in \mathcal{X}$$

and covariance function

$$C(x, x') := \text{Cov}(f(x), f(x')), x, x' \in \mathcal{X}$$

Note 7. Essentially, GP is a distribution defined over functions.

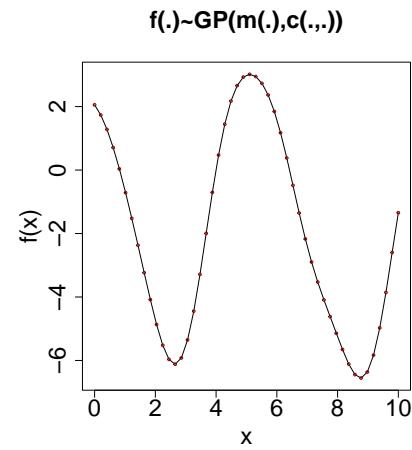
Example 8. One way to simulate a GP realization $f(\cdot) \sim GP(\mu(\cdot), C(\cdot, \cdot))$, with $\mu(x) = 1 + 0.01x$, and $C(x, x') = \exp\left(-\frac{1}{2} \frac{|x-x'|^2}{1.5}\right)$ in the range $x \in [0, 10]$ is given in the following R code.

¹No need to memorize the formulas in (1.2), (1.3), and (1.4). The material in Notes 2, 3, and 4 is given as a motivation for the Gaussian process regression.

```

rm(list=ls())
set.seed(99)
n <- 50 #descretize the problem
x <- seq(from = 0, to = 10, length = n)
mu_x <- matrix(1,n)
cov_x_x <- matrix(nrow = n, ncol = n)
for (i in 1:n) {
  mu_x[i] <- 1+0.01*x[i]
  for (j in 1:n) {
    cov_x_x[i,j]<- 10*exp(-0.5*(x[i]-x[j])^2/1.5)
  }
}
f <- rmvnorm(n = 1, mean = mu_x , sigma = cov_x_x)
plot(x,f,type="l",
  main="f(.)~GP(m(.),c(..))",
  ylab="f(x)",xlab="x",
  cex.axis=2, cex.lab=2, cex.main=2, cex.sub=2)
points(x, f, cex = .5, col = "dark red")

```



Note 9. Consider a function $\eta : \mathcal{X} \rightarrow \mathbb{R}$ with $\eta(x) = \langle \psi(x), w \rangle$ where $\psi(x)$ is a vector of known basis (feature) functions mapping from the input space \mathcal{X} to the feature space \mathcal{F} , and $w \in \mathbb{R}^d$ is an unknown vector a priori following a normal distribution $w \sim N(0, V)$, where the prior mean is set to zero denoting complete uncertainty about the sign of w 's. Then the marginal $\eta(\cdot)$ follows a Normal distribution as a linear transformation of Normal variates with mean $E(\eta(x)) = 0$ and covariance $\text{Cov}(\eta(x), \eta(x')) = \psi(x)^\top V \psi(x')$ for any $x, x' \in \mathcal{X}$. Based on the Kernel trick and Definition 6, we can equivalently specify $\eta(\cdot) \sim GP(0, C(\cdot, \cdot))$ for some corresponding kernel / covariance function $C(x, x') = \psi(x)^\top V \psi(x')$.

Note 10. We introduce the concept of Gaussian process regression in the machine learning framework below.

Note 11. Consider the predictive rule $h(x) = \eta(x)$, and assume that $\eta : \mathcal{X} \rightarrow \mathbb{R}$ with unknown formula (possibly up to a set of properties, we will discuss this later) and $\mathcal{X} \subseteq \mathbb{R}^d$.

Example 12. Figure 2.1a shows a function $\eta(\cdot)$. We pretend that we do not know $\eta(\cdot)$ but we wish to recover it. To recover $\eta(x)$, we collect training data set as in Figure 2.1b.

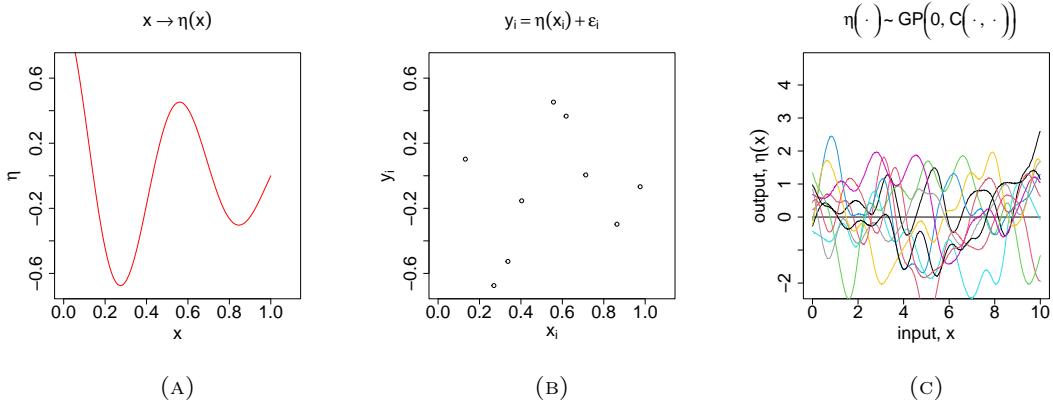


FIGURE 2.1. A toy example. (2.1a) shows the true $\eta(x)$. Fig 2.1b shows the training sample $\{z_i = (x_i, y_i)\}$ s.t. $y_i = \eta(x_i) + \epsilon_i$. Fig 2.1c shows several realisations of the prior $\eta(\cdot) \sim \text{GP}(0, C(\cdot, \cdot))$. R code for the plots/running example is available from https://github.com/georgios-stats/Machine_Learning_and_Neural_Networks_III_Epiphanie_2024/tree/main/Lecture_handouts/code/10.Gaussian_process_regression/plots.R

Note 13. For training purposes, assume there is available a set of observables $\{z_i = (x_i, y_i)\}_{i=1}^n$ whose sampling distribution is such that

$$(2.1) \quad \begin{aligned} y_i &= \eta(x_i) + \epsilon_i, \quad \epsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2), \quad i = 1, \dots, n \\ &\text{or equivalently} \\ y_i | \eta(\cdot), \sigma^2 &\stackrel{\text{iid}}{\sim} N(\eta(x_i), \sigma^2), \quad i = 1, \dots, n \end{aligned}$$

for some unknown $\sigma^2 > 0$. Statistical model (2.1) can result by considering a quadratic loss $\ell(h, z = (x, y)) = \frac{1}{\sigma^2} (h(x) - y)^2$ and sampling distribution with pdf

$$\text{pr}(y | \{x_i, y_i\}) \propto \exp\left(-\sum_{i=1}^n \ell(h(x_i), (x_i, y_i))\right).$$

As $\eta(x)$ is assumed to be unknown, according to the Bayesian paradigm and by taking advantage of Note 9, we assign a GP prior on $\eta(\cdot)$

$$(2.2) \quad \eta(\cdot) \sim \text{GP}(\mu(\cdot|\beta), C(\cdot, \cdot|\phi))$$

where μ is parametrized by unknown β (e.g. $\mu(x|\beta) = x^\top \beta$), and C is parametrized by unknown ϕ (e.g. $C(x, x'|\phi) = \exp\left(-\frac{1}{2\phi} \|x - x'\|_2^2\right)$; radial/Gaussian kernel). Summing up, the Bayesian model

$$\begin{cases} y_i | \eta(\cdot), \sigma^2 & \stackrel{\text{iid}}{\sim} N(\eta(x_i), \sigma^2), \quad i = 1, \dots, n \\ \eta(\cdot) | \beta, \phi & \sim \text{GP}(\mu(\cdot|\beta), C(\cdot, \cdot|\phi)) \end{cases}$$

up to some unknown tuning parameters σ^2 , β , and $\phi > 0$.

Notation 14. From now on, to ease the notation, x , σ^2 , β , and ϕ are suppressed from the conditioning; e.g. we use $C(\cdot, \cdot)$ instead of $C(\cdot, \cdot | \phi)$.

Note 15. A realization of (2.2) (with ref to Fig 2.1c) can be simulated by setting a finite vector $x := (x_1, \dots, x_m)$, computing μ such that $[\mu]_i = \mu(x_i)$ and C such that $[C]_{i,j} = C(x_i, x_j)$ and drawing an m dimensional vector η from the multivariate Normal distribution $N(\mu, C)$. See Fig 2.1c, and the R code provided in the caption.

Note 16. Consider $\eta_* = \eta(X_*)$ where $X_* = (x_{*,1}, x_{*,2}, \dots, x_{*,m})^\top$ is a vector of new inputs of any length $m > 0$. The joint distribution of $(\eta_*, y)^\top$ is

$$(2.3) \quad \begin{pmatrix} \eta_* \\ y \end{pmatrix} \sim N \left(\begin{pmatrix} \mu(X_*) \\ \mu(X) \end{pmatrix}, \begin{pmatrix} C(X_*, X_*) & C(X_*, X) \\ C(X, X_*) & C(X, X) + I\sigma^2 \end{pmatrix} \right)$$

where $C(X, X_*)$ is a Gram matrix over X and X_* such as $[C(X, X_*)]_{i,j} = C(x_i, x_{*,j})$.

Note 17. The conditional distribution of $\eta_* = \eta(X_*)$ given the training sample $\{z_i = (x_i, y_i)\}$, as results from 2.3 (Proposition 44),

$$(2.4) \quad \eta_* | y \sim N(\mu_*(X_*), C_*(X_*, X_*))$$

is a normal distribution, with mean

$$(2.5) \quad \mu_*(X_*) = E(\eta_* | y) = \mu(X_*) + C(X_*, X) (C(X, X) + I\sigma^2)^{-1} (y - \mu(X))$$

at X_* and with covariance function

$$(2.6) \quad C_*(X_*, X_*) = \text{Cov}(\eta_* | y) = C(X_*, X_*) - C(X_*, X) (C(X, X) + I\sigma^2)^{-1} (C(X_*, X))^\top$$

Note 18. Comparing (2.6) with (1.4), the extra $+\sigma^2$ in the first term of the right-hand-side part is because in Note 17 we predict (and are interested in) the “underline pattern only” $\eta_* = \eta(x_*)$ while in Note we predict the “experimental outcome” $y_* = \eta_* + \epsilon$, $\epsilon \sim N(0, \sigma^2)$.

Note 19. Because X_* is of any finite length, and the derivations in Note 17, by definition of GP, the predictive distribution of $\eta(\cdot)$ given the data $\{z_i = (x_i, y_i)\}$ is the Gaussian process

$$(2.7) \quad \eta(\cdot) | \{(x_i, y_i)\} \sim GP(\mu_*(\cdot), C_*(\cdot, \cdot))$$

with mean function and covariance function

$$(2.8) \quad \mu_*(x_*) = \mu(x_*) + C(x_*, X) (C(X, X) + I\sigma^2)^{-1} (y - \mu(X))$$

$$(2.9) \quad C_*(x_*, x'_*) = C(x_*, x'_*) - C(x_*, X) (C(X, X) + I\sigma^2)^{-1} C(X, x'_*)$$

for any points $x_*, x'_* \in \mathcal{X}$. If I consider $X_* = (x_*, x'_*)^\top$, (2.8) results as the first block of (2.5), and (2.9) results as the top off-diagonal block of (2.6).

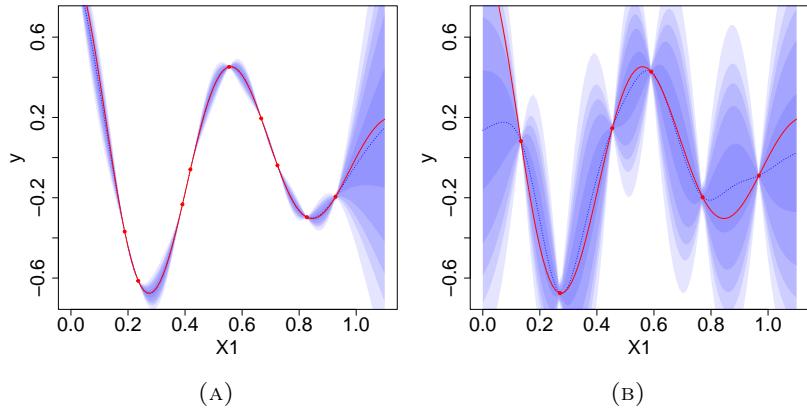


FIGURE 2.2. Predictive GP regressions given different numbers of training data points, for some values of σ^2 , ϕ , and β . Fig 2.2a shows the predictive GP given $n = 10$ data-points. Fig 2.2b shows the predictive GP given $n = 6$ data-points.

Note 20. Note that the posterior expected rule at $x_* \in \mathcal{X}$ is

$$(2.10) \quad \begin{aligned} E(h(x_*)|y) &= E(\eta(x_*)|y) = \mu(x_*) + C(x_*, X) (C(X, X) + I\sigma^2)^{-1} (y - \mu(X)) \\ &= \sum_{i=1}^n \alpha_i C(x_i, x_*), \text{ by assuming that } \mu(\cdot)=0 \end{aligned}$$

where $\alpha = y^\top (C(X, X) + I\sigma^2)^{-1} C(X, x_*)$. This is in accordance to the Representation theorem (Theorem 21 in Handout 7 Kernel methods) with reference to the Bayesian linear regression (Note 2).

Example 21. Given some values for σ^2 , ϕ , and β (see Section σ^2 , ϕ , and β), the predictive GP regression for $\eta(\cdot)$ is represented in Figure 2.2a. The red dots are training data points $\{z_i = (x_i, y_i)\}$. The red line is the predictive GP mean (2.8). The blue shades is the areas between the 95% quantiles of the Normal distribution $N(\mu_*(x_*), C_*(x_*, x_*))$ at each point x_* of the inputs. Note that $N(\mu_*(x_*), C_*(x_*, x_*))$ is just a snapshot of 2.7 at point x_* .

Note 22. (Related to Example 21) Intuitively, we can imagine that the conditional $\eta(\cdot)|y$ in Figure 2.2a results from the marginal/prior $\eta(\cdot)$ in Figure 2.1c by forcing all the lines in Figure 2.1c to pass through the data points in Figure 2.1c. So the more the data points, the smaller the uncertainty around the predictive mean; compare Fig 2.2a using 10 examples and Fig 3.1b using 5 examples.

3. TRAINING (VIA EMPIRICAL BAYES)

Note 23. Recall that the mean and covariance functions in (2.7) depend on tunable parameters σ^2 , ϕ , and β . When the number of training examples is small, the behavior of (2.7) is sensitive to these hyperparameters. In Figure 3.1, there are two instances of GP regression given 6 examples where the tunable parameters are different.

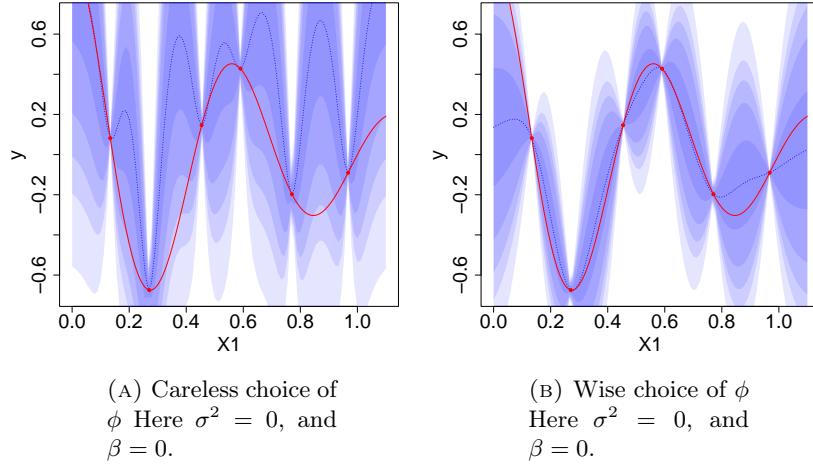


FIGURE 3.1. Sensitivity in the tunable parameters ϕ . Here $\sigma^2 = 0$, and $\beta = 0$.

Note 24. The marginal likelihood $f(y|\sigma^2, \phi, \beta)$ of y given the known parameters σ^2, ϕ, β results from (2.3) as

$$(3.1) \quad y|\sigma^2, \phi, \beta \sim N(\mu(X|\beta), C(X, X|\phi) + I\sigma^2).$$

Note 25. Let $\mu_\beta = \mu(X|\beta)$ and $C_\phi = C(X, X|\phi)$. To learn the unknown hyper-parameters $\theta = (\sigma^2, \phi, \beta)$ according to classical training methods, we can specify the empirical risk function $\hat{R}(\sigma^2, \phi, \beta)$ by using the marginal likelihood in (3.1), as $\hat{R}(\sigma^2, \phi, \beta) = -2 \log(f(y|\sigma^2, \phi, \beta))$ ². Hence, we need to minimize

$$(3.2) \quad \begin{aligned} \left(\hat{\sigma}^2, \hat{\phi}, \hat{\beta} \right) &= \arg \min_{\sigma^2, \phi, \beta} \left(-2 \log \left(N(y | \mu_\beta, C_\phi + I\sigma^2) \right) \right) \\ &= \arg \min_{\sigma^2, \phi, \beta} \left(\underbrace{\log(|C_\phi + I\sigma^2|) + (y - \mu_\beta)^\top (C_\phi + I\sigma^2)^{-1} (y - \mu_\beta)}_{\hat{R}(\sigma^2, \phi, \beta)} \right). \end{aligned}$$

Note 26. (3.2) can be solved via GD (Algorithm 8 Handout 2 Gradient descent), or the Stochastic Gradient (Handout 3: Stochastic gradient descent) as the required gradient can be easily computed as

$$\begin{aligned} \frac{dR}{d\beta_j} &= (C_\phi + I\sigma^2)^{-1} (y - \mu_\beta) \frac{d\mu_\beta}{d\beta_j} \\ \frac{dR}{d\phi_j} &= \text{tr} \left((C_\phi + I\sigma^2)^{-1} \left[\frac{\partial C_\phi}{\partial \phi_j} \right] \right) + (y - \mu)^\top (C_\phi + I\sigma^2)^{-1} \left[\frac{\partial C_\phi}{\partial \phi_j} \right] (C_\phi + I\sigma^2)^{-1} (y - \mu) \\ \frac{dR}{d\sigma^2} &= \text{tr} \left((C_\phi + I\sigma^2)^{-1} \right) + (y - \mu)^\top (C_\phi + I\sigma^2)^{-1} (C_\phi + I\sigma^2)^{-1} (y - \mu) \end{aligned}$$

No need to memorized

²Equivalently Empirical Bayes training procedures

4. EXAMPLES OF COVARIANCE FUNCTIONS

Note 27. The covariance function $C : \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty)$ with $C(x, x')$ describes how much two random variables x, x' change together.

Note 28. Covariance function is a functional parameter of the GP prior (2.2). Different covariance functions represent different properties, hence they impose different prior info in the GP regression model; they are crucial parameters.

Note 29. Any positive definite Kernel as described in Section 4 in Handout 7 Kernel methods can be used as a covariance function. Consequently kernel construction approaches and theories introduced can be used for the covariance functions as well.

Definition 30. Stationary covariance function is called a covariance function $C : \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty)$ whose image can be written as $C(x, x') = C(\|x - x'\|)$ namely, the dependence between any pair of input points x, x' is a function of their distance and only.

Note 31. In 1-D Gaussian process, one way to understand the characteristic length-scale of the process (if this exists) is in terms of the number of upcrossings of a level u . Consider stationary covariance function $C(x, x') = C(\|x - x'\|)$. The expected number of upcrossings $E(N_u)$ of the level u on the unit interval by a zero-mean, stationary, is

$$E(N_u) = \frac{1}{2\pi} \sqrt{\frac{-C''(0)}{C'(0)}} \exp\left(-\frac{u^2}{2C(0)}\right)$$

Note 32. Popular covariance functions are

Gaussian covariance function: given as

$$C(r) = \exp\left(-\frac{1}{2\phi^2}r^2\right)$$

- It is infinitely differentiable, which means that the GP is very smooth.
- The parameter ϕ is called lengthscale.
- The number of upcrossing at level u is $E(N_u) = (2\phi^2)^{-1}$ meaning that smaller ϕ represents more upcrossings, hence represents smaller scale dependences

Exponential Covariance Function: given as

$$C(r) = \exp\left(-\frac{1}{\phi}|r|\right)$$

- It is not differentiable at $r = 0$

Matern Class of Covariance Functions: given as

$$C_\nu(r) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}r}{\phi}\right)^\nu B_\nu\left(\frac{\sqrt{2\nu}r}{\phi}\right)$$

where $B_\nu(\cdot)$ is a modified Bessel functions (description of Bessel functions is out of the scope). Matern covariance function gives the Exponential one for $\nu = 1/2$, and the Gaussian one for $\nu \rightarrow \infty$.

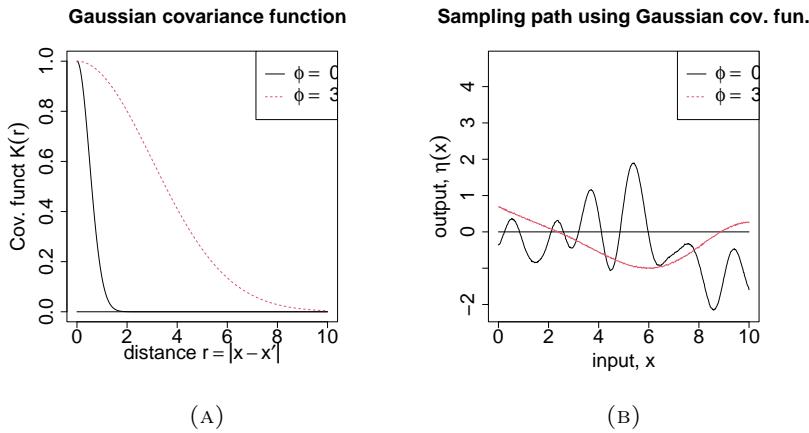


FIGURE 4.1. Investigation of the effect of choosing a different length scale ϕ

Example 33. We investigate the effect of choosing a different length scale ϕ in Figure 4.1. We consider the Gaussian covariance function. Figure 4.1a shows that the smaller the length scale ϕ is the more it focuses on the small scale dependences. Figure 4.1a shows that the smaller the length scale ϕ is the more upcrossings the associated GP has. This agrees with Note 31. Bottom line, it is reasonable for the researcher to specify a GP prior with a covariance with a smaller length scale ϕ if there is some a priori knowledge that the line $\eta(\cdot)$ may have strong small scale variation (eg changes more frequently).

Example 34. We investigate the effect of choosing a different covariance functions: Gaussian, Exponential, Matern 3/2 covariance functions (Note 32) in Figure 4.2. In Fig 4.2a, we observe that (a) covariance function reduces with distance that is closer points are expected to have larger dependence; (b) Gaussian cov gives more weight to small scale dependences (closer points) than the exponential. In Fig 4.2b, we see that the GP using a Gaussian cov is smoother than the GP using an Exponential cov. as expected as the former is infinitely differentiable while the latter is not differentiable. In general, Matern ($\nu = 3/2$)'s behavior is between those of Gaussian and Exponential, as expected.

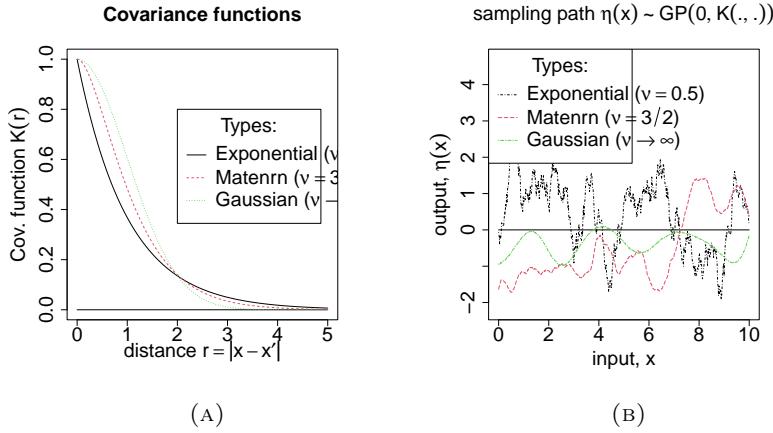


FIGURE 4.2. Investigation of the effect of choosing a different covariance functions: Gaussian, Exponential, Matern 3/2 covariance functions. (4.2a) shows the covariance function $C(x, x')$ against the distance $|x - x'|$. (4.2a) shows the covariance function $C(x, x')$ against the distance $|x - x'|$. (4.2b) shows realizations of GP using different covariance functions.

Note 35. Anisotropic versions of these isotropic covariance functions can be created anisotropy by setting $r(x, x') = (x - x')^\top M(x - x')$ for some positive semi-definite matrix M . If M is diagonal this implements the use of different length-scales on different dimension of inputs. Off-diagonal elements of M implement cross-dimensional dependencies in the inputs.

5. PRACTICAL MATERS

Note 36. One may consider some low degree polynomial form $\mu(x|\beta) = \sum_{j=0}^p x^j \beta_j$, however in this case $\mu(\cdot|\beta)$ and $C(\cdot, \cdot|\phi)$ may compete as $C(\cdot, \cdot|\phi)$ can express such behaviors according to Note 9. For this reason, the usual specification of $\mu(\cdot|\beta)$ in (2.2) is $\mu(x|\beta) = 0$ implying a priori complete uncertainty about the sign of $\eta(x)$ at each x .

Note 37. (Geostatistical model) Thinking as a statistician, one may decompose (2.2) as

$$\eta(\cdot) = \mu(\cdot|\beta) + \xi(\cdot|\phi)$$

where $\mu(\cdot|\beta)$ is modeled as a low degree polynomial (e.g., 2nd degree) representing large scale dependences (see polynomial regression in Term 1), and $\xi(\cdot|\phi) \sim GP(0, C(\cdot, \cdot|\phi))$ representing lower scale dependence by using an appropriate kernel. Seeing the big picture, (2.1) can be re-stated as $y_i = \mu(x_i|\beta) + \xi(x_i|\phi) + \epsilon_i$ where term ϵ_i represents noise (no dependence or so short scale dependence that can be considered as noise in the model), $\xi(x_i|\phi)$ represents low-scale dependence (about nearby inputs points), and $\mu(x_i|\beta)$ represents large scale dependence (about very distant inputs points).

6. PRACTICE, IMPLEMENTATION, AND CODE

Below is some practical examples on the implementation of Gaussian process regression in R programming environment by using the R packages: DiceKriging and DiceOptim.

- Roustant, O., Ginsbourger, D., & Deville, Y. (2012). DiceKriging, DiceOptim: Two R packages for the analysis of computer experiments by kriging-based metamodeling and optimization. *Journal of statistical software*, 51, 1-55.

The examples below were created for the undergraduate programme SURF 2016 at Purdue University (July 8, 2016) however they are suitable to the course for your practice.

Toy example

https://github.com/georgios-stats/Intro_GPR_SURF_2016/blob/master/Numerical_example.ipynb

The Piston Simulation function model in 2D

https://github.com/georgios-stats/Intro_GPR_SURF_2016/blob/master/Practice_2D.ipynb

Practice Catalytic Reaction 5D

https://github.com/georgios-stats/Intro_GPR_SURF_2016/blob/master/Practice_CatalyticReaction_5D_solution.ipynb

Practice Piston 7D

https://github.com/georgios-stats/Intro_GPR_SURF_2016/blob/master/Practice_Piston_7D_solution.ipynb

Practice Robot Arm 8D

https://github.com/georgios-stats/Intro_GPR_SURF_2016/blob/master/Practice_Robot_Arm_8D_solution.ipynb

APPENDIX A. MULTIVARIATE NORMAL DISTRIBUTION³ $x|\mu, \Sigma \sim N_d(\mu, \Sigma)$

Definition 38. A d -dimensional random variable $x \in \mathbb{R}^d$ is said to have a multivariate Normal (Gaussian) distribution, if for every d -dimensional fixed vector $\alpha \in \mathbb{R}^d$, the random variable $\alpha^\top x$ has a univariate Normal (Gaussian) distribution.

Definition 39. We denote the d -dimensional Normal distribution with mean μ and covariance matrix $\Sigma \geq 0$ as $N_d(\mu, \Sigma)$.

Notation 40. The d -dimensional standardized Normal distribution is $N_d(0, I)$.

Proposition 41. Let random variable $x \sim N_d(\mu, \Sigma)$, fixed vector $c \in \mathbb{R}^q$ and fixed matrix $A \in \mathbb{R}^q \times \mathbb{R}^d$. The random vector $y = c + Ax$ has distribution $y \sim N_q(c + A\mu, A\Sigma A^\top)$.

Proposition 42. Let a d -dimensional random vector $x \sim N_{(any)}(\mu, \Sigma)$.

- (1) Let $y = Ax$ and $z = Bx$, where $A \in \mathbb{R}^{q \times d}$ and $B \in \mathbb{R}^{k \times d}$: The vectors $y = Ax$ and $z = Bx$ are independent if and only if $A\Sigma B^\top = 0$.
- (2) Let $x = (x_1, \dots, x_d)^\top$: The x_1, \dots, x_d are mutually independent if and only if the corresponding off diagonal parts of the Σ are zero.

Proposition 43. Any sub-vector of a vector with multivariate Normal distribution has a multivariate Normal distribution.

Proposition 44. [Marginalization & conditioning] Let $x \sim N_d(\mu, \Sigma)$. Consider partition such that

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}; \quad \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}; \quad \Sigma = \begin{bmatrix} \Sigma_1 & \Sigma_{21}^\top \\ \Sigma_{21} & \Sigma_2 \end{bmatrix},$$

where $x_1 \in \mathbb{R}^{d_1}$, and $x_2 \in \mathbb{R}^{d_2}$. Then:

- (1) For the marginal, it is $x_1 \sim N_{d_1}(\mu_1, \Sigma_1)$.
- (2) For the conditional, if $\Sigma_1 > 0$, it is

$$x_2|x_1 \sim N_{d_2}(\mu_{2|1}, \Sigma_{2|1})$$

where

$$(A.1) \quad \mu_{2|1} = \mu_2 + \Sigma_{21}\Sigma_1^{-1}(x_1 - \mu_1) \quad \text{and} \quad \Sigma_{2|1} = \Sigma_2 - \Sigma_{21}\Sigma_1^{-1}\Sigma_{21}^\top$$

Proposition 45. The density function of the d -dimensional Normal distribution with mean μ and covariance matrix Σ , when Σ is symmetric positive definite matrix ($\Sigma > 0$), exists and it is equal to

$$(A.2) \quad f(x) = (2\pi)^{-\frac{d}{2}} \det(\Sigma)^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right)$$

³More detailed material about the Multivariate Normal distribution can be found in the can be found in “Handout 2: Revision in mixture of probability distributions” of the module “Bayesian Statistics III/IV (MATH3341/4031)” Michaelmas term, 2021 available from https://github.com/georgios-stats/Bayesian_Statistics_Michaelmas_2021/blob/main/Lecture_handouts/02_Revision_in_mixture_of_probability_distributions.pdf. The material in this section is just a sub-set of the statements in the referenced handout.