

ELEC-E5430 Large-Scale Data Analysis

Esa Ollila

Contents

1	Introduction	1
1.1	Classification task	2
1.2	Regression (prediction) task	3
1.3	Discussion	4
2	Classification: basic concepts	1
2.1	Optimal classifier	1
2.2	Classification costs and Bayes risk	4
2.3	Predictor and the loss functions	5
2.3.1	Logistic transformation	7
2.3.2	Loss functions for regression	8
2.3.3	Comparison of loss criterions	8
2.4	Performance measures	9
2.4.1	Test accuracy/error rate	10
2.4.2	Other measures	11
2.5	Conclusions	11
3	Decision tree methods	13
3.1	Decision tree	13
3.2	Regression trees	14
3.2.1	Pruning the tree	15
3.3	Classification trees	16
3.3.1	Choosing the variable to split on and the split point	17
3.4	Bagging	18
3.5	Random forests	20
4	Boosting	23
4.1	General ensemble scheme	23
4.2	AdaBoost	24
4.3	Forward Stagewise Additive Modeling	26

4.4	Exponential Loss + Stagewise additive modeling = AdaBoost.M1	29
4.5	Discussion	31
5	Lasso	33
5.1	Big Data Challenges	34
5.2	Penalized/Regularized regression	34
5.3	Ridge regression	35
5.3.1	Computation of the ridge estimator	36
5.3.2	Bias-variance tradeoff	37
5.4	Lasso	38
5.4.1	Geometry of the lasso	39
5.4.2	Lasso solution path	39
5.4.3	Standardizing the features and the penalty	40
5.5	Computation of the lasso solution	41
5.5.1	Subgradient optimality conditions	43
5.5.2	Cyclic coordinate descent	44
5.5.3	Lasso solution path	46
5.5.4	Why CCD works for large-scale data?	47
5.6	Discussion	48
6	Lasso: extensions	49
6.1	Elastic net	49
6.1.1	Computation via the CCD	50
6.2	Generalized lasso	52
6.2.1	Fused lasso	52
6.2.2	Trend filtering	56
6.3	Group lasso	59
6.4	Discussion	60

Chapter 1

Introduction

We have a set of *input* variables (or *features*)

$$X = (X_1, \dots, X_p)$$

that are used to predict the *output* (or *outcome* or *response*) variable Y . In a typical scenario, we have an outcome measurement, usually quantitative (such as a stock price) or categorical (such as disease/no disease), that we wish to predict based on a set of features (such as diet or clinical measurements).

We assume there is an underlying function $f(\cdot)$ that captures the input-output relationship which we would like to estimate. We do not know f , but we get to observe example input-output pairs. This is the so called *supervised learning* problem, where a *training data*

$$\mathcal{T} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$$

is available for estimating f . Often f is parametrized, so we only need to learn a vector parameter that determines the function f .

The training data pair (\mathbf{x}_i, y_i) is assumed to be generated at random via one of the following mechanism:

1. (\mathbf{x}_i, y_i) is a realisation from an unknown $p + 1$ variate joint distribution of (X, Y)
2. each input \mathbf{x}_i is drawn independently from some unknown distribution and we get to observe the pair $(\mathbf{x}_i, f(\mathbf{x}_i) + \varepsilon_i)$, where ε_i represents a random error term.

Many machine learning problems can be posed as classification or regression tasks. These differ in output types and consequently in the prediction tasks (prediction or classification).

We use uppercase letters such as X or Y when referring to an abstract variable or when viewing X as a random variable or random vector. For example $X_j \in \mathbb{R}$ can represent a feature such as *Height of a person* (in a specific population), but $x_j \in \mathbb{R}$ is its observed value. Similarly $X \in \mathbb{R}^p$ can represent a vector of feature variables, and $\mathbf{x} = (x_1, \dots, x_p)^\top \in \mathbb{R}^p$ observation of these variables for one measurement instance.

Thus for observed values (or non-random data points) we use the notations:

- *matrices* are represented by bold uppercase letters; E.g., $\mathbf{X} \in \mathbb{R}^{N \times p}$ denotes a fixed (known) $N \times p$ matrix.
- *vectors* are represented by bold lowercase letters. E.g., $\mathbf{x}_i \in \mathbb{R}^p$ denotes a fixed (known) $p \times 1$ vector. By a vector we always refer to a column vector.

The set of N input (observed) p -vectors \mathbf{x}_i , $i = 1, \dots, N$, in the training data are often collected to an $N \times p$ matrix:

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{Np} \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_N^\top \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_p \end{pmatrix}$$

This convention distinguishes a p -vector of inputs $\mathbf{x}_i \in \mathbb{R}^p$ for the i th observation from the N -vectors \mathbf{x}_j consisting of all the observations on variable (feature) X_j (e.g., "height"). Since all vectors are assumed to be column vectors, the i th row of \mathbf{X} is \mathbf{x}_i^\top , the vector transpose of $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^\top$. We collect the outputs, y_i , into a single vector

$$\mathbf{y} = (y_1, \dots, y_N)^\top.$$

1.1 Classification task

Classification is a problem where two or more populations are known a priori and one or more new observations are classified into one of the known populations (classes) based on the measured characteristics. It is assumed that there are known number $K \geq 2$ classes (or populations or states), the output Y taking values in a finite set $\mathcal{G} = \{\mathcal{G}_1, \dots, \mathcal{G}_K\}$. We would like to predict qualitative (categorical) outputs $Y \in \mathcal{G}$ representing class labels

Classification task can be often presented as a problem of partitioning the input vector space into disjoint (decision) regions. The task is to find the function (*discriminant rule*) $G(X) \mapsto Y$ that maps the inputs most accurately to their corresponding class labels.

It will be more convenient to encode the labels of the input variable Y with numeric values. We often use the convention

$$\mathcal{G} = \{1, \dots, K\}$$

in the general $K > 2$ case and

$$\mathcal{G} = \{1, 2\} \quad \text{or} \quad \mathcal{G} = \{-1, 1\}$$

in the two-class ($K = 2$) case. Such numeric codes are sometimes referred to as *targets*.

Training data \mathcal{T} consists of observations $\mathbf{x}_i \in \mathbb{R}^p$ from one of the K populations and $y_i \in \mathcal{G} = \{1, 2, \dots, K\}$ is the associated known class label of the i^{th} observation. The goal in the classification task is to assign a new observation $\mathbf{x} = (x_1, \dots, x_p)^\top \in \mathbb{R}^p$ to one of the $K \geq 2$ classes as accurately as possible.

A formal definition of a discriminant rule is given below.

Definition 1.1. A *discriminant rule* or *classifier* is a function $G(\mathbf{x}) : \mathbb{R}^p \rightarrow \mathcal{G} = \{1, \dots, K\}$, i.e., a function that takes a data vector and returns a class label. Discriminant rule partitions the input space \mathbb{R}^p into K *decision regions*,

$$\Gamma_k \equiv \Gamma_k(G) = \{\mathbf{x} \in \mathbb{R}^p : G(\mathbf{x}) = k\}$$

which are mutually disjoint and exhaustive sets, i.e.,

$$\Gamma_k \cap \Gamma_j = \emptyset \quad \forall k \neq j \quad \text{and} \quad \bigcup_{k=1}^K \Gamma_k = \mathbb{R}^p.$$

Equivalently, $G(\mathbf{x})$ can be formulated as

$$G(\mathbf{x}) = \arg \max_{k \in \{1, \dots, K\}} G_k(\mathbf{x}) \quad (1.1)$$

where $G_k(\mathbf{x}) \in \mathbb{R}$ is a *discriminant score* of \mathbf{x} for class k , and

$$\Gamma_k = \{\mathbf{x} \in \mathbb{R}^p : G_k(\mathbf{x}) > G_j(\mathbf{x}) \quad \forall k \neq j \in \{1, \dots, K\}\}$$

i.e., class k is chosen if it has the greatest discriminant score (and decide ties by random guesses).

In the binary classification problem ($K = 2$) and dimension $p = 2$, the Figure 1.1 illustrate the division of the sample space \mathbb{R}^p into two disjoint decision regions.

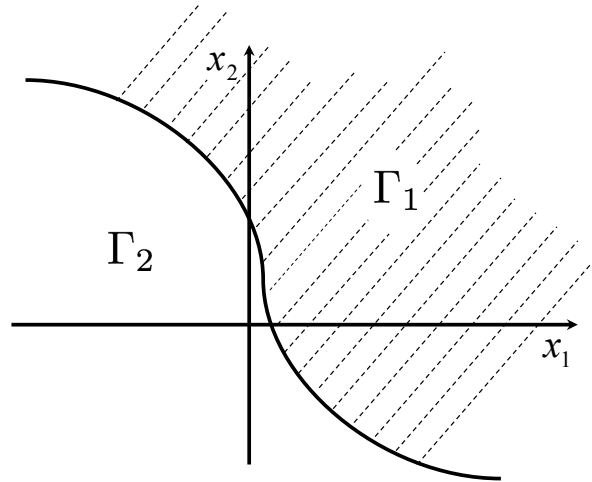


FIGURE 1.1: Two groups $K = 2$ and $\mathbf{x} \in \mathbb{R}^2$ case: a discriminant rule distinguishes \mathbb{R}^2 into two disjoint decision regions Γ_1 and Γ_2 ($\mathbb{R}^2 = \Gamma_1 \cup \Gamma_2$).

1.2 Regression (prediction) task

Refers to the case when we predict quantitative output $Y \in \mathbb{R}$, referred to as *response variable*, given the input $X = (X_1, \dots, X_p)^T$ where X_1, \dots, X_p are referred to as *predictor variables* or *features*. The task is to find the predictor function $f(X) \mapsto Y$ that

maps the inputs $X = (X_1, \dots, X_p)^\top$ most accurately to their corresponding outputs Y (for example by minimizing the MSE for training data: $\sum_{i=1}^N (y_i - f(\mathbf{x}_i))^2$). In high-dimensional cases, feature selection, so selecting significant features and getting rid of noisy features is essential part of the problem.

The regression model is

$$Y = f(X) + \varepsilon,$$

where ε is the zero mean random error term that account for modelling and measurement errors, and the predictor function is

$$f(X) = g\left(\beta_0 + \sum_{j=1}^p X_j \beta_j\right), \quad (1.2)$$

where $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top \in \mathbb{R}^p$ is an unknown vector of *regression coefficients*, $\beta_0 \in \mathbb{R}$ is the *intercept* and g is a fixed *link function* such as linear function $g(x) = x$, yielding the *linear regression model*.

$$f(X) = \beta_0 + \sum_{j=1}^p X_j \beta_j. \quad (1.3)$$

Thus we observe the pairs (\mathbf{x}_i, y_i) where $y_i = f(\mathbf{x}_i) + \varepsilon_i$, where $\{\varepsilon_i\}_{i=1}^N$ are i.i.d. unobserved zero mean random errors. Then using the training data $\mathcal{T} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, the aim is to estimate the unknown regression coefficients $(\beta_0, \boldsymbol{\beta}) \in \mathbb{R} \times \mathbb{R}^p$ for a given function f .

Q: Let response Y be the house price of an apartment in Helsinki. What would be useful predictors X_1, \dots, X_p in the linear regression model?

In terms of the training data, the linear input-output relationship can be represented by a set of N equations

$$\begin{aligned} y_1 &= \beta_0 + x_{11}\beta_1 + \dots + x_{1p}\beta_p + \varepsilon_1 \\ &\vdots \\ y_N &= \beta_0 + x_{N1}\beta_1 + \dots + x_{Np}\beta_p + \varepsilon_N \end{aligned}$$

where ε_i -s represent i.i.d. random variables from the noise distribution ε . These can be expressed using matrix-vector notations as

$$\begin{aligned} \mathbf{y} &= \beta_0 \mathbf{1} + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \\ &= \beta_0 \mathbf{1} + \boldsymbol{\beta}_1 x_1 + \dots + \boldsymbol{\beta}_p x_p + \boldsymbol{\varepsilon}. \end{aligned} \quad (1.4)$$

1.3 Discussion

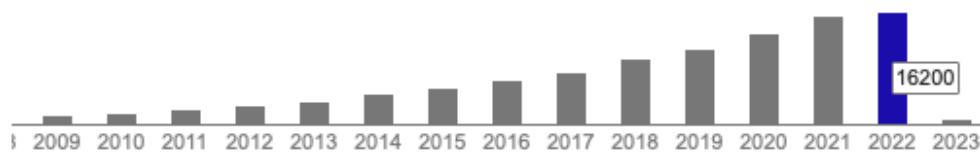
The aim of these chapters are to provide a brief introduction to methods that are proven to be powerful in supervised learning tasks. Most data analysis languages such as python have put a lot of effort to provide powerful and computationally efficient implementations of these methods. We focus on 3 powerful ideas in supervised learning: a) Decision trees, Bagging and Random Forests, b) Boosting, and c) Lasso regression and its extensions.

Paper	citations	citations / year
Lasso [Tibshirani, 1996]	48872	1879
Decision Trees [Breiman et al., 1984]	58493	1551
Random forests [Breiman, 2001]	102244	4868
Boosting [Freund and Schapire, 1997]	24957	998
Gradient Boosting [Friedman, 2001]	20429	972

TABLE 1.1: *Methods of the chapters. (Citation count: Jan 8, 2023)*

Citations time series for random forest and lasso (Jan 8, 2023):

Cited by 102244

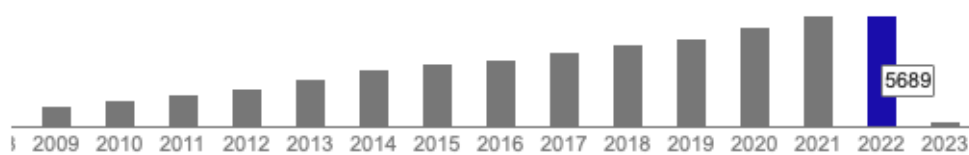


Random forests

L Breiman - Machine learning, 2001

[Cited by 102244](#) [Related articles](#) [All 83 versions](#)

Cited by 48872



Regression shrinkage and selection via the lasso

R Tibshirani - Journal of the Royal Statistical Society: Series B ..., 1996

[Cited by 48863](#) [Related articles](#) [All 49 versions](#)

Chapter 2

Classification: basic concepts

This chapter focuses on building the basic concepts of classification and the optimal Bayesian discriminant rule. The chapter covers Sections 10.5 and Section 10.6 from [Hastie et al., 2009] as well as some selected topics from Chapter 2, such as Section 2.4.

2.1 Optimal classifier

We assume that the joint $(p + 1)$ -variate distribution of input-output pair (Y, X) is known. The class label $Y \in \mathcal{G} = \{1, \dots, K\}$ is a discrete random variable with a probability mass function (p.m.f.)

$$\pi_k = \Pr(Y = k) = \begin{cases} \text{"probability that a randomly selected} \\ \text{observation is from class } k \end{cases}$$

These are the *a priori class probabilities* ($\sum_i \pi_i = 1$) which are assumed to be known. The prior probabilities can be significantly different between the classes.

For simplicity of exposition, assume that the p -variate feature vector contains measurements on a continuous p -variate random vector $X = (X_1, \dots, X_p)$. The class conditional probability density function (p.d.f.)

$$f_{X|Y}(\mathbf{x} | k), \quad k \in \mathcal{G} \quad (2.1)$$

is the distribution of X when it is from class k . Note also that

$$\Pr(X \in \mathcal{X} | Y = k) = \int_{\mathcal{X}} f_{X|Y}(\mathbf{x} | k) d\mathbf{x}.$$

The data conditional probability, i.e., the *a posteriori class probabilities*, are

$$p_k(\mathbf{x}) = \Pr(Y = k | X = \mathbf{x}) = \frac{f_{X|Y}(\mathbf{x} | k) \Pr(Y = k)}{f_X(\mathbf{x})} = \frac{f_{X|Y}(\mathbf{x} | k) \pi_k}{\sum_{k=1}^K f_{X|Y}(\mathbf{x} | k) \pi_k} \quad (2.2)$$

Note that $p_k(\mathbf{x}) \in [0, 1]$ and

$$\sum_{k \in \mathcal{G}} p_k(\mathbf{x}) = 1. \quad (2.3)$$

It is customary to take logarithm of p.d.f.-s as it often simplifies expressions. We express the log-posterior by

$$\log p_k(\mathbf{x}) = \ln f_{X|Y}(\mathbf{x} | k) + \ln \pi_k, \quad (2.4)$$

where we ignore the irrelevant additive constant ($= -\ln f_X(\mathbf{x})$) that does not depend on k .

The *classification loss* or *0-1 loss* is defined as

$$L_{0/1}(y, G(\mathbf{x})) = 1_{\{y \neq G(\mathbf{x})\}} \quad (2.5)$$

which is equal to 1 if the classifier G misclassifies (\mathbf{x}, y) , and 0 otherwise. Since we assume a full knowledge of the joint distribution, we can design an optimal classifier, which is the one that minimizes the risk (probability of error).

Definition 2.1. The (Bayes) *risk* of a discriminant rule G is

$$r(G) = \Pr(G(X) \neq Y) = \mathbb{E}[1_{\{Y \neq G(X)\}}] \quad (2.6)$$

and it equals the probability that the rule makes an error (or the expected classification loss).

The risk measures how the estimated rule performs on average. It can also be written as

$$r(G) = \sum_{k \in \mathcal{G}} \mathbb{E}_X[1_{\{G(X) \neq k\}} | Y = k] \pi_k = \sum_{k \in \mathcal{G}} \Pr(X \notin \Gamma_k(G) | Y = k) \pi_k \quad (2.7)$$

where $\Gamma_k(G)$ denote the classification regions (cf. [Definition 1.1](#)). Its empirical version, computed on the training sample,

$$\hat{r}_N(G) = \frac{1}{N} \sum_{i=1}^N L_{0/1}(y_i, G(\mathbf{x}_i)) = \frac{1}{N} \sum_{i=1}^N 1_{\{y_i \neq G(\mathbf{x}_i)\}} \quad (2.8)$$

is called *empirical risk*.

The Bayes classifier is a rule that maximizes the posterior probability (2.2).

Definition 2.2. The *Bayes classifier* G^* assigns \mathbf{x} to class k^* having maximum a posteriori probability, $G^*(\mathbf{x}) = k^*$, where¹

$$\begin{aligned} k^* &= \arg \max_{k \in \mathcal{G}} p_k(\mathbf{x}) = \arg \max_{k \in \mathcal{G}} f_{X|Y}(\mathbf{x} | k) \pi_k \\ &= \arg \max_{k \in \mathcal{G}} \{\ln f_{X|Y}(\mathbf{x} | k) + \ln \pi_k\}, \end{aligned} \quad (2.9)$$

or in other words

$$p_{k^*}(\mathbf{x}) \geq p_j(\mathbf{x}) \Leftrightarrow \pi_{k^*} f_{X|Y}(\mathbf{x} | k^*) \geq \pi_j f_{X|Y}(\mathbf{x} | j) \quad \forall k^* \neq j.$$

¹Since the marginal distribution $f_X(\mathbf{x})$ of X (i.e., the unconditional likelihood over all populations that we could observe \mathbf{x}) is irrelevant constant that just scales the posterior probabilities in (2.2), and therefore, we can ignore it in (2.9)

Bayes rule is also the optimal classification rule:

Theorem 2.1. *The Bayes rule G^* is optimal in the sense of minimizing the error rate (risk), i.e., it verifies*

$$r(G) \geq r(G^*)$$

where $G(\cdot)$ can be any classifier.

Remark 2.1. In the binary classification problem ($K = 2$, $\mathcal{G} = \{0, 1\}$), the Bayes classifier can be simply written as

$$G^*(\mathbf{x}) = \begin{cases} \text{Class 0,} & \text{if } p_1(\mathbf{x}) - p_0(\mathbf{x}) < 0 \\ \text{Class 1,} & \text{if } p_1(\mathbf{x}) - p_0(\mathbf{x}) > 0 \end{cases} \quad (2.10)$$

$$= \begin{cases} \text{Class 0,} & \text{if } p_1(\mathbf{x}) - \frac{1}{2} < 0 \\ \text{Class 1,} & \text{if } p_1(\mathbf{x}) - \frac{1}{2} > 0. \end{cases} \quad (2.11)$$

The case of ties, so when an observation falls in the decision boundary

$$f(\mathbf{x}) = p_1(\mathbf{x}) - p_0(\mathbf{x}) = p_1(\mathbf{x}) - 1/2 = 0$$

can be handled by a coin flip. Note that in the last identity we used the property that $p_0(\mathbf{x}) + p_1(\mathbf{x}) = 1$. Thus, encoding the classes as $\mathcal{G} = \{-1, 1\}$ instead of $\mathcal{G} = \{0, 1\}$ allows us to write the Bayes rule concisely as

$$G^*(\mathbf{x}) = \text{sign}[f(\mathbf{x})] \quad \text{with} \quad f(\mathbf{x}) = p_1(\mathbf{x}) - \frac{1}{2}. \quad (2.12)$$

Remark 2.2. When $K = 2$ and $\mathcal{G} = \{0, 1\}$, we may express the decision regions as

$$\Gamma_0^* = \left\{ \mathbf{x} : \frac{f_{X|Y}(\mathbf{x} | 0)}{f_{X|Y}(\mathbf{x} | 1)} > \frac{\pi_1}{\pi_0} \right\} \quad \text{and} \quad \Gamma_1^* = \left\{ \mathbf{x} : \frac{f_{X|Y}(\mathbf{x} | 0)}{f_{X|Y}(\mathbf{x} | 1)} < \frac{\pi_1}{\pi_0} \right\}$$

and handling ties as random guessing. The detection rule can thus be expressed as

$$L(\mathbf{x}) = \frac{f_{X|Y}(\mathbf{x} | 0)}{f_{X|Y}(\mathbf{x} | 1)} \underset{1}{\overset{0}{\geq}} \frac{\pi_1}{\pi_0}$$

and notice that $L(\mathbf{x})$ is the *likelihood ratio*. If the a priori probabilities are identical, $\pi_0 = \pi_1 = 1/2$, then the decision regions are completely based on comparing the likelihoods of \mathbf{x} for each class.

Example 2.1. Assume $K = 2$ classes with class conditional distributions following exponential distributions:

$$\begin{aligned} X|Y = 0 &\sim \text{Exp}(\lambda_0) \quad \text{and} \quad X|Y = 1 \sim \text{Exp}(\lambda_1) \\ f_{X|Y}(x | 0) &= \lambda_0 \exp\{-\lambda_0 x\} \quad \text{and} \quad f_{X|Y}(x | 1) = \lambda_1 \exp\{-\lambda_1 x\}, \quad x > 0 \end{aligned}$$

where $\lambda_k > 0$ is the rate parameter, $\lambda_k^{-1} = \mathbb{E}[X | Y = k]$, $k \in \{0, 1\}$. Without loss of generality, assume $\lambda_0 < \lambda_1$ (this can be always done by simply relabelling the classes).

- Derive the classification region Γ_0^* (that minimize the Bayes risk) in the case of uniform prior probabilities ($\pi_0 = \pi_1 = 1/2$).
- Calculate the risk (so the probability of an error) $r(G^*) = \Pr(Y \neq G^*(X))$ of this optimal Bayes classifier when $\lambda_0 = 1$ and $\lambda_1 = 3$.

If time permits, we go through this example on lectures... ■

2.2 Classification costs and Bayes risk

Discriminant rule should yield as few misclassifications as possible on average. In some applications one may also want to take into account possible costs caused by misclassifications.

Quantifying the consequences of the decisions can be formally expressed via a *cost function*:

$$C(k, j) : \mathcal{G} \times \mathcal{G} \rightarrow \mathbb{R}$$

which quantifies the cost of misclassifying an observation into class j when its true class is k . Commonly, one assumes that

$$C(k, k) = 0 \quad \text{and} \quad C(k, j) > 0, \quad \forall k \neq j \in \mathcal{G}$$

i.e., the cost is zero for correct classification and non-zero when an observation is incorrectly classified.

The simplest cost is the *uniform cost*

$$C(k, j) = \mathbf{1}_{\{k \neq j\}} = \begin{cases} 1, & k \neq j \\ 0, & k = j \end{cases} \quad (2.13)$$

as it gives unit cost to all errors. Our goal is then to find a rule $G(\mathbf{x})$ that minimizes the *expected cost of misclassification*, defined as

$$\text{ECM}(G) = \mathbb{E}_{X,Y} [C(Y, G(X))] \quad (2.14)$$

$$= \sum_{k=1}^K \mathbb{E}_X [C(k, G(X)) \mid Y = k] \pi_k, \quad (2.15)$$

where $\mathbb{E}_X [C(k, G(X)) \mid Y = k]$ is the expected cost of misclassification, when the observation is from class k . Notice that the Bayes risk $r(G)$ coincides with ECM based on uniform costs.

Suppose we have two classes. The (expected) cost of classifying an observation from class 0 to class 1 is

$$\mathbb{E}_X [C(0, G(X)) \mid Y = 0] = \int C(0, G(\mathbf{x})) f_{X|Y}(\mathbf{x} \mid 0) d\mathbf{x} = C(0, 1) \int_{\Gamma_1(G)} f_{X|Y}(\mathbf{x} \mid 0) d\mathbf{x}$$

and similarly

$$\mathbb{E}_X [C(1, G(X)) \mid Y = 1] = C(1, 0) \int_{\Gamma_0(G)} f_{X|Y}(\mathbf{x} \mid 1) d\mathbf{x}.$$

The ECM in (2.14) is

$$\pi_0 \cdot C(0, 1) \int_{\Gamma_1(G)} f_{X|Y}(\mathbf{x} \mid 0) d\mathbf{x} + \pi_1 \cdot C(1, 0) \int_{\Gamma_0(G)} f_{X|Y}(\mathbf{x} \mid 1) d\mathbf{x}.$$

We then have the following result.

Theorem 2.2. *The discriminant rule $G^*(\mathbf{x})$ that minimizes the ECM is based on discriminant scores*

$$\begin{aligned} G_0(\mathbf{x}) &= \ln f_{X|Y}(\mathbf{x} | 0) + \ln \pi_0 + \ln C(0, 1), \\ G_1(\mathbf{x}) &= \ln f_{X|Y}(\mathbf{x} | 1) + \ln \pi_1 + \ln C(1, 0), \end{aligned}$$

where the decision regions are

$$\begin{aligned} \Gamma_0^* &= \{\mathbf{x} : G_0(\mathbf{x}) \geq G_1(\mathbf{x})\} = \left\{ \mathbf{x} : \frac{f_{X|Y}(\mathbf{x} | 0)}{f_{X|Y}(\mathbf{x} | 1)} \geq \frac{C(1, 0)}{C(0, 1)} \cdot \frac{\pi_1}{\pi_0} \right\}, \\ \Gamma_1^* &= \{\mathbf{x} : G_0(\mathbf{x}) < G_1(\mathbf{x})\} = \left\{ \mathbf{x} : \frac{f_{X|Y}(\mathbf{x} | 0)}{f_{X|Y}(\mathbf{x} | 1)} < \frac{C(1, 0)}{C(0, 1)} \cdot \frac{\pi_1}{\pi_0} \right\}. \end{aligned}$$

The optimal discriminant rule $G^*(\mathbf{x})$ is actionable if we know perfectly the conditional class p.d.f.'s $f_{X|Y}(\mathbf{x} | 0)$ and $f_{X|Y}(\mathbf{x} | 1)$, prior probabilities π_0 and π_1 and misclassification costs $C(0, 1)$ and $C(1, 0)$, or their ratios:

$$P = \frac{\pi_1}{\pi_0}, \quad C = \frac{C(1, 0)}{C(0, 1)}, \quad \text{and} \quad L(\mathbf{x}) = \frac{f_{X|Y}(\mathbf{x} | 0)}{f_{X|Y}(\mathbf{x} | 1)}. \quad (2.16)$$

Note that $L(\mathbf{x})$ is the likelihood ratio (*cf.* [Remark 2.2](#)). Choosing uniform costs, so $C(0, 1) = C(1, 0) = 1$, one obtain the Bayes rule which minimizes the risk $r(G)$.

Example 2.2. Suppose it is known that it is twice as costly to assign an observation from class 1 to class 0 than vice versa and that approximately 20% of observations belong to class 1. When an observation \mathbf{x} receives values $f_{X|Y}(\mathbf{x} | 0) = 0.3$ and $f_{X|Y}(\mathbf{x} | 1) = 0.4$, then is it classified to class 1 or class 2? ■

2.3 Predictor and the loss functions

Consider the binary class problem with output variable $Y \in \mathcal{G} = \{-1, 1\}$. Then a classifier $G(\cdot)$ produces a prediction taking one of the two values, -1 or 1 , and can be represented as

$$G(\mathbf{x}) = \text{sign}[f(\mathbf{x})] \quad (2.17)$$

for some function $f : \mathbb{R}^p \rightarrow \mathbb{R}$, called the predictor function. We can write the misclassification loss in (2.5) as

$$L_{0/1}(y, f(\mathbf{x})) = 1_{\{y \text{sign}[f(\mathbf{x})] \neq 1\}} = 1_{\{yf(\mathbf{x}) < 0\}} \quad (2.18)$$

where $m = yf(\mathbf{x})$ is called the *margin* of (y, \mathbf{x}) . The 1-0 loss in (2.18) assigns unit penalty for negative margin values, and no penalty at all for positive ones. Consequently, the risk (2.6) and the empirical risk (2.8) can also be expressed as:

$$r(f) = \mathbb{E}[1_{\{Yf(X) < 0\}}] \quad (2.19)$$

$$r_N(f) = \frac{1}{N} \sum_{i=1}^N 1_{\{y_i f(\mathbf{x}_i) < 0\}}. \quad (2.20)$$

name	$L(y, f(\mathbf{x}))$	$f^*(\mathbf{x})$
<i>misclassification loss</i>	$1_{\{yf(\mathbf{x}) < 0\}}$	$p(\mathbf{x}) - 1/2$
<i>exponential loss</i>	$\exp(-yf(\mathbf{x}))$	$\frac{1}{2} \log \frac{p(\mathbf{x})}{1 - p(\mathbf{x})}$
<i>binomial deviance</i>	$\log(1 + e^{-2yf(\mathbf{x})})$	$\frac{1}{2} \log \frac{p(\mathbf{x})}{1 - p(\mathbf{x})}$
<i>support vector (hinge loss)</i>	$[1 - yf(\mathbf{x})]_+$	$p(\mathbf{x}) - 1/2$

TABLE 2.1: Popular loss functions for classification as well as the associated optimal predictor function solving (2.21). Notation $[y]_+$ means positive part of y , i.e., $[y]_+ = y1_{\{y > 0\}} = \max(y, 0)$. We use shorthand notation $p(\mathbf{x})$ for $p_1(\mathbf{x}) = \Pr(Y = 1 \mid X = \mathbf{x})$.

Thus the risk empirical $r_N(f)$ is simply the average of negative margins. An observation (\mathbf{x}_i, y_i) with positive margin $y_i f(\mathbf{x}_i) > 0$ is classified correctly whereas those with negative margins $y_i f(\mathbf{x}_i) < 0$ are misclassified.

It would then be natural to find f that minimizes the empirical risk (2.8). However, this problem turns out to be NP-complete, meaning that no polynomial-time algorithm is believed to exist to solve it. Thus different loss function $L(y, f(\mathbf{x})) : \mathcal{G} \times \mathbb{R} \rightarrow \mathbb{R}$ are used instead of the misclassification loss $L_{0/1}(y, f(\mathbf{x}))$. In the two-class case, $L(y, f(\mathbf{x}))$ will be commonly a function of the margin $yf(\mathbf{x})$. In the regression set-up, where $f(\mathbf{x})$ expresses the dependence of output $y \in \mathbb{R}$ on input $\mathbf{x} \in \mathbb{R}^p$, the predictor function can be e.g., a linear regression fit, $f(\mathbf{x}) = \mathbf{x}^\top \beta$.

The predictor function $f(\mathbf{x})$ that defines the classifier in (2.17) is chosen as the function that minimizes the (conditional) expected loss

$$f^*(\mathbf{x}) = \arg \min_{f(\mathbf{x})} \mathbb{E}_Y [L(Y, f(\mathbf{x})) \mid X = \mathbf{x}] \quad (2.21)$$

This is the optimal population rule, assuming knowledge of the conditional distribution $Y \mid X = \mathbf{x}$. Table 2.1 specifies the optimal predictors function for each considered loss function, where we have used shorthand notation $p(\mathbf{x})$ for $p_1(\mathbf{x}) = \Pr(Y = 1 \mid X = \mathbf{x})$. The goal in (2.21) is to produce positive margins as frequently as possible. Naturally, any loss criterion used to derive $f(\mathbf{x})$ should penalize negative margins more heavily than positive ones since positive margin observations are already correctly classified.

For misclassification loss, the classifier is the Bayes classifier whose predictor function is $f^*(\mathbf{x}) = p(\mathbf{x}) - 1/2$ as pointed out in Remark 2.1. For exponential loss function the optimal rule equals one-half of log odds. As we shall see later, the exponential loss is related to AdaBoost classifier and $f^*(\mathbf{x})$ is the minimizer of the population version of the AdaBoost criterion. For binomial deviance loss, one obtains the same population minimizer $f^*(\mathbf{x})$ as for exponential loss function. For support vector loss, the optimal predictor function $f^*(\mathbf{x})$ is the same as for Bayes classifier. In this sense, support vector loss can be preferred. Comparison of loss functions are given in subsection 2.3.3.

Note that $1_{\{x \leq 0\}} \leq e^{-x}$ and thus we notice from (2.20) that the empirical risk can

be upper bounded by

$$\frac{1}{N} \sum_{i=1}^N e^{-y_i f(\mathbf{x}_i)}.$$

Thus AdaBoost classifier that aims at minimizing the exponential loss can be interpreted as method that minimizes the upper bound of the empirical risk (training error). In fact, binomial deviance and SVM loss have the same interpretation as will become clear from [Figure 2.1](#).

2.3.1 Logistic transformation

The binomial deviance loss function (*cf.* [Table 2.1](#)) equals the negative log-likelihood function of the binomial distribution, called the *deviance*, using the logit transformation defined below.

Definition 2.3. In 2-class problem, let $p(\mathbf{x}) = \Pr(Y = 1 \mid X = \mathbf{x})$. Then define the symmetric *logistic transformation* as

$$f(\mathbf{x}) = \frac{1}{2} \log \frac{p(\mathbf{x})}{1-p(\mathbf{x})} \Leftrightarrow p(\mathbf{x}) = \frac{e^{f(\mathbf{x})}}{e^{f(\mathbf{x})} + e^{-f(\mathbf{x})}} = \frac{1}{1 + e^{-2f(\mathbf{x})}} \quad (2.22)$$

i.e., f equals half of the log-odds ratio.

Symmetric logistic transformation is related to conventional logistic or *sigmoid function*, commonly used to transform values on $(-\infty, \infty)$ into numbers on $(0, 1)$, defined as

$$\sigma(z) = \frac{e^z}{1 + e^z} = \frac{1}{1 + e^{-z}} \Leftrightarrow \sigma^{-1}(z) = \log \frac{z}{1-z}. \quad (2.23)$$

Thus $p(\mathbf{x}) = \sigma(2f(\mathbf{x}))$ with $f(\mathbf{x}) = \frac{1}{2} \sigma^{-1}(p(\mathbf{x}))$.

In the two class case, we may model a random variable $Y' \in \{0, 1\}$ as being generated from a Bernoulli distribution $\text{Ber}(p(\mathbf{x}))$, where $p(\mathbf{x})$ is defined in (2.22). The conditional probability mass function of $Y'|X = \mathbf{x}$ is thus

$$f_{Y'|X}(y'|\mathbf{x}) = p(\mathbf{x})^{y'} (1 - p(\mathbf{x}))^{1-y'}, \quad y' \in \{0, 1\}.$$

The Bernoulli log-likelihood function is then

$$l(y', p(\mathbf{x})) = y' \log p(\mathbf{x}) + (1 - y') \log(1 - p(\mathbf{x})) \quad (2.24)$$

which is also sometimes referred to as *cross entropy loss*. The *Binomial deviance*, i.e., the negative log-likelihood, $-l(y, p(\mathbf{x}))$, can be written using the relationship of $f(\mathbf{x})$ and $p(\mathbf{x})$ in (2.22) and output encoding $y' = (y + 1)/2$, as

$$-l(y, f(\mathbf{x})) = \log(1 + e^{-2yf(\mathbf{x})}), \quad y \in \{-1, 1\}. \quad (2.25)$$

Note that $\exp(-yf(\mathbf{x}))$ itself is not a proper log-likelihood, as it does not equal the log of any probability mass function on plus or minus 1.

2.3.2 Loss functions for regression

In the regression the response is non-categorical, $Y \in \mathbb{R}$, and different loss functions are used. The most common loss functions are the *squared loss* (or L_2 -loss)

$$L(y, f(\mathbf{x})) = \frac{1}{2}(y - f(\mathbf{x}))^2 \quad (2.26)$$

or the *absolute loss* (or L_1 -loss)

$$L(y, f(\mathbf{x})) = |y - f(\mathbf{x})|. \quad (2.27)$$

They can also be used for binary-classification. In this case, they can be parametrized by margin value $yf(x)$ since for $y \in \{-1, 1\}$ one has that

$$\begin{aligned} |y - f(\mathbf{x})| &= |1 - yf(\mathbf{x})|, \\ (y - f(\mathbf{x}))^2 &= (1 - yf(\mathbf{x}))^2 = 1 - 2yf(\mathbf{x}) + (yf(\mathbf{x}))^2. \end{aligned}$$

However, these functions are non-monotone functions of the margin value which is not optimal for classification setting as described in the next section.

For squared-error loss, the optimal predictor function $f^*(\mathbf{x})$ is

$$\begin{aligned} f^*(\mathbf{x}) &= \arg \min_{f(\mathbf{x})} \frac{1}{2} \mathbb{E}_Y [(Y - f(\mathbf{x}))^2 \mid X = \mathbf{x}] = \frac{1}{2} \mathbb{E}[Y \mid X = \mathbf{x}] \\ &= p(\mathbf{x}) - 1/2. \end{aligned}$$

Thus we see that the population predictor function for the squared-error loss function is the same as for optimal Bayes classifier that minimizes the risk. However, squared-error loss function fails to penalize negative margins.

2.3.3 Comparison of loss criterions

Figure 2.1 shows the misclassification, exponential, binomial deviance, support vector and squared-error loss function as a function of the margin $m = y \cdot f(\mathbf{x})$. Both the exponential, deviance and support vector loss can be viewed as monotone decreasing continuous approximations (surrogates) to misclassification loss or convex upper bounds for the misclassification error. Misclassification loss is non-differentiable and also non-convex as a function of the margin m and hence not suitable for boosting which essentially uses functional gradient descent for optimization.

The difference between them is in degree of penalizing negative margins. Binomial deviance can be considered to be *more robust* than exponential loss as it assigns relatively less influence on observations with large negative margins. Unlike exponential loss, it grows linearly as the margin value m tends to $-\infty$. Hence its performance is not as much affected if there are misspecification of the class labels in the training data. We also notice from Figure 2.1 that the squared-error loss is not a good surrogate for misclassification error.

When robustness is a concern, squared-error loss for regression and exponential loss for classification are not the best criteria from a statistical perspective. However, they both lead to the elegant modular boosting algorithms in the context of forward stagewise additive modeling.

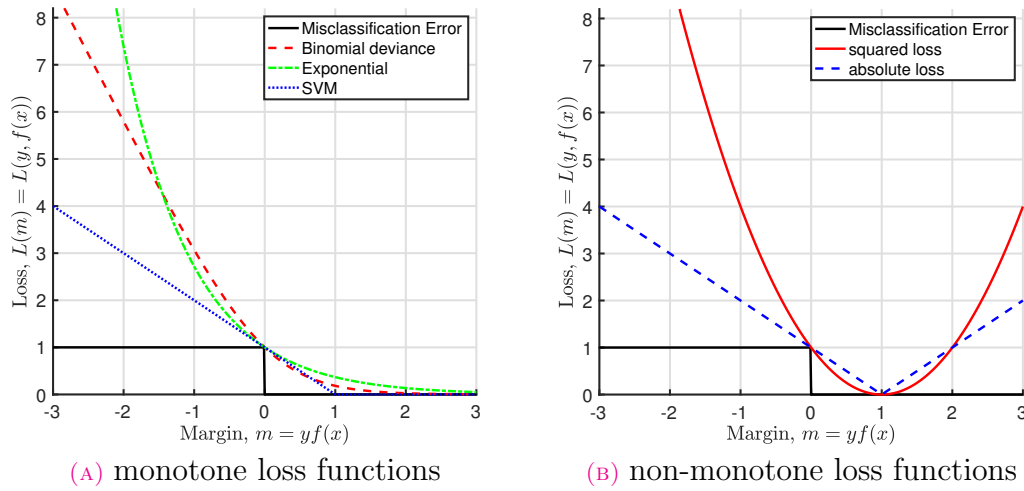


FIGURE 2.1: Loss functions $L(y, f(\mathbf{x}))$ for two-class classification is expressed as a function of margin $m = yf(\mathbf{x})$ as $L(y, f(\mathbf{x})) = L(m)$. The response is $y = \pm 1$; the prediction is $f(\mathbf{x})$, with class prediction $\text{sign}[f(\mathbf{x})]$. The losses $L(m)$ in (a) are misclassification: $1_{\{m < 0\}}$; exponential: $\exp(-m)$; binomial deviance: $\log(1 + \exp(-2m))$; and support vector: $(1 - m)_+$ and in (b) squared error: $(1 - m)^2$; and absolute $|1 - m|$. Each function has been scaled so that it passes through the point $(0, 1)$.

2.4 Performance measures

Given the estimated rule $\hat{G}(\mathbf{x})$, its risk (cf. Definition 2.1), $r(\hat{G}) = \Pr(\hat{G}(X) \neq Y)$ is difficult to calculate and hence one computes its empirical version. Consider binary classification problem and encode the classes as negative ('-'), e.g., a negative result on covid-19 test, and positive ('+'), e.g., a positive result on covid-19 test. The *training error rate* or *empirical risk* is then defined by

$$r_N(\hat{G}) = \frac{1}{N} \sum_{i=1}^N 1_{\{\hat{G}(\mathbf{x}_i) \neq y_i\}} = \frac{\text{FP} + \text{FN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

where $N = \text{TP} + \text{TN} + \text{FP} + \text{FN}$, and

FN (aka *miss-detection*)
 = # of observations of class + that are misclassified
 FP (aka *false alarm*)
 = # of observations of class - that are misclassified

where FP and FN (TP and TN) stand for false (true) positive and negative, respectively. These numbers are often represented via 2×2 *confusion matrix*:

		Predicted class		sum
		+	-	
True class	+	TP	FN	TP + FN
	-	FP	TN	FP + TN

Usually normalized confusion matrix is reported which gives the numbers as percentages (per total number of instances in each class, i.e., figures are “row-normalized”). The *true positive rate* (TRP) and *true negative rate* (TNR) are

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad \text{and} \quad \text{TNR} = \frac{\text{TN}}{\text{FP} + \text{TN}}.$$

Training error rate is far too optimistic and underestimates the true error rate (risk). This is because the data used to compute the discriminant rule is also used to evaluate it. Therefore, it is customary to report test error rate (explained next) which is as easy to calculate and does not require any distributional assumptions. Moreover, confusion matrices are reported for test data (not for training data).

2.4.1 Test accuracy/error rate

The *test error rate* (TER) is computed as follows:

1. *Hold-out*: Randomly split the available data into training set and test set using some split ratio^a
2. Training set is used to construct the classifier \hat{G} and the test set is used to evaluate a performance measure on a test set.^b
3. Repeat steps 1 and 2 (using another random split), e.g., 100 times^c
4. TER is computed as the average of obtained error rates while the test set accuracy is computed as $1 - \text{TER}$.

^ae.g., ratio 2 : 1 implies that $(2/3)^{rd}$ of observations in each class are used to build the training set and the remaining $(1/3)^{rd}$ are left for test set.

^bThe most commonly used performance measure is error rate, so the proportion of observations misclassified in the test set

^cthis is sometimes called repeated K -fold cross validations

In python, you can use `train_test_split` available at scikit-learn package for splitting the data to train and test sets.

Sometimes, the data size may not be sufficient for splitting the data into sufficiently large train and test set. Then less effective cross-validated TER can be computed as follows: the data set \mathcal{T} is split into K disjoint subsets. For each subset of data, say \mathcal{T}_i , train on all but \mathcal{T}_i , then calculate the error rate on data \mathcal{T}_i that was left out. Finally average the obtained errors rates.

Sometimes also stratified sampling is used in step 1 (hold-out): it may be helpful to sample so that proportions of $+/-$ observations is maintained in training/test data splits. This is the case especially when the data is *unbalanced* (i.e., the number of instances in classes vary significantly). Moreover, when data is unbalanced, then it is better to use some other performance measure than error rate, for example, the F1-score.

2.4.2 Other measures

Test error rate is not everything! If a dataset is unbalanced, the overall accuracy is not representative of the true performance of a classifier.

For example, consider an unbalanced binary classification problem with $\#'-' = 990$ and $\#'+ ' = 10$. Then simply labelling everything $' -'$ yields 99% accuracy and has no false alarms. But, this classifier misses all positive class observations, so has 100% miss rate.

Hence it is customary to report also precision and recall (which are often reported as a pair):

- $Recall = \frac{TP}{TP + FN}$ equals true positive rate (TPR)
 $= \text{\#found true objects} / \text{\#all true objects}$

Recall is also sometimes called *Sensitivity*.

- $Precision = \frac{TP}{TP + FP}$ emphasizes TP-s and FP-s
 $= \text{\#found true objects} / \text{\#found all objects}$

Other commonly reported measures are F1-scores and specificity. The *F1-score* is computed as harmonic mean of precision and recall metrics, yielding

$$F1\text{-score} = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}.$$

Given its trade-off nature, the F1-score is quite commonly used and better adapted than the accuracy to unbalanced scenarios. There are also other metrics commonly reported such as

$$Specificity = \frac{TN}{TN + FP}. \quad (\text{equals true negative rate (TNR)})$$

We can generalize above measures to the multi-class case by summarizing over the rows and columns of the confusion matrix. Assuming that the confusion matrix, denoted by $\mathbf{M} = (M_{ij})$ is oriented as above, so each row corresponds to specific value for the "truth", we calculate

$$\text{precision}_i = \frac{M_{ii}}{\sum_j M_{ji}} \quad \text{and} \quad \text{recall}_i = \frac{M_{ii}}{\sum_j M_{ij}}.$$

Thus precision is the fraction of cases where the classifier correctly predicted class i out of all cases for which the algorithm predicted i (correctly and incorrectly). Recall is the fraction of cases where the classifier correctly predicted i out of all of the cases which are labelled as i .

2.5 Conclusions

The purpose of this chapter was to introduce the basic concepts of classification such as risk and the optimal Bayes classifier, different loss functions, and measures of accuracy.

Although the concepts were illustrated in binary classification problem, they generalize in straightforward manner to multi-class case which we omit due to lack of time.

Chapter 3

Decision tree methods

Tree-based methods partition the feature space into a set of rectangles, and then fit a simple model (like a constant) in each one. They are conceptually simple yet powerful. We first describe a popular method for tree-based regression and classification called *CART* [Breiman et al., 1984]. Since the set of splitting rules used to segment the predictor space can be summarized in a tree, these types of approaches are known as *decision-tree* methods. Good reading of these methods are Section 8.7 (Bagging), Section 9.2 (Tree-based methods) and Sections 15.2-15.3 (Random forests) of Hastie et al. [2009].

3.1 Decision tree

Decision trees can be applied to both regression (considered in [Section 3.2](#)) and classification (considered in [Section 3.3](#)) problems. We consider a regression problem with continuous response Y and inputs X_1 and X_2 , each taking values in the unit interval. [Figure 3.1](#) illustrates a simple *recursive binary tree* :

- split the space into 2 regions, and model the response by mean of Y in each region. Choose the variable and split-point to achieve the best fit.
- one or both of these regions are split into two more regions, and this process is continued, until some stopping rule is applied.

The partition in [Figure 3.1](#) is obtained in 3 stages:

1. Split at $X_1 = t_1$.
 2. Split the region $X_1 \leq t_1$ at $X_2 = t_2$ and $X_1 > t_1$ at $X_1 = t_3$.
 3. Split the region $X_2 > t_2$ and $X_1 \leq t_1$ at $X_1 = t_4$.
 4. Split the region $X_1 > t_3$ at $X_2 = t_5$.
- \Rightarrow yields a partition of (X_1, X_2) -space into regions R_1, R_2, \dots, R_6 .

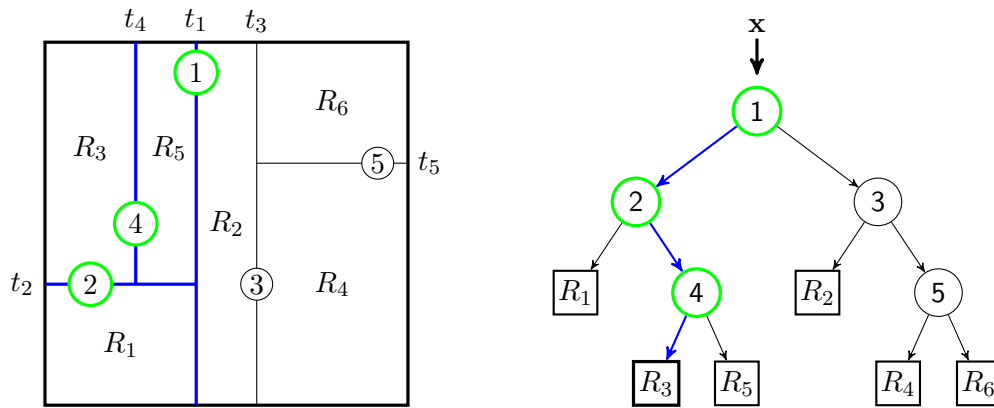


FIGURE 3.1: Left panel shows partition of a 2-dimensional feature space by recursive binary splitting, as used in CART, applied to some fake data. Right panel shows the tree corresponding to the partition. The test data \mathbf{x} fed into the tree belong to the region R_3 .

The corresponding regression model predicts Y with a constant c_m in region R_m , that is,

$$\hat{f}(X) = c_m \sum_{m=1}^6 \mathbf{1}_{\{(X_1, X_2) \in R_m\}}.$$

The *terminal nodes* or *leaves* of the tree correspond to regions R_1, R_2, \dots, R_6 . The points along the tree where the predictor space is split are referred to as *internal nodes*.

3.2 Regression trees

Given the data of inputs and a response, $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$, with $\mathbf{x}_i^\top = (x_{i1}, \dots, x_{ip})$, the algorithm needs to decide:

1. *What is the best split?* (which variables and split points)
2. *How to split the variables*, i.e., what topology (shape) the tree should have? E.g., binary tree or multiway?
3. *How large to grow the tree?*

Suppose first that we have a partition R_1, R_2, \dots, R_M , and we model the response as a constant c_m in each region:

$$f(\mathbf{x}) = c_m \sum_{m=1}^M \mathbf{1}_{\{\mathbf{x} \in R_m\}}.$$

If we adopt as our criterion minimization of sum of squares $\sum (y_i - f(\mathbf{x}_i))^2$, it is easy to see that the best \hat{c}_m is

$$\hat{c}_m = \text{ave}(y_i | \mathbf{x}_i \in R_m) = \text{average of } y_i \text{ in region } R_m.$$

Finding the best binary partition in terms of minimum sum of squares is generally computationally infeasible.

Hence one proceeds with a following *greedy algorithm*:

1. Starting with all of the data, split variable j at a split point s , and define the pair of half-planes:

$$R_1(j, s) = \{\mathbf{x} \mid x_j \leq s\} \quad \text{and} \quad R_2(j, s) = \{\mathbf{x} \mid x_j > s\}$$

2. Determine the best splitting variable j and a split point s that solves

$$\min_{j,s} \left[\min_{c_1} \sum_{\mathbf{x}_i \in R_1(j,s)} (y_i - c_1)^2 + \min_{c_2} \sum_{\mathbf{x}_i \in R_2(j,s)} (y_i - c_2)^2 \right]$$

where for any j and s , the inner minimization is solved by

$$\hat{c}_1 = \text{ave}(y_i \mid \mathbf{x}_i \in R_1(j, s)) \quad \text{and} \quad \hat{c}_2 = \text{ave}(y_i \mid \mathbf{x}_i \in R_2(j, s))$$

3. Having found the best split, we partition the data into the two resulting regions and repeat the splitting process on each of the two regions. This process is continued until stopping criterion is met.

For each splitting variable, the determination of the split point s can be done very quickly and hence by scanning through all of the inputs, determination of the best pair (j, s) is feasible.

Q: What is a good stopping criterion - i.e., **how large we should grow the tree?** How do we know if we should split the tree node?

If we grow a tree deep enough, we end up in massive overfitting as we can usually fit the training data perfectly. The *tree size*, M , i.e., the number of terminal nodes, is a tuning parameter governing the model's complexity, and the optimal tree size should be adaptively chosen from the data, e.g., via pruning.

3.2.1 Pruning the tree

We define a *subtree* $T \subset T_0$ to be any tree that can be obtained by *pruning* T_0 , that is, collapsing any number of its internal (non-terminal) nodes. We index terminal nodes (leaves) by m , with node m representing region R_m . Let $|T|$ denote the number of terminal nodes in T , and define

$$\begin{aligned} N_m &= \#\{\mathbf{x}_i \in R_m\}, & (= \text{node size}) \\ \hat{c}_m &= \frac{1}{N_m} \sum_{\mathbf{x}_i \in R_m} y_i, \\ Q_m(T) &= \frac{1}{N_m} \sum_{\mathbf{x}_i \in R_m} (y_i - \hat{c}_m)^2. \end{aligned} \tag{3.1}$$

The *cost complexity criterion* is then defined as

$$C_\alpha(T) = \sum_{m=1}^{|T|} N_m Q_m(T) + \alpha |T|,$$

where $\alpha \geq 0$ is a tuning parameter that governs the tradeoff between tree size and its goodness of fit to the data:

- Large α results in smaller trees T (and vice versa).
- When $\alpha = 0$ the solution is the full tree T_0 .

The **cost-complexity pruning** proceeds as follows:

1. Grow a large tree T_0 , stopping the splitting process only when some minimum node size (say 5) is reached.
2. For each α , find the subtree T_α that minimizes $C_\alpha(T)$:

$$T_\alpha = \arg \min_{T \subseteq T_0} \sum_{m=1}^{|T|} N_m Q_m(T) + \alpha |T|,$$

using *weakest-link pruning*¹:

- 2a) successively collapse the internal node that produces the smallest per-node increase in $\sum_m N_m Q_m(T)$
 - 2b) continue until one reaches the single-node (root) tree
 - 2c) this produces a sequence of subtrees which contains T_α .
3. Estimation of α is achieved by five- or tenfold cross-validation: $\hat{\alpha}$ minimizes the cross-validated sum of squares.
 4. Final tree is $T_{\hat{\alpha}}$.

3.3 Classification trees

If the task is classification, so $Y \in \{1, 2, \dots, K\}$, the only changes needed in the tree algorithm pertain to the criteria for splitting nodes and pruning the tree.

It is obvious that the node impurity measure $Q_m(T)$ in (3.1) needs to be redefined. Suppose there are N_m observations in a node m that represents a region R_m . Let $(\hat{p}_{m,1}, \dots, \hat{p}_{m,K})$ denote the frequencies, that is,

$$\hat{p}_{m,k} = \frac{1}{N_m} \sum_{\mathbf{x}_i \in R_m} 1_{\{y_i=k\}} := \begin{array}{l} \text{proportion of class } k \\ \text{observations in node } m \end{array}.$$

We classify the observations in node m to class

$$k(m) = \arg \max_k \hat{p}_{m,k},$$

i.e., the majority class in node m .

Commonly used options for node impurity are listed in Table 3.1. Notice that minimizing the Gini-index will favour pure nodes, which is why it often the favoured option

¹For each α one can show that there is a unique smallest subtree T that minimizes $C_\alpha(T)$.

Misclassification:	$\text{Err}(m) = 1 - \hat{p}_{m,k(m)}$
Gini index:	$\text{Gini}(m) = \sum_{k \neq k'} \hat{p}_{m,k} \hat{p}_{m,k'} = 1 - \sum_{k=1}^K \hat{p}_{m,k}^2$ - maximized when $\hat{p}_{m,k} = 1/K$ with value $1 - 1/K$ - minimized when all cases belong to a single class.
Entropy/deviance:	$H(m) = - \sum_{k=1}^K \hat{p}_{m,k} \log \hat{p}_{m,k}$ - maximized when $\hat{p}_{m,k} = 1/K$ with value $\log K$. - minimized when one class has no cases in it.

TABLE 3.1: Different measures $Q_m(T)$ of node impurity

for node impurity. For two classes, if p is the proportion in the second class, these three measures are:

$$1 - \max(p, 1 - p), \quad 2p(1 - p) \quad \text{and} \quad -p \log p - (1 - p) \log(1 - p),$$

respectively. They are shown in Figure 3.2. Either Gini index $\text{Gini}(\cdot)$ or entropy $H(\cdot)$ are favoured when growing the tree:

- they are differentiable (better suited for numerical optimization)
- they are more sensitive to changes in the node probabilities.

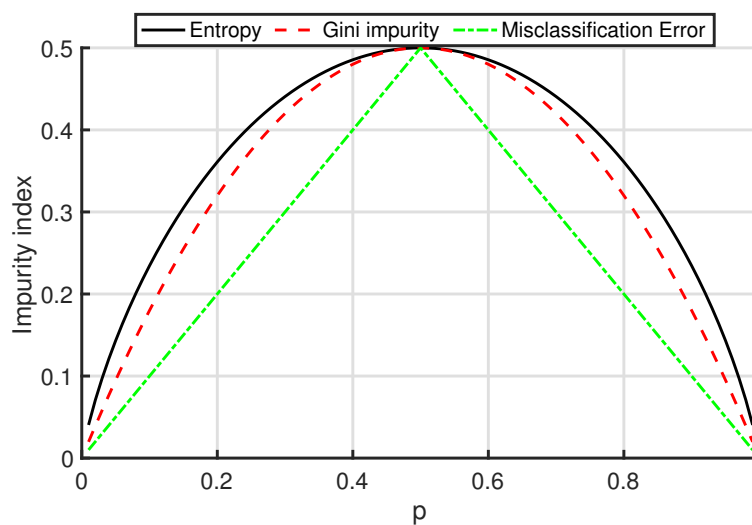


FIGURE 3.2: Node impurity measures in the binary classification problem. The cross-entropy has been scaled by $1/2$.

3.3.1 Choosing the variable to split on and the split point

Suppose the parent node m is split into J new child nodes m_j , $j = 1, \dots, J$. Each child has a count n_j and let

$$(\hat{p}_{m_j,1}, \dots, \hat{p}_{m_j,K}), \quad j = 1, \dots, J$$

denote the vector of class frequencies at child node m_j . The parent node has then a count $n = \sum_j n_j$ cases and the *gain* in Gini index is computed as

$$\text{Gain}(G, m \rightarrow (m_j)_1^J) = G(m) - \sum_{j=1}^J \frac{n_j}{n} G(m_j).$$

(Similarly we can compute the gain for $H(m)$ or $\text{Err}(m)$.)

One can then choose the best splitting node as the one that produces the biggest gain in G . Similarly a continuous variable can be discretized and then choose the splitting point as the point in the grid that produces the largest gain.

3.4 Bagging

A binary decision tree takes a majority vote within each cell of a partition of the feature space that has appearance as illustrated in Figure 3.1. Since the partition of the feature space is rough (highly non-smooth regardless of how large the tree is grown) and highly unstable/non-robust (i.e., the partition is learned with a lot of variance and a small change in the inputs can lead to large change in the output).

Ensemble methods aim to tackle these deficiencies of the decision trees. An *ensemble* (or committee) of classifiers is a classifier build upon some combination of learners. An ensemble method combines the following simple steps

1. Generate a number of classifiers $G_1(\cdot), \dots, G_M(\cdot)$ based on the same data set but typically using some degree of randomness (e.g., in the selection of the features, or in the selection of the observations from the full data set).
2. Aggregate the classifiers to make the final prediction (take a majority vote).
3. In the case of ties (e.g., some classes have equal number of votes), then take a random pick of the classes.

Even though none of the classifier performs particularly well at its own, the aggregated result can be good. This is due to the wisdom of the masses.

Bagging aka *bootstrap aggregating* introduces randomness via bootstrap samples. A *bootstrap sample* $\mathcal{T}^* = \{(y_i^*, \mathbf{x}_i^*)\}_{i=1}^N$ is formed from the original sample by resampling **with replacement** from the original training data $\mathcal{T} = \{(y_i, \mathbf{x}_i)\}_{i=1}^N$. Suppose we have a learning algorithm (e.g., a tree), and suppose we have generated \mathcal{T}_b^* , for $b = 1, \dots, B$ bootstrap samples. Let G_b be the classifier when applied to the bootstrap sample \mathcal{T}_b^* . The predicted class of bagging classifier is then a majority vote over all base classifiers G_1, \dots, G_B or the one with highest mean probability estimate across the base classifiers. The latter option is nowadays the default choice while majority voting is used only when the base classifier does not provide the class probability estimates. Thus one computes

regression: $\hat{f}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B G_b(\mathbf{x})$ is the prediction of $\mathbf{x} \in \mathbb{R}^p$.

classification: Predicted class label for $\mathbf{x} \in \mathbb{R}^p$ is determined using majority voting:

$$\hat{G}(\mathbf{x}) = \arg \max_k \sum_{b=1}^B 1_{\{G_b(\mathbf{x})=k\}}$$

or mean probability:

$$\hat{G}(\mathbf{x}) = \arg \max_k \sum_{b=1}^B \hat{\text{Pr}}_b(Y = k | X = \mathbf{x})$$

where $p^{(b)}(\mathbf{x}) = \hat{\text{Pr}}_b(Y = k | X = \mathbf{x})$ is the probability prediction for k th class obtained by $G_b(\cdot)$.

Algorithm 3.1 describes the Bagging procedure when base classifier is a tree. Note that the algorithm is general, and not in any means restricted to using trees. Basically, you can replace tree by any classifier.

Algorithm 3.1: Bagging algorithm for classification or regression using trees

Input : $\mathcal{T} = \{(y_i, \mathbf{x}_i^\top)\}_{i=1}^N$ (training data), B (number of bootstrap samples)

Output : Ensemble of trees $\{T_b\}_{b=1}^B$

```

1  for  $b = 1$  to  $B$  do
2    Draw a bootstrap sample  $\mathcal{T}_b^*$  of size  $N$  from the training data.
3    Build a tree  $T_b(\cdot)$  on  $\mathcal{T}_b^*$ 
    // To make a prediction at a new point  $\mathbf{x}$ :
4  if classification task then
5    if not majority voting then
6       $\hat{G}(\mathbf{x}) = k$  having largest mean probability  $\sum_{b=1}^B p_k^{(b)}(\mathbf{x})$ 
      // above  $p_k^{(b)}(\mathbf{x}) = \hat{\text{Pr}}_b(Y = k | X = \mathbf{x})$  based on  $b$ th tree  $T_b(\cdot)$ 
7    else
8       $\hat{G}(\mathbf{x}) =$  majority vote over  $\{T_b(\mathbf{x})\}_{b=1}^B$ 
9  else
10    $\hat{f}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B T_b(\mathbf{x})$ .
```

An important feature of bagging (as well as of random forests) is its use of *out-of-bag* (OOB) samples: We can use cases not included in the Bootstrap sample (out-of-bag cases) as a test set. For each observation $\mathbf{z}_i = (\mathbf{x}_i, y_i)$, construct its bagging prediction by averaging only those trees which are formed using bootstrap samples in which \mathbf{z}_i did not appear. Here averaging means either averaging the probability predictions or the class predictions in the case of majority voting. If an observation has no votes yet² one assigns the observation randomly to one of the classes. Ties can occur (namely, the case that two or more of the classes have same number of votes) when majority voting is used; in such case one again assigns the observation randomly to one of the classes. An OOB error estimate is almost identical to that obtained by N-fold cross-validation or to test error. Once the OOB error stabilizes, the training can be terminated. Thus OOB provides a simple means for choosing a good choice for B .

²this can happen only in the beginning of iterations. When iterations are ten or more, all observations usually have at least one prediction.

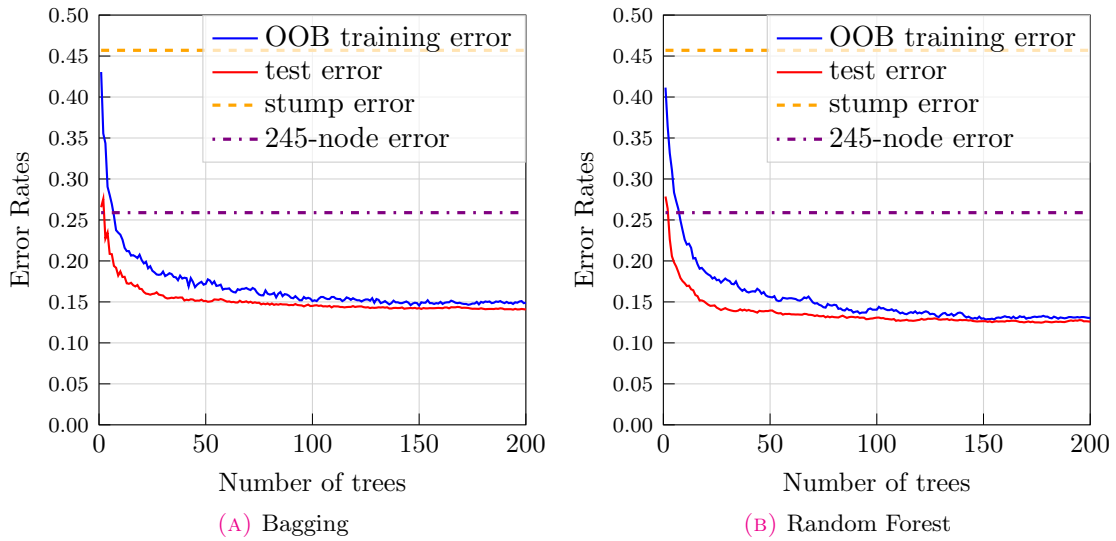


FIGURE 3.3: Results for simulated data of [Example 3.1](#). Error rates of (a) Bagging and (b) Random Forest (using $d = 2$ and $n_{min} = 3$) as a function of the number of trees. Also shown are the test error rate for a single stump, and a 245-node classification tree.

Example 3.1. We consider the case that the features X_1, \dots, X_{10} are standard independent Gaussian, and the deterministic target Y is defined by

$$Y = \begin{cases} 1 & , \text{if } \sum_{j=1}^{10} X_j^2 > \chi_{10}^2(0.5) \\ -1 & , \text{otherwise} \end{cases},$$

where $\chi_{10}^2(0.5)$ is the median of a chi-squared random variable with 10 degrees of freedom (sum of squares of 10 standard Gaussians). We generate 2000 training observations $\{\mathbf{x}_i, y_i\}_{i=1}^N$ and 10,000 test observations. Hence we will have approximately 1000 (resp. 5000) cases in each class in the training (resp. test) data set.

Figure 3.3a illustrates the power of Bagging with threes. Applying stump (2 node tree) alone to the training data set yields a very poor test error rate of 46.0%, compared to 50.0% achieved by random guessing. Bagging achieves 14.05% test error rate which is significantly better than a single large classification tree (error rate 25.89%). Here we used the mean probability for class prediction. Figure shows the OOB misclassification error compared to the test error. Although 200 trees are averaged here, it appears from the plot that about 100 would be sufficient. The confusion matrix is shown in [Figure 3.4](#). ■

3.5 Random forests

A random forest is an ensemble of decision trees where each decision tree is independently randomized. Bagging with decision trees is a simple example of a random forest.

Random forest as such suffer from the fact that bootstrap samples are highly correlated. As a consequence, different trees tend to select the same features and lead to

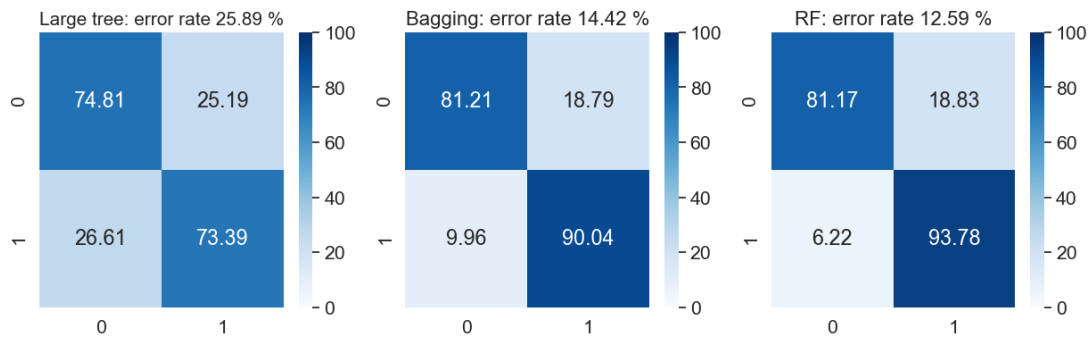


FIGURE 3.4: Confusion matrices of Example 3.1 and Example 3.2. Large tree refers to a large tree with 245 nodes.

(too) similar partitions of the feature space. Ideally, one would like to have partitions that are more independent. A simple cure is to incorporate *random feature selection*:

1. Generate classifiers by choosing random subsets of features and construct a decision tree on just the selected features. Random subsets are considered at each internal node so that the search space for each split is smaller.
2. Combine the approach with bagging.

Since the obtained partitions are less correlated, the obtained final aggregate prediction has reduced variance. Due to the reduced space at each split, the individual trees are built much faster than in bagging. As a rule of thumb one can use $d = \sqrt{p}$ features. The developer of random forests, Leo Breimann, called the random forests as the best "off-the-shelf" method for classification.³ The algorithm give in Algorithm 3.2 follows the original publication Breiman [2001] but slight modifications to it are possible. Important parameters are d , i.e., the number of randomized features in each split) and n_{min} , i.e., minimum node size for ending splitting of nodes.

Example 3.2. We consider the same data set as in Example 3.1 and Figure 3.3b illustrates the results of random forest with threes. The parameters used were $d = 2$ (number of features used for splitting) and $n_{min} = 3$ (minimum node size) and $B = 200$. Random forest achieves error rate of 12.40% with 200 bags while bagging alone achieved 14.05% test error rate. Confusion matrix is shown in Figure 3.4.

Figure also shows the OOB training error compared to the test error. Although 200 trees are averaged, it appears that about 100 may suffice. The confusion matrices are shown in Figure 3.4b. ■

³An "off-the-shelf" method is one that can be directly applied to the data without requiring a great deal of time consuming data preprocessing or careful tuning of the learning procedure.

Algorithm 3.2: RandomForest algorithm for classification or regression

Input : $\mathcal{T} = \{(y_i, \mathbf{x}_i)\}_{i=1}^N$ (training data), B (number of bootstrap samples), d (number of features in each split), n_{min} (minimum node size)

Output : Ensemble of trees $\{T_b\}_{b=1}^B$

```

1  for  $b = 1$  to  $B$  do
2    Draw a bootstrap sample  $\mathcal{T}_b^*$  of size  $N$  from the training data.
    // Build a tree  $T_b(\cdot)$  on  $\mathcal{T}_b^*$ 
3    for each terminal node in tree do
4      Select  $d$  variables at random from the  $p$  variables.
5      Pick the best variable/split-point among the  $d$  variables.
6      Split the node into two daughter nodes on the selected feature until the
      minimum node size  $n_{min}$  is reached

    // To make a prediction at a new point  $x$ :
7    if classification task then
8      if not majority voting then
9         $\hat{G}(\mathbf{x}) = k$  having largest mean probability  $\sum_{b=1}^B p_k^{(b)}(\mathbf{x})$ 
        // above  $p_k^{(b)}(\mathbf{x}) = \hat{\text{Pr}}_b(Y = k \mid X = \mathbf{x})$  based on  $b$ th tree  $T_b(\cdot)$ 
10     else
11        $\hat{G}(\mathbf{x}) =$  majority vote over  $\{T_b(\mathbf{x})\}_{b=1}^B$ 
12  else
13     $\hat{f}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B T_b(\mathbf{x})$ .
```

Chapter 4

Boosting

Boosting is one of the most important recent developments in classification methodology. Boosting is an ensemble method that combines the outputs of many "weak" classifiers (such as simple trees) to produce a more powerful *committee* vote. A weak learner is a learning algorithm capable of producing classifiers with probability of error strictly (but not considerably) smaller than that of random guessing. A strong learner on the other hand is such that it is able to yield classifiers with arbitrarily small error probability given a sufficient amount of training data. Boosting is a fundamental concept in statistical learning as the idea is rather simple and can be modified and extended to many problems.

The first simple boosting procedure was developed by Schapire [1990] who showed that a weak learner could always improve its performance by training two additional classifiers on filtered versions of the input data stream. Later Freund and Schapire [1996, 1997] developed more adaptive and realistic AdaBoost, short for Adaptive Boosting, algorithm that combines many weak learners simultaneously. The techniques provably improved or "boosted" the performance of a single classifier by producing a majority vote of similar classifiers. The developers, Yoav Freund and Robert Schapire, won the 2003 Gödel Prize for their work on AdaBoost.

4.1 General ensemble scheme

Consider a weak learner or base procedure which constructs a function $G(\mathbf{x})$ based on input data:

$$\text{input } \{(\mathbf{x}_i, y_i)\}_{i=1}^N \longrightarrow \boxed{\text{weak learner}} \longrightarrow G(\cdot)$$

Boosting applies sequentially the weak classification/regression algorithm to repeatedly modified versions of the data. As an output one gets a sequence of weak learners $\{G_m(\mathbf{x})\}_{m=1}^M$.

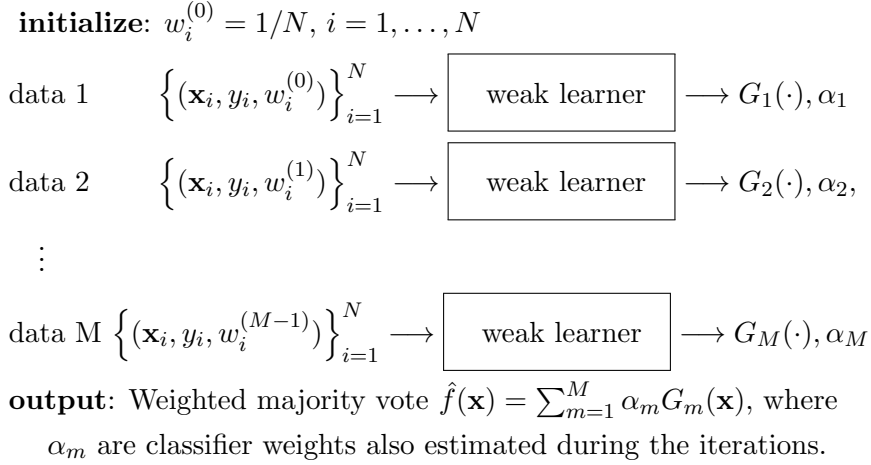
The data modifications at each boosting step consist of applying weights w_1, \dots, w_N to each of the training observations $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$. Boosting can be applied to a regres-

sion or classification algorithm that accepts case weights, e.g.,

$$\underset{f}{\text{minimize}} \sum_{i=1}^N w_i L(y_i, f(\mathbf{x}_i)).$$

The predictions from all of them are combined through a weighted majority vote.

General boosting ensemble scheme:



For example, in the case of binary classification problem, the final prediction is

$$G(\mathbf{x}) = \text{sign}(\hat{f}(\mathbf{x})) = \text{sign} \left(\sum_{m=1}^M \alpha_m G_m(\mathbf{x}) \right), \quad (4.1)$$

where the *classifier weights*, $\alpha_1, \dots, \alpha_M$, computed by the boosting algorithm, are designed so that more accurate classifier $G_m(\mathbf{x})$ in the sequence obtains a larger weight and thus has a higher influence on $G(\mathbf{x})$. The *data weights* $w_i^{(m)}$, $i = 1, \dots, N$ at each boosting iteration m depends on the accuracy of the previous classifiers, allowing the algorithm to focus its attention on those samples that are still incorrectly classified.

4.2 AdaBoost

The AdaBoost algorithm is the most well known boosting algorithm and it was originally developed for binary classification problem. Here the base classifier $G(\mathbf{x})$ attains values in $\{-1, 1\}$. [Algorithm 4.1](#) shows the details of the *AdaBoost.M1* algorithm also known as the *discrete AdaBoost* [Friedman et al., 2000]. It is called discrete since the base classifier $G_m(\mathbf{x})$ returns a discrete class label.

The most important tuning parameter of AdaBoost is the number of weak learners M (or iterations) used in the aggregation. AdaBoost and other boosting procedures are quite resistant to overfitting when increasing the number of iterations M . This has been observed empirically, and is illustrated in [Example 4.1](#) below. Nevertheless,

Algorithm 4.1: AdaBoost.M1 (alias Discrete AdaBoost) for binary classification

Initialize: $w_i^{(0)} = 1/N$, $i = 1, \dots, N$

1 **for** $m = 1$ **to** M **do**

2 Fit a classifier $G_m(\mathbf{x}) \in \{-1, 1\}$ to the training data using weights $w_i^{(m-1)}$.

3 Compute the weighted error rate

$$\text{err}_m = \frac{\sum_{i=1}^N w_i^{(m-1)} \mathbf{1}_{\{y_i \neq G_m(\mathbf{x}_i)\}}}{\sum_{i=1}^N w_i^{(m-1)}}$$

4 Compute $\alpha_m = \log((1 - \text{err}_m)/\text{err}_m)$.

5 Update the weights $w_i^{(m)} = w_i^{(m-1)} \cdot \exp[\alpha_m \cdot \mathbf{1}_{\{y_i \neq G_m(\mathbf{x}_i)\}}]$, $i = 1, \dots, N$.

Output : $G(\mathbf{x}) = \text{sign}(\hat{f}(\mathbf{x}))$ and $\hat{f}(\mathbf{x}) = \sum_{m=1}^M \alpha_m G_m(\mathbf{x})$

clear overfitting does occur for some datasets. Although early stopping is not necessary, AdaBoost and also other boosting algorithms are overfitting eventually, and early stopping, i.e., choosing a valid number of iterations M is necessary.

Some key points of AdaBoost.M1 are:

- Initially, $w_i = 1/N$, so at first step one trains the classifier on the data in the usual manner.
- For each successive iteration $m = 2, 3, \dots, M$ the observation weights are individually modified and the classification algorithm is reapplied to the weighted observations.
- At step m , observations that were misclassified by weak learner $G_{m-1}(\cdot)$ induced at the previous step have their weights *increased*.
- Thus, as iterations proceed, observations that are difficult to classify correctly receive ever-increasing influence.

A generalization of AdaBoost.M1 was proposed in Freund and Schapire [1996], and was further studied by Schapire and Singer [1999], that uses real-valued confidence-rated predictions rather than the $\{-1, 1\}$ of Discrete AdaBoost. Algorithm 4.2 presents this more general version of AdaBoost, referred to as *AdaBoost.R* (or *real AdaBoost*) in which the weak learner returns a class probability estimate $p^{(m)}(\mathbf{x}) = \hat{\text{Pr}}_{\mathbf{w}}(Y = 1 | \mathbf{x}) \in [0, 1]$ at m^{th} iteration. Its contribution to the final classifier is one half the logit-transform of this probability estimate. In real AdaBoost, the base classifier returns a real-valued prediction, namely, a probability estimate of the class. Thus AdaBoost.R can use classifiers that can compute class prediction probability. For example, the predicted class probability for a decision tree is the fraction of samples of the same class in the terminal leaf.

The default algorithm in scikit-learn's `AdaBoostClassifier` is AdaBoost.R. In case you wish to use AdaBoost.M1, then choose option `algorithm='SAMME'`.

Example 4.1. Consider the same test and training data set as in Example 3.1. The performance of AdaBoost when the weak classifier is a stump (two terminal node

Algorithm 4.2: AdaBoost.R (alias Real AdaBoost)

Initialize: $w_i^{(0)} = 1/N$, $i = 1, \dots, N$

- 1 **for** $m = 1$ **to** M **do**
- 2 Fit a classifier $G_m(\mathbf{x})$ to obtain a class probability estimate
 $p^{(m)}(\mathbf{x}) = \hat{\text{Pr}}_{\mathbf{w}}(Y = 1|\mathbf{x}) \in [0, 1]$, using weights $w_i^{(m-1)}$ on the training data
- 3 Compute $f_m(\mathbf{x}) = \frac{1}{2} \log p^{(m)}(\mathbf{x}) / (1 - p^{(m)}(\mathbf{x})) \in \mathbb{R}$.
- 4 Compute $w_i^{(m)} = w_i^{(m-1)} \exp\{-y_i f_m(\mathbf{x}_i)\}$, $i = 1, 2, \dots, N$ and renormalize so
that $\sum_i w_i^{(m)} = 1$.

Output : $G(\mathbf{x}) = \text{sign}(\hat{f}(\mathbf{x}))$ and $\hat{f}(\mathbf{x}) = \sum_{m=1}^M f_m(\mathbf{x})$

classification tree) is shown in Figure 4.1a,b. Applying stump alone to the training data set yields a very poor test error rate of 46.0%, compared to 50.0% achieved by random guessing. However, as boosting iterations proceed the error rate steadily decreases, reaching 10.25% after $M = 600$ iterations when using AdaBoost.M1 and 5.63% using AdaBoost.R. Hence, boosting this weak classifier reduces its prediction error rate by a factor of 5 and 9, respectively. Both methods also outperform a single large classification tree (error rate 24.55%).

Figure 4.1c,d shows the performance when the weak learner is an 8-node decision tree. Initially, error rates for boosting eight-node trees decrease much more rapidly than for stumps. However, the error rates quickly level off and improvement is very slow after about 100 iterations. The error rate after 600 iteration is 6.86% and 7.19%, respectively. If we allow the AdaBoost.M1 with stumps to continue iterations, we can notice that after roughly 2000 iterations, it has achieved the same test error rate as AdaBoost.M1 using 8-node trees with $M = 600$ iterations. We also notice that (for this data set) boosting outperforms both random forest and bagging.

The confusion matrix of Adaboost based on stumps and $M = 600$ iterations are shown in Figure 4.2. These can be compared with Figure 3.4. ■

4.3 Forward Stagewise Additive Modeling

Consider a *basis function expansion* of the form

$$f(\mathbf{x}) = \sum_{m=1}^M \beta_m b(\mathbf{x}; \gamma_m)$$

where β_m , $m = 1, 2, \dots, M$ are the expansion coefficients and $b(\mathbf{x}; \gamma) \in \mathbb{R}$ are (usually simple) "basis" functions of the multivariate argument \mathbf{x} , characterized by a set of parameters γ .

Typically these models are fit by minimizing empirical risk

$$\underset{\{\beta_m, \gamma_m\}_1^M}{\text{minimize}} \sum_{i=1}^N L\left(y_i, \sum_{m=1}^M \beta_m b(\mathbf{x}_i; \gamma_m)\right). \quad (4.2)$$

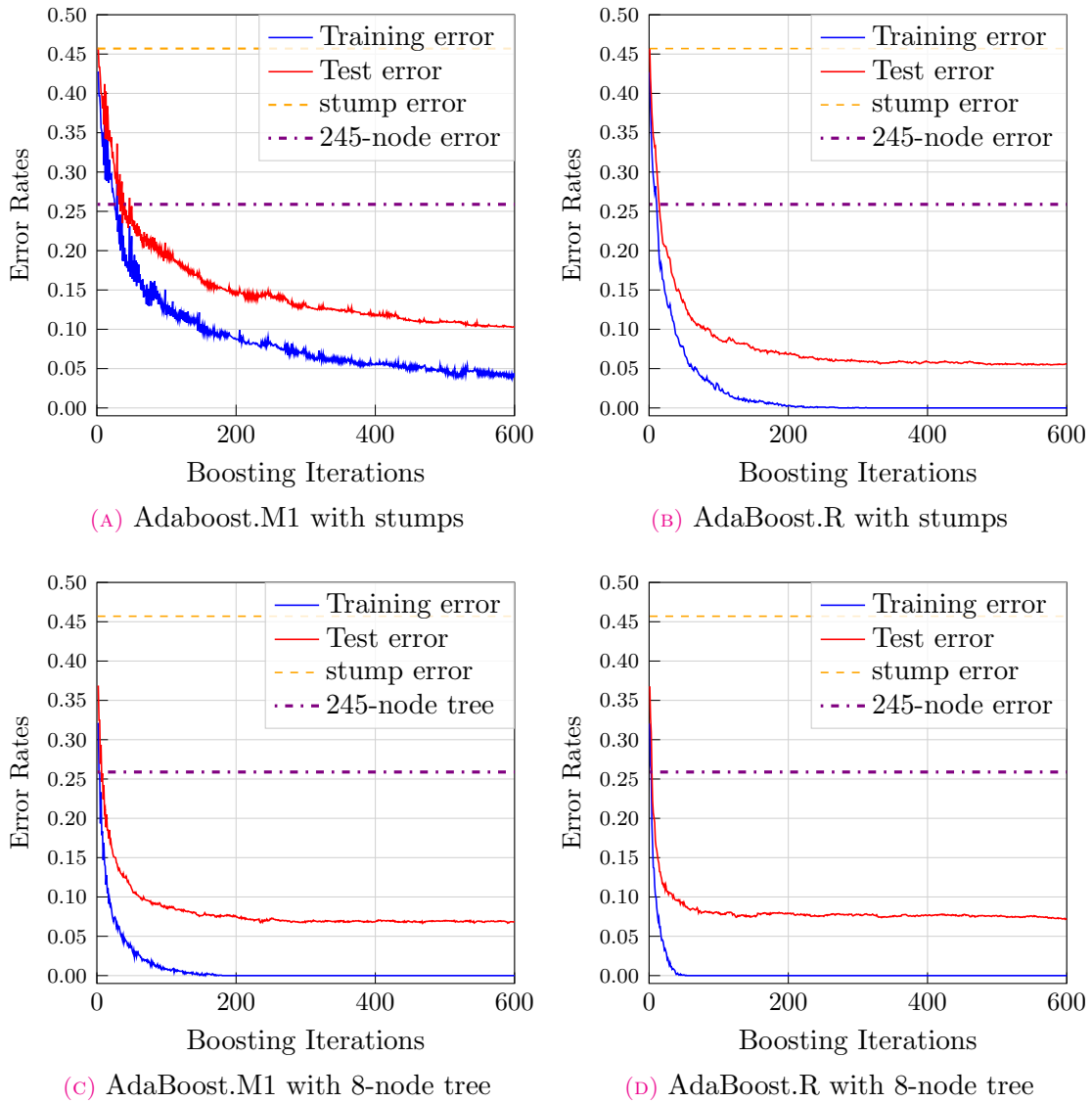


FIGURE 4.1: Results for simulated data of *Example 4.1*. Error rates as a function of the number of boosting iterations when using AdaBoost.M1 and AdaBoost.R when the weak learner is (a),(b) a stump and (c),(d) eight-node trees. Also shown are the test error rate for a single stump, and a 244-node classification tree

Solving this problem is often computationally infeasible. However, a simple approximate solution can be developed when it is feasible to rapidly solve the subproblem of fitting just a single basis function:

$$\underset{\beta, \gamma}{\text{minimize}} \sum_{i=1}^N L(y_i, \beta b(\mathbf{x}_i; \gamma)).$$

This is the idea in **forward stagewise modeling**, which *approximates* the solution to (4.2) by sequentially adding new basis functions to the expansion without adjusting the parameters and coefficients of those that have already been added.

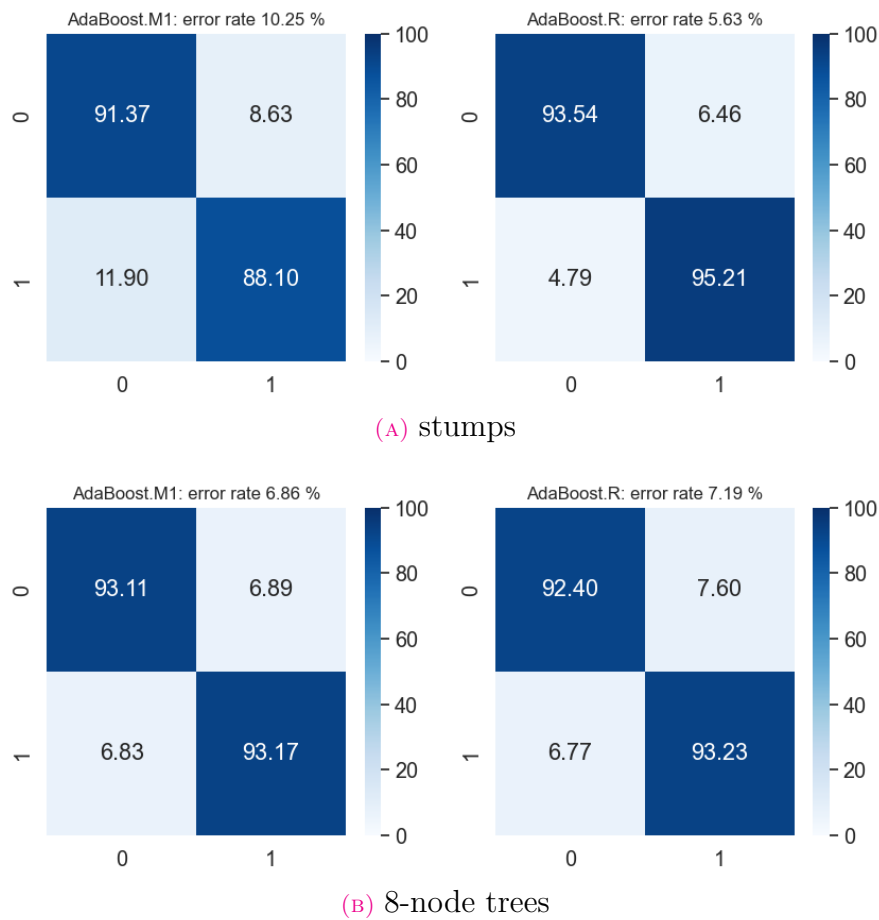


FIGURE 4.2: Confusion matrices for AdaBoost.M and .R with stumps and 8 node trees ($M = 600$).

Algorithm 4.3: Forward Stagewise Additive Modeling

Initialize: $f_0(\mathbf{x}) = 0$

1 **for** $m = 1$ **to** M **do**

2 Compute

$$(\beta_m, \gamma_m) = \arg \min_{\beta, \gamma} \sum_{i=1}^N L(y_i, f_{m-1}(\mathbf{x}_i) + \beta b(\mathbf{x}_i; \gamma))$$

3 Set $f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \beta_m b(\mathbf{x}; \gamma_m)$.

This procedure, detailed in [Algorithm 4.3](#), proceeds as follows:

- At each iteration m , one solves for the optimal basis function $b(\mathbf{x}; \gamma_m)$ and corresponding coefficient γ_m to add to the current expansion $f_{m-1}(\mathbf{x})$.
- This produces $f_m(\mathbf{x})$, and the process is repeated. Previously added terms are not modified.

For squared-error loss one has

$$\begin{aligned} L(y_i, f_{m-1}(\mathbf{x}_i) + \beta b(\mathbf{x}_i; \gamma)) &= (\underbrace{y_i - f_{m-1}(\mathbf{x}_i)}_{= r_{im}} - \beta b(\mathbf{x}_i; \gamma))^2 \\ &= (r_{im} - \beta b(\mathbf{x}_i; \gamma))^2. \end{aligned} \quad (4.3)$$

Thus, for squared-error loss, the term $\beta_m b(\mathbf{x}; \gamma_m)$ that best fits the current residuals $\{r_{im}\}_{i=1}^N$ is added to the expansion at each step.

4.4 Exponential Loss + Stagewise additive modeling = AdaBoost.M1

Next we show that in binary classification problem, AdaBoost.M1 ([Algorithm 4.1](#)) is equivalent to forward stagewise additive modeling ([Algorithm 4.3](#)) that uses the exponential loss

$$L(y, f(\mathbf{x})) = \exp(-yf(\mathbf{x})) \quad (4.4)$$

and where the individual classifiers $G_m(\mathbf{x}) \in \{-1, 1\}$ take the role of the basis functions $b(\mathbf{x}; \gamma)$.

When $L(\cdot, \cdot)$ is the exponential loss in (4.4), the classifier G_m and corresponding coefficient β_m to be added at each step solve

$$\begin{aligned} (\beta_m, G_m) &= \arg \min_{\beta, G} \sum_{i=1}^N \exp\{-y_i(f_{m-1}(\mathbf{x}_i) + \beta G(\mathbf{x}_i))\} \\ &= \arg \min_{\beta, G} \sum_{i=1}^N w_i^{(m)} \exp(-\beta y_i G(\mathbf{x}_i)), \end{aligned} \quad (4.5)$$

where

$$w_i^{(m)} = \exp(-y_i f_{m-1}(\mathbf{x}_i)) \quad (4.6)$$

can be regarded as a weight that is applied to each observation.

The solution to (4.5) can be obtained in two steps. First, keeping β fixed, we can solve for G :

$$G_m = \arg \min_G \sum_{i=1}^N w_i^{(m)} \mathbf{1}_{\{y_i \neq G(\mathbf{x}_i)\}} \quad (4.7)$$

which is the classifier that minimizes the weighted error rate in predicting y . This follows by noting that the criterion in (4.5) may be written as

$$\begin{aligned} & \sum_{i=1}^N w_i^{(m)} \exp(-\beta y_i G(\mathbf{x}_i)) \\ &= e^{-\beta} \sum_{y_i=G(\mathbf{x}_i)} w_i^{(m)} + e^{\beta} \sum_{y_i \neq G(\mathbf{x}_i)} w_i^{(m)} \\ &= (e^{\beta} - e^{-\beta}) \sum_{i=1}^N w_i^{(m)} \mathbf{1}_{\{y_i \neq G(\mathbf{x}_i)\}} + e^{-\beta} \sum_{i=1}^N w_i^{(m)}. \end{aligned}$$

Then, after we plug in this G_m into (4.5) and solve for β we obtain

$$\beta_m = \frac{1}{2} \log \frac{1 - \text{err}_m}{\text{err}_m}, \quad (4.8)$$

where

$$\text{err}_m = \frac{\sum_{i=1}^N w_i^{(m)} \mathbf{1}_{\{y_i \neq G_m(\mathbf{x}_i)\}}}{\sum_{i=1}^N w_i^{(m)}} \quad (4.9)$$

is the minimized weighted error rate (cf. Line 3 of Algorithm 4.1).

The approximation is then updated

$$f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \beta_m G_m(\mathbf{x})$$

which causes the weights in (4.6) for the next iteration to be

$$\begin{aligned} w_i^{(m+1)} &= w_i^{(m)} \cdot e^{-\beta_m y_i G_m(\mathbf{x}_i)} \\ &= w_i^{(m)} \cdot e^{\alpha_m \mathbf{1}_{\{y_i \neq G_m(\mathbf{x}_i)\}}} \cdot e^{-\beta_m} \end{aligned} \quad (4.10)$$

where $\alpha_m = 2\beta_m$ is exactly the quantity defined at line 4 of Algorithm 4.1. In obtaining (4.10) we used that

$$-y G(\mathbf{x}) = 2 \cdot \mathbf{1}_{\{y \neq G(\mathbf{x})\}} - 1.$$

The constant $e^{-\beta_m}$ in (4.10) multiplies all weights by the same value, so it can be ignored. Thus (4.10) is exactly equivalent to line 5 of Algorithm 4.1.

Remark 4.1. Recall from Table 2.1 that the optimum predictor function $f^*(\mathbf{x})$ that minimizes the exponential loss is one half of log odds. Thus we can interpret the additive expansion $\hat{f}(\mathbf{x}) = \sum_{m=1}^M \alpha_m G_m(\mathbf{x})$ of AdaBoost as an estimate of log-odds of $p(\mathbf{x}) = \Pr(Y = 1|X = \mathbf{x})$. Note that it is not $\frac{1}{2} \times$ times log-odds since the actual

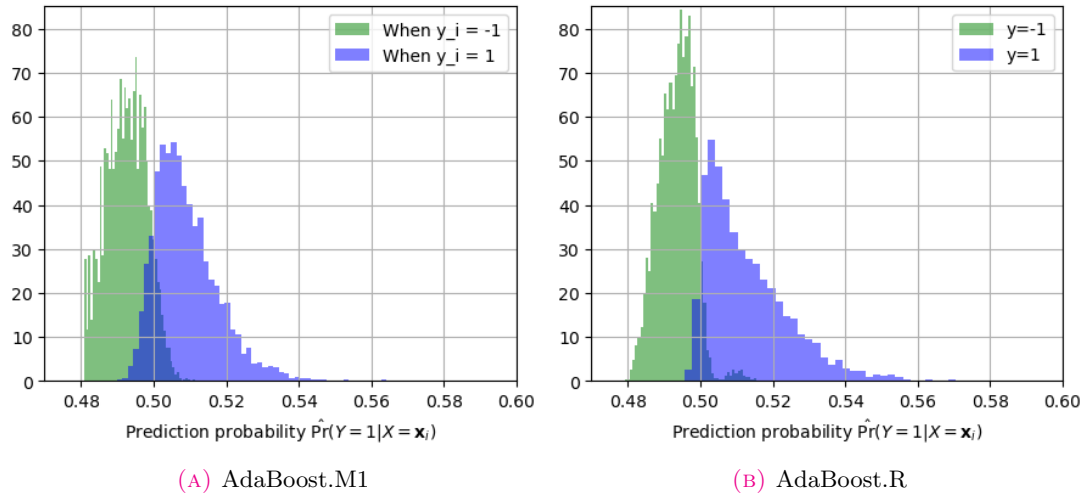


FIGURE 4.3: Results for simulated data of [Example 4.1](#). Density histograms of prediction probabilities $\hat{p}(\mathbf{x}_i) = \Pr(Y = 1 \mid X = \mathbf{x}_i)$ for inputs \mathbf{x}_i in the test data. The density on light blue (resp. green) are test data cases with $y_i = 1$ (resp. $y_i = -1$).

multiplier is $\beta_m = \alpha_m/2$ and the scaling by $1/2$ was ignored by convenience. This justifies using its sign as the classification rule in AdaBoost.M1. Moreover, due to (2.22) it allows us to define probability estimates $\hat{p}(\mathbf{x}) = \Pr(Y = 1 \mid X = \mathbf{x})$ simply as

$$\hat{p}(\mathbf{x}) = \sigma(\hat{f}(\mathbf{x})) = \frac{1}{1 + e^{-\hat{f}(\mathbf{x})}}$$

where

$$\hat{f}(\mathbf{x}) = \sum_{m=1}^M \alpha_m G_m(\mathbf{x}) / \sum_{m=1}^M \alpha_m$$

is weighted mean predictions in the ensemble. This is what `.predict_proba` returns when using the `AdaBoostClassifier` with option `algorithm='SAMME'`. [Figure 4.3](#) shows the normalized histograms of predictions probabilities $\hat{p}(\mathbf{x})$ for test data. The density on light blue (resp. red) are observations with $y_i = 1$ (resp. $y_i = -1$). These we would like to have prediction probabilities > 0.5 (resp. < 0.5). Thus we would like the histogram to be located above (resp. below) 0.5 vertical line. If the histograms do not overlap, then one has 0% test error. As can be noted, AdaBoost.M1 has more unwanted overlaps than AdaBoost.R.

4.5 Discussion

We only touched a surface of vast literature on boosting in this chapter. There are tons of extensions of the boosting principle. Here we mention only three most well-known methods that your future boss may expect you to know when you apply to data / ML scientist position:

- *LogitBoost* [Friedman et al., 2000] fits an additive logistic regression models by stagewise optimization of the Bernoulli log-likelihood. Instead of minimizing the

exponential loss as AdaBoost (which is an approximation of the Bernoulli log-likelihood), it minimizes the Bernoulli log-likelihood directly.

- *Gradient Boosting Machine (GBM)* developed by Friedman [2001] is a more general statistical framework for boosting that is based on its connection to **functional gradient descent (FGD)**, allowing to interpret boosting as a method for function estimation.
- *XGBoost* [Chen and Guestrin, 2016] is highly popular optimized distributed gradient boosting algorithm that builds on GBM framework. In mid 2010-s it was in many winning ML architectures in kaggle competitions.

Chapter 5

Lasso

Sparse regression methods have become increasingly important since data are measured on large number of variables (features). In such high-dimensional regression problems, the regression model is often *ill-posed*, i.e., the number of inputs p exceeds the number of outputs N

The popularity of sparse regression methods took off after the introduction of the lasso (Least Absolute Shrinkage and Selection Operator) regression [Tibshirani, 1996]. Its influence has been immerse in natural sciences, engineering and related fields. Lasso is the topic of this chapter.

We will denote the sample mean of responses and predictors variables by

$$\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i \quad \text{and} \quad \bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_p)^\top,$$

where $\bar{x}_j = \frac{1}{N} \sum_{i=1}^N x_{ij}$, for $j = 1, \dots, p$, is the sample mean of the j^{th} predictor variable.

In the computations of the lasso solution, it is convenient to use the centered responses/predictors,

$$\mathbf{y}_c = \mathbf{H}\mathbf{y} = \mathbf{y} - \bar{y}\mathbf{1} = \begin{pmatrix} y_1 - \bar{y} \\ \vdots \\ y_N - \bar{y} \end{pmatrix} \quad (5.1)$$

and

$$\begin{aligned} \mathbf{X}_c &= \mathbf{H}\mathbf{X} = \mathbf{X} - \mathbf{1}\bar{\mathbf{x}}^\top \\ &= \begin{pmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \cdots & x_{1p} - \bar{x}_p \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & \cdots & x_{2p} - \bar{x}_p \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} - \bar{x}_1 & x_{N2} - \bar{x}_2 & \cdots & x_{Np} - \bar{x}_p \end{pmatrix} \end{aligned} \quad (5.2)$$

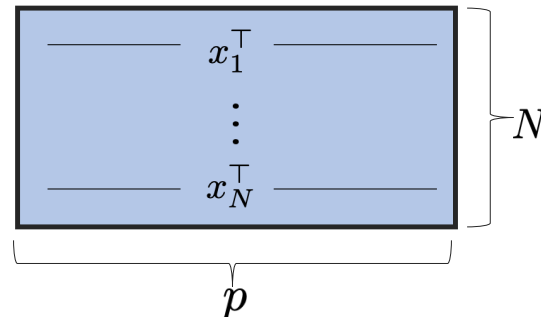
respectively, where

$$\mathbf{H} = \mathbf{I} - (1/N)\mathbf{1}\mathbf{1}^\top$$

is the *centering matrix* ($\mathbf{H}^2 = \mathbf{H}$, $\mathbf{H}^\top = \mathbf{H}$). We write the inner product in \mathbb{R}^N as $\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^N x_i y_i$.

5.1 Big Data Challenges

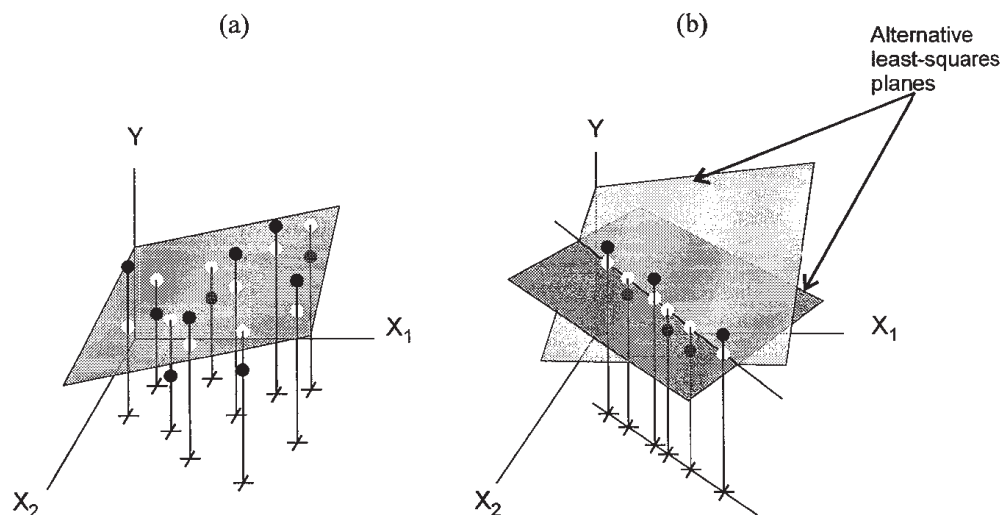
Often the predictor matrix \mathbf{X} is flat-and-long ($p > N$ or $N \approx p$):



Such data configuration causes many challenges:

- there are infinitely many least-squares (LS) solutions ($p > N$) or solution is subject to a large variance ($N \approx p$) \Rightarrow introducing some bias to the estimate is beneficial if the reduction in variance is considerable (bias-variance tradeoff).
- model complexity vs parsimony (interpretation): among the large # of predictors, we would like to identify the ones that exhibit the strongest effects \Rightarrow sparse $\hat{\beta}$ is desired.

Another issue is *multicollinearity*, i.e., very large correlations between predictors. This can cause huge variance, since the LSE $\hat{\beta}$ 'explodes' when $\text{cond}(\mathbf{X}^\top \mathbf{X})^{-1}$ is very large (see [Theorem 5.2](#)). If \mathbf{X} is not full rank ($\text{rank}(\mathbf{X}) < \min\{N, p\}$), the LSE is not unique and there are infinitely many solutions as is illustrated in the figure below.



5.2 Penalized/Regularized regression

How to solve the challenges above? A commonly used approach is to use regularization/penalization techniques, that is, *regularize* β_j 's, i.e., control how large the coefficients

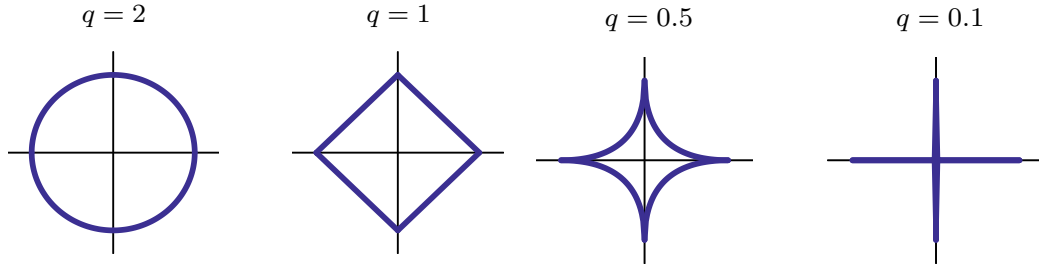


FIGURE 5.1: The constraint regions of unit $(\ell_q)^q$ balls: $\sum_{j=1}^p |\beta_j|^q \leq 1$, which is convex for $q \geq 1$ and non-convex for $0 \leq q < 1$.

cients are allowed to grow.

A general penalized regression problem is posed as

$$\min_{\beta_0, \beta} \{L(\beta_0, \beta) + \lambda P(\beta)\}, \quad (5.3)$$

where

- $L(\beta_0, \beta) : \mathbb{R} \times \mathbb{R}^p \rightarrow \mathbb{R}_0^+$ is the *criterion function* (data fidelity term) that depends on the data $\mathcal{T} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$.
- *Penalty function*: $P : \mathbb{R}^p \rightarrow \mathbb{R}_0^+$ penalizes large values of β , and can (when suitably chosen) enforce sparse solutions.
- *Penalty parameter* $\lambda > 0$ controls the trade-off between the two terms (data fidelity vs sparsity).

Note that penalization is not applied to intercept $\beta_0 \Rightarrow$ we do not assume that the regression plane should pass through or close to the origin.

One can pose the same problem in constrained (regularized) form:

$$\min_{\beta_0, \beta} L(\beta_0, \beta) \text{ subject to } P(\beta) \leq t,$$

where $L(\beta_0, \beta)$ and $P(\beta)$ are as earlier, and $t > 0$ is a regularization parameter that bounds the magnitude of the regression coefficients. For convex $L(\beta_0, \beta)$ and $P(\beta)$, the regularized and the penalized formulations are equivalent (1-to-1). A popular penalty function is the $(\ell_q)^q$ -norm $P(\beta) = \|\beta\|_q^q$ which is illustrated in [Figure 5.1](#). Recall that the ℓ_q -norm is defined as

$$\|\beta\|_q = \sqrt[q]{\sum_{j=1}^p |\beta_j|^q}.$$

5.3 Ridge regression

An older idea of regularization (prior to lasso) in the regression model is *ridge regression* (RR) [Hoerl and Kennard, 1970]. It uses residual sum of squares (RSS) as its criterion

function and the squared ℓ_2 -norm as its penalty function, and solves

$$\underset{\beta_0 \in \mathbb{R}, \beta \in \mathbb{R}^p}{\text{minimize}} \quad \frac{1}{2} \| \mathbf{y} - \beta_0 \mathbf{1} - \mathbf{X}\beta \|_2^2 + \lambda \|\beta\|_2^2, \quad (5.4)$$

where $\lambda > 0$ is the penalty parameter. The equivalent constrained form of (5.4) is

$$\underset{\beta_0 \in \mathbb{R}, \beta \in \mathbb{R}^p}{\text{minimize}} \quad \frac{1}{2} \| \mathbf{y} - \beta_0 \mathbf{1} - \mathbf{X}\beta \|_2^2 \quad \text{subject to} \quad \|\beta\|_2^2 \leq t,$$

where the shrinkage parameter t is 1-to-1 with λ .

Theorem 5.1. *The optimization problem (5.4) is strictly convex and has a unique minimizer for $\lambda > 0$ given by*

$$\begin{aligned} \hat{\beta}_{\text{RR},0}(\lambda) &= \bar{y} - \bar{\mathbf{x}}^\top \hat{\beta}_{\text{RR}}(\lambda), \\ \hat{\beta}_{\text{RR}}(\lambda) &= (\mathbf{X}_c^\top \mathbf{X}_c + \lambda \mathbf{I})^{-1} \mathbf{X}_c^\top \mathbf{y}_c, \end{aligned} \quad (5.5)$$

where \mathbf{y}_c and \mathbf{X}_c are centered response and feature matrix.

Proof. At lectures, if time permits. □

Discussion: What can we learn from (5.5)? Compare to LSE ($\lambda = 0$), or using different large/small values of λ . Also think of the case that features are highly correlated or that the feature matrix \mathbf{X}_c is orthonormal, $\mathbf{X}_c^\top \mathbf{X}_c = \mathbf{I}$.

Often the benefits of Ridge regression are most striking when predictors are correlated.

5.3.1 Computation of the ridge estimator

How to compute the ridge regression estimator? One easy way is by using augmented data set. Denote

$$\mathbf{X}_\lambda = \begin{pmatrix} \mathbf{X}_c \\ \sqrt{\lambda} \mathbf{I}_p \end{pmatrix} \quad \text{and} \quad \mathbf{y}_\lambda = \begin{pmatrix} \mathbf{y}_c \\ \mathbf{0}_{p \times 1} \end{pmatrix}$$

that is, we append p zeros to \mathbf{y}_c and a scaled $p \times p$ identity matrix to \mathbf{X}_c . Then observe that

$$\hat{\beta}_{\text{RR}}(\lambda) = (\mathbf{X}_\lambda^\top \mathbf{X}_\lambda)^{-1} \mathbf{X}_\lambda^\top \mathbf{y}_\lambda$$

and hence $\hat{\beta}_{\text{RR}}(\lambda)$ is nothing but a LS fit of \mathbf{X}_λ to \mathbf{y}_λ .

Remark 5.1. The RR estimator is not equivariant under scaling of variables. This also means that there is ambiguity in the specification of λ ; namely choice of $\lambda = 1$ would give a different estimator when variable X_1 is measured in the units of (say) kilograms instead of pounds. Due to this, the columns of \mathbf{X} are often standardized, i.e., the predictors are scaled so that they have a standard deviation equal to 1. The RR solution is then computed and retransformed back to the original scale. Same holds for lasso explained later. Standardization is especially important when predictors are on different scales and different types (categorical, continuous, etc).

5.3.2 Bias-variance tradeoff

Recall that mean squared error (MSE) is expressed as $\text{MSE} = \text{variance} + (\text{bias})^2$, and in multiparameter problems it reads as

$$\begin{aligned}\text{MSE}(\hat{\beta}) &= \mathbb{E}[(\hat{\beta} - \beta)(\hat{\beta} - \beta)^\top] \\ &= \text{cov}(\hat{\beta}) + \text{bias}(\hat{\beta}) \cdot [\text{bias}(\hat{\beta})]^\top\end{aligned}\quad (5.6)$$

where $\text{bias}(\hat{\beta}) = \mathbb{E}[\hat{\beta}] - \beta$ is the bias and

$$\text{cov}(\hat{\beta}) = \mathbb{E}[(\hat{\beta} - \mathbb{E}[\hat{\beta}])(\hat{\beta} - \mathbb{E}[\hat{\beta}])^\top]$$

denotes the covariance matrix of an estimator $\hat{\beta} \in \mathbb{R}^p$ of β .

Theorem 5.2. *If error terms ε_i 's in the linear regression model (1.4) are i.i.d. with zero mean and variance $\text{var}(\varepsilon_i) = \sigma^2$, then the bias and the covariance matrix of the ridge regression estimator is*

$$\text{bias}[\hat{\beta}_{\text{RR}}(\lambda)] = -\lambda \mathbf{R}_\lambda \beta \quad \text{and} \quad \text{cov}(\hat{\beta}_{\text{RR}}(\lambda)) = \sigma^2 \mathbf{R}_\lambda (\mathbf{X}_c^\top \mathbf{X}_c) \mathbf{R}_\lambda,$$

where $\mathbf{R}_\lambda = (\mathbf{X}_c^\top \mathbf{X}_c + \lambda \mathbf{I})^{-1}$.

Proof. At lectures, if time permits. □

It follows from Theorem 5.2 and (5.6) that the MSE is

$$\text{MSE}(\hat{\beta}_{\text{RR}}(\lambda)) = \mathbf{R}_\lambda \{\sigma^2 \mathbf{X}_c^\top \mathbf{X}_c + \lambda^2 \beta \beta^\top\} \mathbf{R}_\lambda.$$

Furthermore, in the model with no intercept, Theorem 5.2 also holds but then the centered \mathbf{X}_c is replaced with the non-centered predictor matrix \mathbf{X} . For $\lambda = 0$, the RR estimator equals the LSE, and then Theorem 5.2 gives the classic result:

$$\mathbb{E}[\hat{\beta}_{\text{LS}}] = 0 \quad \text{and} \quad \text{cov}(\hat{\beta}_{\text{LS}}) = \sigma^2 (\mathbf{X}_c^\top \mathbf{X}_c)^{-1},$$

where one needs to assume that $N > p$ and \mathbf{X}_c is full rank.

The total MSE is defined as

$$\text{tr}\{\text{MSE}(\hat{\beta})\} = \sum_{i=1}^p \text{MSE}(\hat{\beta}_i).$$

For ridge regression estimator, the total MSE becomes

$$\text{tr}\{\text{MSE}(\hat{\beta}_{\text{RR}}(\lambda))\} = \underbrace{\sum_{i=1}^p \text{var}(\hat{\beta}_{\text{RR},i}(\lambda))}_{\downarrow \text{ as } \lambda \uparrow} + \underbrace{\sum_{i=1}^p \text{bias}^2(\hat{\beta}_{\text{RR},i}(\lambda))}_{\uparrow \text{ as } \lambda \uparrow}.$$

Furthermore, there always exists λ such that the total MSE of $\hat{\beta}_{\text{RR}}(\lambda)$ is smaller than the total MSE of the LSE $\hat{\beta}_{\text{LS}}$. This is a surprising result: it states that even if the model is exactly correct and follows the exact distribution we specify, we can always obtain a better estimator by shrinking towards zero

Example 5.1. Consider the simple linear model with a single predictor ($p = 1$):

$$y_i = x_i\beta + \varepsilon_i, \quad i = 1, \dots, N,$$

Furthermore, write $\mathbf{x} = (x_1, \dots, x_N)^\top$ for the vector collecting all measurements on the feature. Assume feature is standardized such that $\mathbf{x}^\top \mathbf{x} = \sum_{i=1}^N x_i^2 = 1$ and the error terms are i.i.d. with mean $\mathbb{E}[\varepsilon_i] = 0$ and variance $\text{var}(\varepsilon_i) = \sigma^2$. Then the LSE is $\hat{\beta}_{\text{LS}} = \mathbf{x}^\top \mathbf{y}$ and the RR estimator is

$$\hat{\beta}_{\text{RR}}(\lambda) = (\mathbf{x}^\top \mathbf{x} + \lambda)^{-1} \mathbf{x}^\top \mathbf{y} = \frac{\hat{\beta}_{\text{LS}}}{1 + \lambda}.$$

Based on Theorem 5.2, the MSE of $\hat{\beta}_{\text{RR}}(\lambda)$ is

$$\text{MSE}(\hat{\beta}_{\text{RR}}(\lambda)) = \frac{\sigma^2 + \lambda^2 \beta^2}{(1 + \lambda)^2}.$$

The optimal penalty parameter λ^* that gives the minimum MSE is

$$\lambda^* = \arg \min_{\lambda} \text{MSE}(\hat{\beta}_{\text{RR}}(\lambda)) = \frac{\sigma^2}{\beta^2}. \quad (5.7)$$

The minimum MSE for the optimal RR estimator $\hat{\beta}^* = \hat{\beta}_{\text{RR}}(\lambda^*)$ is then

$$\text{MSE}(\hat{\beta}^*) = \frac{\sigma^2 + (\lambda^*)^2 \beta^2}{(1 + \lambda^*)^2} = \sigma^2 \frac{1}{1 + \frac{\sigma^2}{\beta^2}} \quad (5.8)$$

which is smaller than the MSE of LSE, $\text{MSE}(\hat{\beta}_{\text{LS}}) = \sigma^2$, for all value of $\beta \in \mathbb{R}$. Figure 5.2 illustrates the MSE curve as a function of λ in the case that $\beta = 1$ and $\sigma^2 = 0.1$. ■

5.4 Lasso

Lasso [Tibshirani, 1996] solves (5.3) using normalized RSS as its criterion function and ℓ_1 -norm as the penalty term:

$$(\hat{\beta}_0(\lambda), \hat{\beta}(\lambda)) = \arg \min_{\beta_0 \in \mathbb{R}, \beta \in \mathbb{R}^p} \frac{1}{2N} \|\mathbf{y} - \beta_0 \mathbf{1} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1 \quad (5.9)$$

where the penalty parameter $\lambda > 0$:

- controls the (bias-variance) tradeoff between the penalty and minimization of the sum of squared residuals (fit).
- the bigger the λ the greater is the amount of shrinkage. Some of the coefficients can be shrunk all the way to zero.

One may also express lasso in equivalent constrained form

$$\underset{\beta_0 \in \mathbb{R}, \beta \in \mathbb{R}^p}{\text{minimize}} \quad \frac{1}{2N} \|\mathbf{y} - \beta_0 \mathbf{1} - \mathbf{X}\beta\|_2^2 \quad \text{subject to} \quad \|\beta\|_1 \leq t,$$

where the shrinkage parameter t is 1-to-1 with λ .

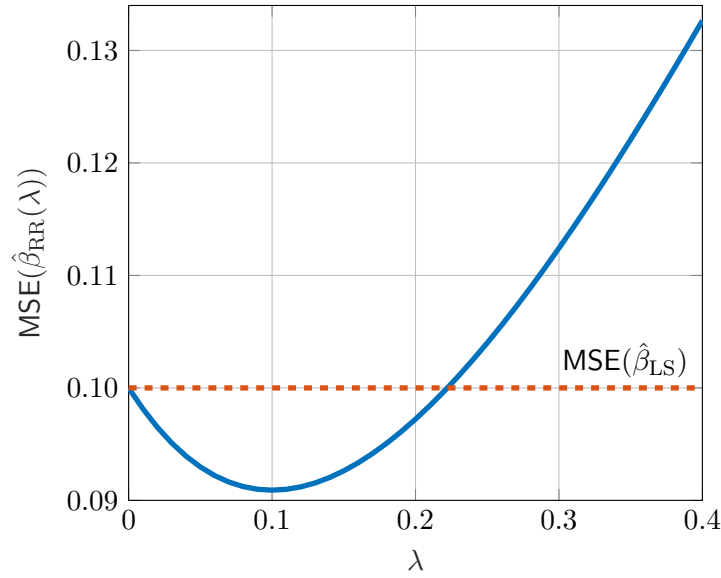


FIGURE 5.2: Plot of MSE of Ridge regression estimator $\hat{\beta}_{RR}(\lambda)$ as a function of λ of Example 5.1. Red dotted line displays the MSE of LSE.

5.4.1 Geometry of the lasso

The geometry of lasso and ridge regression is illustrated in Figure 5.3, where the optimization landscape is pictured in the regression model with no intercept ($\beta_0 = 0$) and $p = 2$ predictors. Note that the equicontours of the RSS criterion are ellipsoids that are centered (minimized) at the LSE. Now the lasso solution can be found as the point where the equicontours first touch the edge of the constraint region (the ℓ_1 -ball). The ridge and the lasso constraints (ℓ_1 and ℓ_2 -balls) are illustrated with the diamond and the circle centered at the origin. Explain from this figure, why lasso yields sparse estimates, but ridge regression does not.

5.4.2 Lasso solution path

In practice, the goal is to solve the complete lasso solution path, i.e., to solve $\hat{\beta}(\lambda)$ for a large range of λ values. To achieve this, most software compute the lasso solutions at a grid points of penalty values,

$$\begin{cases} [\lambda] = \{\lambda_0, \dots, \lambda_L\}, & \lambda_0 > \lambda_1 > \dots > \lambda_L, \\ & \lambda_0 = \max_j \frac{|\langle \mathbf{x}_j, \mathbf{y} \rangle|}{N} \end{cases} \quad (5.10)$$

where λ_0 above is the smallest penalty parameter λ for which the zero solution is obtained, i.e., $\hat{\beta}(\lambda_0) = \mathbf{0}$, but $\hat{\beta}(\lambda) \neq \mathbf{0}$ for $\lambda < \lambda_0$.

The sequence $\{\lambda_i\}$ is often chosen to be equispaced on log-scale:

$$\lambda_L = \epsilon \lambda_0 \quad \text{and} \quad \lambda_j = \epsilon^{j/L} \lambda_0 = \epsilon^{1/L} \lambda_{j-1}, \quad (5.11)$$

where $\epsilon = 10^{-4}$ and $L = 100$ are used as default values.

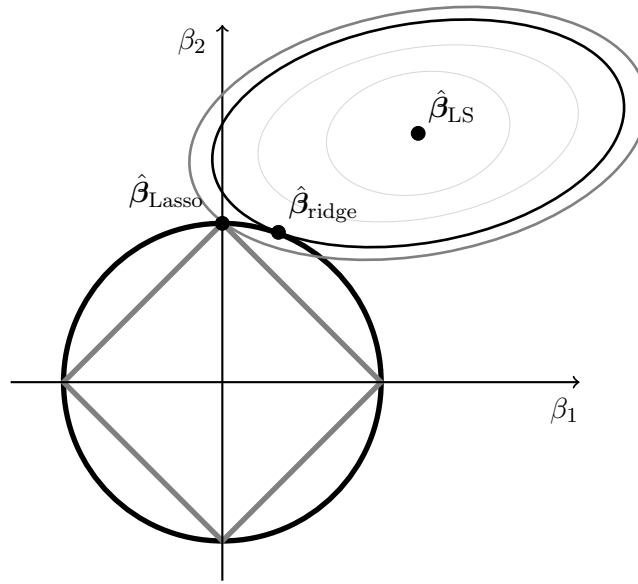


FIGURE 5.3: Penalized least squares solutions for the lasso and ridge regression.

Example 5.2. We consider the benchmark prostate cancer data set ($N = 97, p = 8$) used in many text-books, e.g., Hastie et al. [2009]. The goal is to predict the logarithm of the prostate specific antigen (psa) measurement, $\log(\text{lpsa})$, for men who are about to receive a radical prostatectomy, as a function of the number of clinical measures used. The study had a total of $N = 97$ observations of male patients aged from 41 to 79 years. The predictor variables are the logarithm of the cancer volume (lcavol), the logarithm of the prostate weight (lweight), the age of the patient (age), the logarithm of the benign prostatic hyperplasia amount (lbph), the presence or absence of seminal vesicle invasion (svi), the logarithm of the capsular penetration (lbph), the Gleason grade (gleason), and the percent Gleason grade 4 or 5 (pgg45).

We fit a linear model to the lpsa after first standardizing the predictors to have zero mean and unit variance. Furthermore, the data set is split into a training set of size 67 and a test set of size 30. Figure 5.4 depicts the lasso solution path. ■

5.4.3 Standardizing the features and the penalty

Consider expressing the feature j on different scale, so a transformation of a sort:

$$x_{ij}^* = \frac{x_{ij}}{a_j}, \quad i = 1, \dots, N. \quad (5.12)$$

Then note that

$$x_{ij}\beta_j = \frac{x_{ij}}{a_j} \cdot a_j\beta_j = x_{ij}^*\beta_j^*,$$

where $\beta_j = \frac{\beta_j^*}{a_j}$. In general, if we divide each feature vector $\mathbf{x}_j \in \mathbb{R}^N$ by a_j , e.g., to standardize the feature or to express it in different scale (such as when expressing cm in mm), we divide the transformed estimate $\hat{\beta}_j^*$ by a_j to obtain original regression estimate $\hat{\beta}_j$ on the original scale. It is important when running the lasso *to standardize*

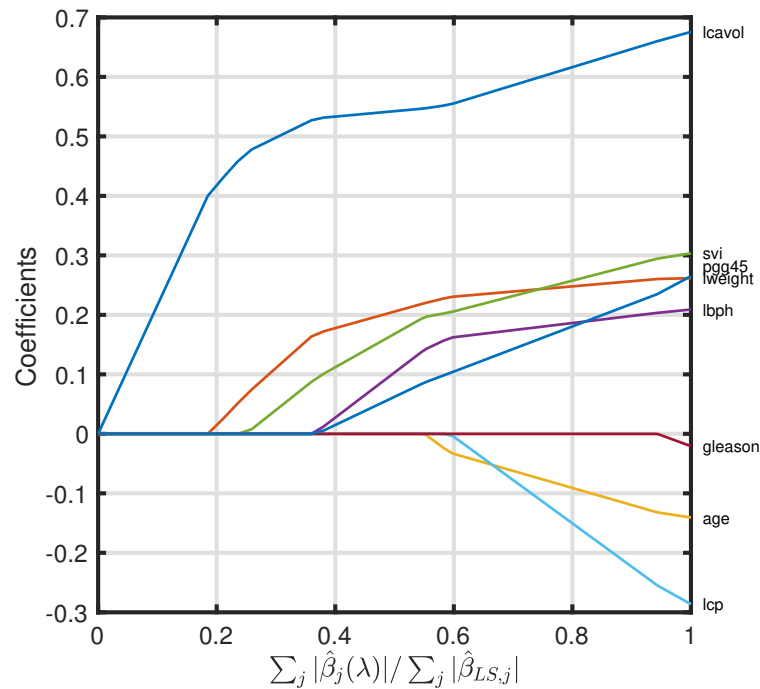


FIGURE 5.4: Lasso solution path where coefficients are plotted versus the (normalized) shrinkage parameter $s = t / \|\hat{\beta}_{LS}\|_1$, where $t = \sum_j |\hat{\beta}_j(\lambda)|$.

the features on same scale as otherwise the lasso solutions would depend on the units. Thus standardization is a basic preprocessing step in CCD.

Different software can use different standardization. For example, scikit-learn uses $\|\mathbf{x}_j\|^2 = 1$, while I (as well as you in HW3) like to standardize such that $\|\mathbf{x}_j\| = N$. This appears more natural as it is equivalent to stating that sample standard deviation of the feature variable is equal to 1. The former standardization applied in sklearn can also make elements too small when N is large.

5.5 Computation of the lasso solution

Difficulties arise when solving the problem of the form:

$$\underset{\beta_0 \in \mathbb{R}, \beta \in \mathbb{R}^p}{\text{minimize}} \quad L(\beta_0, \beta) + \lambda \sum_{j=1}^p P(\beta_j).$$

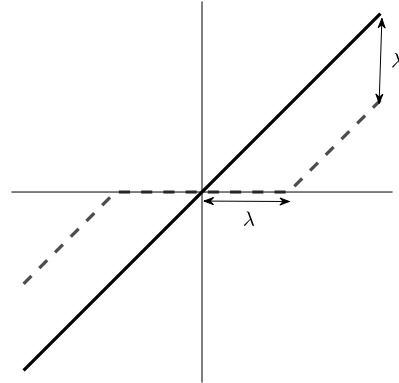
are

- ✗ *Non-smoothness:* objective function is not differentiable at $\beta_j = 0$ e.g., when using lasso penalty $P(\beta_j) = |\beta_j|$.
- ✗ *non-convexity:* e.g., if P is non-convex (we discuss this later).
- ✗ *High-dimensionality:* p can large or huge... ($p > 10^6$)

Cyclic coordinate descent algorithm offers a scalable method (when implemented carefully) to compute the lasso solution path.

Soft-thresholding operator is a basic building block for computing sparse regression estimates.

$$\begin{aligned}\mathcal{S}_\lambda(x) &= \text{sign}(x)(|x| - \lambda)_+ \\ &= \begin{cases} x - \lambda & \text{if } x > \lambda \\ 0 & \text{if } |x| \leq \lambda \\ x + \lambda & \text{if } x < -\lambda, \end{cases}\end{aligned}$$



Here $(t)_+$ denotes the positive part of $t \in \mathbb{R}$: $= t$ if $t > 0$ and 0 otherwise, and $\text{sign}(x)$ is the sign function, i.e., $\text{sign}(x) = +1, -1, 0$ if $x > 0, < 0, = 0$. ST-operator arises from these two results:

Theorem 5.3. (a) Given $y \in \mathbb{R}$, one has that

$$\begin{aligned}\hat{\beta}(\lambda) &= \arg \min_{\beta \in \mathbb{R}} \frac{1}{2}(y - \beta)^2 + \lambda|\beta| \\ &= \mathcal{S}_\lambda(y)\end{aligned}$$

(b) In the single predictor ($p = 1$) case, lasso has closed-form solution:

$$\begin{aligned}\hat{\beta}(\lambda) &= \arg \min_{\beta \in \mathbb{R}} \frac{1}{2N} \sum_{i=1}^N (y_i - \beta x_i)^2 + \lambda|\beta| \\ &= \mathcal{S}_\lambda\left(\frac{\langle \mathbf{x}, \mathbf{y} \rangle}{N}\right)\end{aligned}$$

where the predictor $\mathbf{x} = (x_1, \dots, x_N)^\top$ is standardized such that $\|\mathbf{x}\|^2 = \mathbf{x}^\top \mathbf{x} = N$.

Proof. At lectures, if time permits. □

ST-operator is used in CCD algorithm for computing the lasso solution. Now recall that a proximal operator (or a proximal map) of convex function h can be defined as

$$\text{prox}_h(\mathbf{z}) = \min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{z} - \beta\|_2^2 + h(\beta).$$

Note: It follows from the definition that

$$\text{prox}_{\alpha h}(\mathbf{z}) = \min_{\beta \in \mathbb{R}^p} \frac{1}{2\alpha} \|\mathbf{z} - \beta\|_2^2 + h(\beta).$$

Thus it follows from [Theorem 6.1a](#), and from separability of lasso penalty that proximal operator of $\lambda\|\beta\|_1$ is

$$\text{prox}_{\lambda\|\cdot\|_1}(\mathbf{z}) = \min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{z} - \beta\|_2^2 + \lambda\|\beta\|_1 = \mathcal{S}_\lambda(\mathbf{z}) \quad (5.13)$$

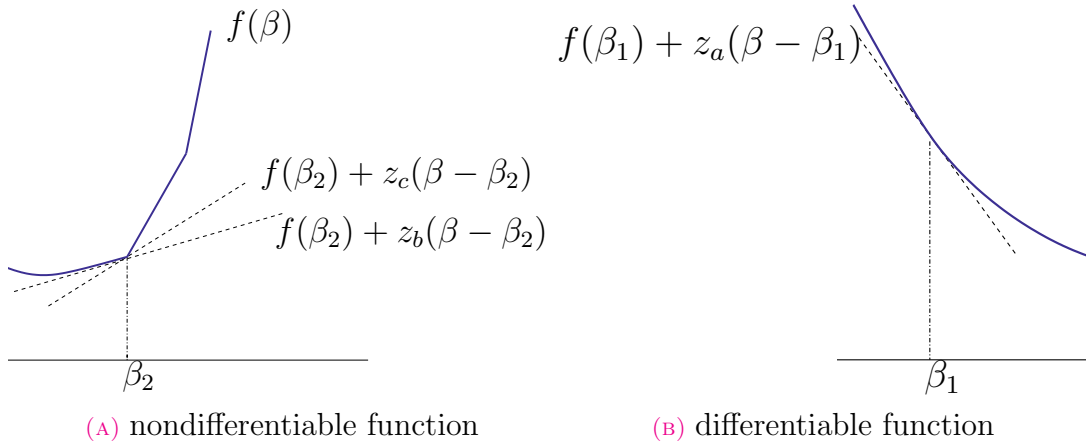


FIGURE 5.5: Subdifferential of a function, [Hastie et al., 2015, Figure 5.3].

5.5.1 Subgradient optimality conditions

Recal from Prof. Vorobyov's lecture notes that for a convex and differentiable function $f : \mathbb{R}^p \rightarrow \mathbb{R}$ it holds that

$$f(\beta') \geq f(\beta) + \nabla f(\beta)^\top (\beta' - \beta) \quad \forall \beta' \in \mathbb{R}^p.$$

A vector $\mathbf{z} \in \mathbb{R}^p$ is called a subgradient of f at β if

$$f(\beta') \geq f(\beta) + \mathbf{z}^\top (\beta' - \beta) \quad \forall \beta' \in \mathbb{R}^p,$$

i.e., linear approximation always underestimates f . The set of all subgradients of f at β is called the *subdifferential*, denoted $\partial f(\beta)$. For convex f , such \mathbf{z} always exists. If f is differentiable at β , then $\mathbf{z} = \nabla f(\beta)$ uniquely. See Figure 5.5.

Then recall (cf. Section 3.2 of Prof. Vorobyov's lecture notes) that for any convex subdifferentiable function $D : \mathbb{R}^p \rightarrow \mathbb{R}$, it holds that

$$\hat{\beta} = \arg \min_{\beta} D(\beta) \quad \Leftrightarrow \quad \mathbf{0} \in \partial D(\hat{\beta}). \quad (5.14)$$

Subdifferential of $|\beta_j|$ is

$$\partial |\beta_j| = \begin{cases} \text{sign}(\beta_j), & \text{for } \beta_j \neq 0 \\ [-1, +1] & \text{for } \beta_j = 0 \end{cases}$$

Consider lasso problem (without intercept):

$$\underset{\beta}{\text{minimize}} \left\{ D(\beta) = \frac{1}{2N} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1 \right\}$$

A necessary and sufficient condition for $\hat{\beta}$ to be the minimizer is then that it solves the *zero subgradient equations*:

$$\partial \left(\frac{1}{2N} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1 \right) \ni \mathbf{0}.$$

Thus $\hat{\beta}$ is a lasso solution iff

$$-\frac{1}{N}\mathbf{x}_j^\top(\mathbf{y} - \mathbf{X}\hat{\beta}) + \lambda s_j = 0, \quad j = 1, \dots, p,$$

where s_j is an element of subdifferential of $|\beta_j|$ evaluated at $\hat{\beta}_j$, so a number verifying $|s_j| \leq 1$. This means that the solution verifies

$$\frac{1}{N}\mathbf{x}_j^\top(\mathbf{y} - \mathbf{X}\hat{\beta}) = \begin{cases} \lambda \text{sign}(\hat{\beta}_j), & \text{if } \hat{\beta}_j \neq 0 \\ \lambda s_j, & \text{if } \hat{\beta}_j = 0 \end{cases}. \quad (5.15)$$

Thus, once you have an algorithm that has found solution, you can always make a sanity check that the condition (5.15) hold. The condition (5.14) also provides a means to derive the proximal operators for many penalties.

5.5.2 Cyclic coordinate descent

Assume now that the penalized objective function,

$$D(\beta_0, \beta_1, \dots, \beta_p) = L(\beta_0, \beta_1, \dots, \beta_p) + \lambda \sum_{j=1}^p P(\beta_j)$$

is such that $L(\beta_0, \beta)$ is convex and differentiable and $P(\cdot)$ is convex (but not necessarily differentiable). *Cyclic Coordinate descent (CCD)* updates β_j by minimizing D in this coordinate while keeping others fixed:

$$\beta_j \leftarrow \arg \min_{\beta_j \in \mathbb{R}} D(\hat{\beta}_0, \dots, \hat{\beta}_{j-1}, \beta_j, \hat{\beta}_{j+1}, \dots, \hat{\beta}_p),$$

and repeatedly cycles through the coefficients one at a time ($j = 0, 1, \dots, p$) until convergence. Tseng [2001] showed that any limit point of CCD is a minimizer of D .

The benefits of CCD are:

- ✓ CCD is a simple algorithm and very easy to implement.
- ✓ Useful and general method for cases, when the single parameter (i.e., one coordinate at a time) problem is easy to solve.
- ✓ Can be used to compute the whole lasso path.

The lasso objective function is separable in coordinates:

$$D(\beta_0, \beta) = \frac{1}{2N} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{k \neq j} x_{ik} \beta_k - x_{ij} \beta_j \right)^2 + \lambda |\beta_j| + \lambda \sum_{k \neq j} |\beta_k|. \quad (5.16)$$

Update of the intercept is (when holding β_j fixed at their current estimates $\hat{\beta}_j$)

$$\begin{aligned}
 \hat{\beta}_0 &\leftarrow \frac{1}{N} \sum_{i=1}^N (y_i - \sum_j x_{ij} \hat{\beta}_j) \\
 &= \frac{1}{N} \sum_{i=1}^N (y_i - \hat{\beta}_0 - \sum_j x_{ij} \hat{\beta}_j + \hat{\beta}_0) \\
 &= \hat{\beta}_0 + \frac{1}{N} \sum_{i=1}^N \hat{r}_i,
 \end{aligned} \tag{5.17}$$

where

$$\hat{r}_i = y_i - \hat{\beta}_0 - \sum_{j=1}^p x_{ij} \hat{\beta}_j,$$

are the current full residuals before the update.

As a preprocessing step, we need to standardize the predictors (the columns \mathbf{x}_j of matrix \mathbf{X} such that $\mathbf{x}_j^\top \mathbf{x}_j = N$) holds. After running the algorithm, we then need to rescale the found regression estimates back to the original scale. This is discussed later.

When optimizing for β_j (holding β_k fixed at their current estimates $\hat{\beta}_k$, $k \neq j$, $j \geq 1$) the last term in (5.16) is a constant and we solve

$$\begin{aligned}
 \hat{\beta}_j &\leftarrow \arg \min_{\beta_j} \frac{1}{2N} \sum_{i=1}^N (y_i - \hat{\beta}_0 - \sum_{k \neq j} x_{ik} \hat{\beta}_k - x_{ij} \beta_j)^2 + \lambda |\beta_j| \\
 &= \arg \min_{\beta_j} \frac{1}{2N} \sum_{i=1}^N (y_i - \hat{\beta}_0 - \sum_k x_{ik} \hat{\beta}_k + x_{ij} \hat{\beta}_j - x_{ij} \beta_j)^2 + \lambda |\beta_j| \\
 &= \arg \min_{\beta_j} \frac{1}{2N} \sum_{i=1}^N (\hat{r}_i + x_{ij} \hat{\beta}_j - x_{ij} \beta_j)^2 + \lambda |\beta_j| \\
 &= \mathcal{S}_\lambda \left(\hat{\beta}_j + \frac{1}{N} \langle \mathbf{x}_j, \hat{\mathbf{r}} \rangle \right),
 \end{aligned} \tag{5.18}$$

where the last identity follows by applying [Theorem 6.1b](#) and $\hat{\mathbf{r}} = (\hat{r}_1, \dots, \hat{r}_N)^\top$ is a vector of current residuals (before the update).

The update for β_j in (5.18) is of the form

$$\text{"new estimate"} \leftarrow \mathcal{S}_\lambda(\text{"current estimate"} + \text{"correction"})$$

and it is in fact just soft-thresholding the conventional coordinate descent update. CCD algorithm is applying soft-thresholding update repeatedly in a cyclical manner, updating the coordinates of $\hat{\beta}$ and the residual vector $\hat{\mathbf{r}}$ vector as it proceeds. Each coordinate update in (5.18) requires computing $\langle \mathbf{x}_j, \hat{\mathbf{r}} \rangle$ and update $\hat{\mathbf{r}} \leftarrow \hat{\mathbf{r}} + (\hat{\beta}_j^{\text{old}} - \hat{\beta}_j) \mathbf{x}_j$ which is of $O(N)$ flops.

Remark 5.2. Cyclic updates of the intercept (5.17) in the CCD algorithm can be omitted if one simply runs the CCD algorithm (assuming a model with no intercept term, $\beta_0 = 0$) but for centered data $\mathbf{X}_c, \mathbf{y}_c$. This is because the optimal solution $\hat{\beta}(\lambda)$

for the centered data $\mathbf{X}_c, \mathbf{y}_c$ is the same as for uncentered data \mathbf{X}, \mathbf{y} (for a model with intercept). Thus the intercept can be calculated in the last stage, simply as $\hat{\beta}_0(\lambda) = \bar{y} - \bar{\mathbf{x}}^\top \hat{\boldsymbol{\beta}}(\lambda)$. This is done in line 5 of [Algorithm 5.1](#). Note also that this procedure is analogous to result of [Theorem 5.1](#) for ridge regression.

The basic form of CCD algorithm for lasso (named `ccdlasso`) tabulated in [Algorithm 5.2](#) assumes that there is no intercept in the model and presumes that the features are standardized so that $\|\mathbf{x}_j\|^2 = N$ holds. The general lasso algorithm implemented e.g., in scikit-learn is given in [Algorithm 5.1](#). When computing the solution it first centers the inputs and outputs, and standardizes the predictors. Then it uses the CCD algorithm to compute the solution (line 3) for regression coefficients. Then it transforms the found solution to the original scale of the data, and computes the intercept (if it is in the model). Naturally, Line 1 (centering) and line 5 (computing the intercept) can be omitted in [Algorithm 5.1](#) if the intercept is not in the linear regression model in the first place.

Algorithm 5.1: lasso algorithm that computes the lasso solution using CCD algorithm in a model with intercept

Input : response $\mathbf{y} \in \mathbb{R}^N$, predictor matrix $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_p) \in \mathbb{R}^{N \times p}$,
penalty parameter $\lambda > 0$, initial estimate (warm start) $\hat{\boldsymbol{\beta}}_{\text{init}} \in \mathbb{R}^p$.

1 Center the inputs and outputs:

$$\mathbf{x}_j \leftarrow \mathbf{x}_j - \bar{x}_j \mathbf{1} \quad \text{and} \quad \mathbf{y} \leftarrow \mathbf{y} - \bar{y} \mathbf{1}$$

$$(j = 1, \dots, p), \text{ where } \bar{x}_j = \frac{1}{N} \sum_{i=1}^N x_{ij} \text{ and } \bar{y} = \frac{1}{N} \sum_{i=1}^N y_i.$$

2 Standardize the feature vectors:

$$\mathbf{x}_j \leftarrow \mathbf{x}_j / s_j \quad \text{for } j = 1, \dots, p$$

$$\text{where } s_j = \|\mathbf{x}_j\|_2 / \sqrt{N}.$$

3 $\hat{\boldsymbol{\beta}}(\lambda) \leftarrow \text{ccdlasso}(\mathbf{y}, \mathbf{X}, \lambda, \hat{\boldsymbol{\beta}}_{\text{init}})$

4 Transform the regression coefficient back to the original scale:

$$\hat{\beta}_j(\lambda) \leftarrow \hat{\beta}_j(\lambda) / s_j \quad \text{for } j = 1, \dots, p$$

5 $\hat{\beta}_0(\lambda) = \bar{y} - \bar{\mathbf{x}}^\top \hat{\boldsymbol{\beta}}(\lambda)$ // Compute the intercept

Output : $(\hat{\beta}_0(\lambda), \hat{\boldsymbol{\beta}}(\lambda)) \in \mathbb{R} \times \mathbb{R}^p$, solution of (5.9)

5.5.3 Lasso solution path

The lasso solution path is computed over a grid $[\lambda]$ in (5.10) of penalty parameter values. The algorithm starts from λ_0 that yields all zeros solution and then goes to next (smaller) value on the grid and uses the previous estimate as a warm start. The pseudo-code for the algorithm, called *pathwise coordinate descent* [Friedman et al., 2007] algorithm, is detailed in [Algorithm 5.3](#).

Algorithm 5.2: `ccdlasso` computes lasso solution for standardized predictors in a model with no intercept.

Input : response $\mathbf{y} \in \mathbb{R}^N$, predictor matrix $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_p) \in \mathbb{R}^{N \times p}$, warm start $\hat{\boldsymbol{\beta}}_{\text{init}} \in \mathbb{R}^p$, penalty parameter $\lambda > 0$. Predictors are standardized such that $\mathbf{x}_j^\top \mathbf{x}_j = N$ holds for $j = 1, \dots, p$.

Initialize: Maximum number of iterations, e.g., $I_{\max} = 10^4$; Convergence threshold, e.g., $\delta = 10^{-4}$

- 1 Set $\hat{\mathbf{r}} \leftarrow \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{init}}$, $\hat{\boldsymbol{\beta}}^{\text{old}} \leftarrow \hat{\boldsymbol{\beta}}_{\text{init}}$
- 2 **for** $i = 1, \dots, I_{\max}$ **do**
- 3 **for** $j = 1$ **to** p **do**
- 4 $\hat{\beta}_j \leftarrow \mathcal{S}_\lambda(\hat{\beta}_j + \frac{1}{N}\langle \mathbf{x}_j, \hat{\mathbf{r}} \rangle)$
- 5 $\hat{\mathbf{r}} \leftarrow \hat{\mathbf{r}} + (\hat{\beta}_j^{\text{old}} - \hat{\beta}_j)\mathbf{x}_j$
- 6 **if** $\|\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^{\text{old}}\|_2 / \|\hat{\boldsymbol{\beta}}\|_2 < \delta$ **then**
- 7 **break**
- 8 $\hat{\boldsymbol{\beta}}^{\text{old}} \leftarrow \hat{\boldsymbol{\beta}}$

Output : $\hat{\boldsymbol{\beta}}(\lambda) = \hat{\boldsymbol{\beta}}$, the minimizer of $\frac{1}{2N}\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_1$

Algorithm 5.3: `lassopath`: pathwise coordinate descent lasso algorithm.

Input : $\mathbf{y} \in \mathbb{R}^N$, $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_p) \in \mathbb{R}^{N \times p}$, L (grid size)

Initialize: $\epsilon = 10^{-3}$

- 1 Compute line 1 and line 2 of [Algorithm 5.1](#) (that is, center the predictors and responses and standardize the predictors)
- 2 Generate the grid: $\lambda_0 = \max_j \frac{1}{N}|\langle \mathbf{x}_j, \mathbf{y} \rangle|$, $\lambda_j = \epsilon^{j/L}\lambda_{j-1}$, $j = 1, \dots, L$.
- 3 $\hat{\boldsymbol{\beta}}_{\text{init}} \leftarrow \hat{\boldsymbol{\beta}}(\lambda_0) \leftarrow \mathbf{0}$
- 4 **for** $j = 1, \dots, L$ **do**
- 5 $\lambda \leftarrow \lambda_j$
- 6 $\hat{\boldsymbol{\beta}}(\lambda) \leftarrow \text{ccdlasso}(\mathbf{y}, \mathbf{X}, \hat{\boldsymbol{\beta}}_{\text{init}} \lambda)$
- 7 $\hat{\boldsymbol{\beta}}_{\text{init}} \leftarrow \hat{\boldsymbol{\beta}}(\lambda)$
- 8 Compute line 4 of [Algorithm 5.1](#) for each $\lambda \in [\lambda]$ to transform the obtained regression estimates back to original scale and compute the intercepts $\hat{\beta}_0(\lambda)$ as in line 5 of [Algorithm 5.1](#)

Output : $(\hat{\beta}_0(\lambda), \hat{\boldsymbol{\beta}}(\lambda)) \in \mathbb{R} \times \mathbb{R}^p$ for $\lambda \in [\lambda]$, solutions of (5.9)

5.5.4 Why CCD works for large-scale data?

CCD can be implemented with smart tricks (as in `GLMnet`):

✓ For large λ , most coordinates that are zero never become non-zero.

⇒ **active set** strategy updates active predictors (i.e., nonzero coefficients) until convergence and then check other variables. See Tibshirani and et al. [2012] for details.

- ✓ **warm starts:** move from large λ to smaller, using solutions at previous λ as initial value for next λ .
- ✓ CCD is easy to extend to generalized linear models (GLM)
- ✗ Coding in lower-level language (C++/Fortran) is necessary due to iterative nature of CCD.

GLMnet which uses Fortran and tricks above in its CCD implementation is very fast.

5.6 Discussion

Benefits of the lasso regression are:

- ✓ Penalty (smart choice of λ) offers an *automated variable selection* (lasso performs estimation and variable selection simultaneously).
- ✓ lasso does variable selection and shrinkage; *ridge only shrinks*.
- ✓ Depicting the whole lasso solution path (as λ grow) informs us when variables drop-out from the model.
- ✓ Works for underdetermined systems ($p > n$) which occur commonly in many applications.
- ✓ Growing importance of sparse representations and modelling.

Shortcoming of the lasso regression are:

- ☹ Lasso solution is unique when the columns of \mathbf{X} are in *general position** and $\lambda > 0$. This holds true even when $N \leq p$
- ☹ When $N \leq p$, the number of nonzero coefficients in any lasso solution is at most N .
- ☹ Lasso *does not cope well with multicollinearity*; the coefficient paths tend to be erratic and can show wild behavior. Moreover, if \mathbf{X} is not full column rank, solution is not unique, and there can be infinitely many solutions.
- ☹ Lasso *does not select groups* \Rightarrow one may wish to select the whole group if one variable amongst them is selected
- ☹ Lasso ignores also possible structured sparsity (e.g., block sparsity, smoothness, etc).

*Columns $\{\mathbf{x}_j\}_{j=1}^p$ of \mathbf{X} are in general position if any affine subspace $\mathbb{L} \subset \mathbb{R}^n$ of dimension $k < N$ contains at most $k + 1$ elements of the set $\{\pm \mathbf{x}_1, \pm \mathbf{x}_2, \dots, \pm \mathbf{x}_p\}$.

Chapter 6

Lasso: extensions

6.1 Elastic net

One of the shortcomings of the lasso was that it *does not select groups*. Ridge regression on the other hand does not have this problem and can handle well correlated predictors.

In some problems, as in microrarray studies, one has that $p \gg N$ and groups of genes in the same biological pathway tend to be expressed (or not) together, and hence measures of their expression tend to be strongly correlated. One can think of these genes as forming a group. Elastic net [Zou and Hastie, 2005] is an extension of lasso that was designed for such problems. It aims at making a compromise between the ridge and the lasso penalties.

The EN penalty is a convex combination of ridge and lasso penalties, defined as

$$P_{\text{EN}}(\beta; \alpha) = \alpha \|\beta\|_1 + \frac{1}{2}(1 - \alpha) \|\beta\|_2^2 = \sum_{i=1}^p \left[\alpha |\beta_i| + \frac{1}{2}(1 - \alpha) \beta_i^2 \right],$$

where $\alpha \in [0, 1]$ is a *tuning parameter* that can be varied:

- $\alpha = 1$ corresponds to lasso regression
- $\alpha = 0$ corresponds to ridge regression

Note that α can be set on subjective grounds or using cross-validation scheme on a (coarse) grid of α values. Elastic net solves the following convex optimization problem

$$\underset{(\beta_0, \beta) \in \mathbb{R} \times \mathbb{R}^p}{\text{minimize}} \left\{ \frac{1}{2N} \sum_{i=1}^N (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2 + \lambda \left[\alpha \|\beta\|_1 + \frac{1}{2}(1 - \alpha) \|\beta\|_2^2 \right] \right\} \quad (6.1)$$

where $\lambda \geq 0$ is the penalty parameter. In constrained form the problem (6.1) can be expressed as

$$\underset{(\beta_0, \beta) \in \mathbb{R} \times \mathbb{R}^p}{\text{minimize}} \quad \frac{1}{2N} \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2 \quad \text{s.t.} \quad \frac{1}{2}(1 - \alpha) \|\beta\|_2^2 + \alpha \|\beta\|_1 \leq t$$

where $t > 0$ denotes constraint/regularization parameter (1-to-1 with λ). The penalty $P_{\text{EN}}(\beta; \alpha)$ is non-differentiable $\forall \alpha \in (0, 1)$ at any point where at least one coordinate

β_j is equal to zero, but it is strictly convex. Thus a unique minimizer for problem (6.1) exists $\forall \lambda > 0$ also when $p > N$ or in case of multicollinearity.

The benefits of EN are

- removes the limitation on the number of selected variables
- encourages grouping effect
- stabilizes the coefficient paths

The EN penalty in $p = 3$ is illustrated in Figure 6.1. Use the figure to explain: (a) what causes the grouping effect? (b) Does EN still give sparse solution?

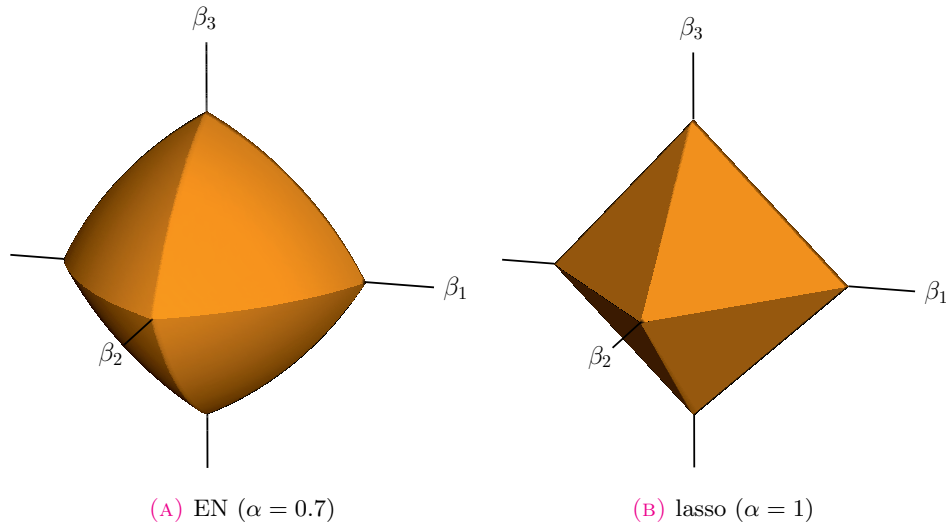


FIGURE 6.1: Illustration of EN ja lasso penalties in $p = 3$. From [Hastie et al., 2015, Figure 4.2].

6.1.1 Computation via the CCD

Cyclic coordinate descent (CCD) algorithm can be used as in lasso using the following result:

Theorem 6.1. (a) Given $y \in \mathbb{R}$, one has that

$$\begin{aligned}\hat{\beta}(\lambda, \alpha) &= \arg \min_{\beta \in \mathbb{R}} \frac{1}{2}(y - \beta)^2 + \lambda\alpha|\beta| + \frac{\lambda(1 - \alpha)}{2}\beta^2 \\ &= \frac{\mathcal{S}_{\lambda\alpha}(y)}{1 + \lambda(1 - \alpha)}\end{aligned}$$

(b) In the single predictor ($p = 1$) case, EN criterion has closed-form solution:

$$\begin{aligned}\hat{\beta}(\lambda, \alpha) &= \arg \min_{\beta \in \mathbb{R}} \frac{1}{2N} \sum_{i=1}^N (y_i - \beta x_i)^2 + \lambda\alpha|\beta| + \frac{\lambda(1 - \alpha)}{2}\beta^2 \\ &= \frac{\mathcal{S}_{\lambda\alpha}\left(\frac{\langle \mathbf{x}, \mathbf{y} \rangle}{N}\right)}{1 + \lambda(1 - \alpha)}\end{aligned}\tag{6.2}$$

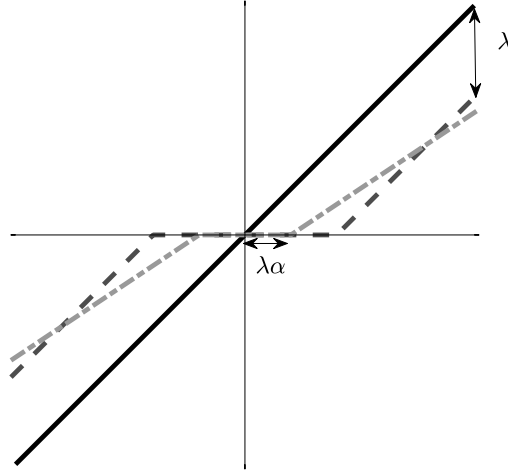


FIGURE 6.2: The EN solution $\hat{\beta}(\lambda, \alpha)$. The black solid line shows the case when $\lambda = 0$ (no penalization). The dashed line depicts the special case of lasso ($\alpha = 1$), while the dash-dotted line corresponds to the case of $\alpha = 0.5$

where the predictor $\mathbf{x} = (x_1, \dots, x_N)^\top$ is standardized such that $\|\mathbf{x}\|^2 = \mathbf{x}^\top \mathbf{x} = N$).

Proof. At lectures, if time permits. □

Note that when $\lambda = 0$, then (6.2) reduces to LSE $\hat{\beta}_{\text{LS}} = \mathbf{x}^\top \mathbf{y} / N$. Equation (6.2) sparks several remarks. First, note that increasing the penalty λ will shrink the EN estimator $\hat{\beta}(\lambda, \alpha)$ linearly towards zero and when $\lambda\alpha$ exceeds the magnitude of the LSE $|\hat{\beta}_{\text{LS}}|$, the EN estimator becomes zero. Second, the sign of the EN estimator is the same as the sign of the LSE, only the magnitude is shrunk. Figure 6.2 illustrates the effect of extra shrinkage of EN soft-thresholding operator for the case of $\alpha = 0.5$, versus the lasso soft threshold ($\alpha = 1$).

The CCD update for j^{th} coefficient is

$$\hat{\beta}_j \leftarrow \frac{\mathcal{S}_{\alpha\lambda}(\hat{\beta}_j + \frac{1}{N} \langle \mathbf{x}_j, \hat{\mathbf{r}} \rangle)}{1 + \lambda(1 - \alpha)} \quad (6.3)$$

where $\hat{\mathbf{r}}$ represents the current residual and $\mathcal{S}_\lambda(x) = \text{sign}(x)(|x| - \lambda)_+$ denotes the soft-thresholding operator as earlier. Note that (6.3) reduces to familiar lasso CCD update when $\alpha = 1$. As in lasso, the predictors are standardized ($\mathbf{x}_j^\top \mathbf{x}_j = N$). Hence the only thing that changes in `cdlasso` algorithm is step 4 which is replaced by the update (6.3). The subgradient optimality condition is now

$$\frac{1}{N} \mathbf{x}_j^\top (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}) - \lambda(1 - \alpha) \hat{\beta}_j = \begin{cases} \lambda\alpha \text{sign}(\hat{\beta}_j), & \text{if } \hat{\beta}_j \neq 0 \\ \lambda\alpha s_j, & \text{if } \hat{\beta}_j = 0 \end{cases}, \quad (6.4)$$

where s_j is a number verifying $|s_j| \leq 1$.

Algorithm 6.1: `ccden` computes elastic net solution for standardized predictors in a model with no intercept.

Input : Data $\mathbf{y} \in \mathbb{R}^N$, $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_p) \in \mathbb{R}^{N \times p}$, warm start $\hat{\boldsymbol{\beta}}_{\text{init}} \in \mathbb{R}^p$, parameters $\lambda > 0$, $\alpha \in [0, 1]$. Predictors are standardized such that $\mathbf{x}_j^\top \mathbf{x}_j = N$ holds for $j = 1, \dots, p$.

Initialize: Maximum number of iterations, e.g., $I_{\max} = 10^4$; Convergence threshold, e.g., $\delta = 10^{-4}$

- 1 Set $\hat{\mathbf{r}} \leftarrow \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{init}}$, $\hat{\boldsymbol{\beta}}^{\text{old}} \leftarrow \hat{\boldsymbol{\beta}}_{\text{init}}$
- 2 **for** $i = 1, \dots, I_{\max}$ **do**
- 3 **for** $j = 1$ **to** p **do**
- 4 $\hat{\beta}_j \leftarrow \frac{\mathcal{S}_{\alpha\lambda}(\hat{\beta}_j + \frac{1}{N}\langle \mathbf{x}_j, \hat{\mathbf{r}} \rangle)}{1 + \lambda(1 - \alpha)}$
- 5 $\hat{\mathbf{r}} \leftarrow \hat{\mathbf{r}} + (\hat{\beta}_j^{\text{old}} - \hat{\beta}_j)\mathbf{x}_j$
- 6 **if** $\|\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^{\text{old}}\|_2 / \|\hat{\boldsymbol{\beta}}\|_2 < \delta$ **then**
- 7 **break**
- 8 $\hat{\boldsymbol{\beta}}^{\text{old}} \leftarrow \hat{\boldsymbol{\beta}}$

Output : $\hat{\boldsymbol{\beta}}(\lambda, \alpha) = \hat{\boldsymbol{\beta}}$, minimizer of $\frac{1}{2N}\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda\{\alpha\|\boldsymbol{\beta}\|_1 + \frac{1}{2}(1 - \alpha)\|\boldsymbol{\beta}\|^2\}$

6.2 Generalized lasso

In many problems features (predictors) X_1, X_2, \dots, X_p can be ordered in some meaningful way, but lasso ignores such ordering. For example, neighboring coefficient values β_j can be piecewise constant over neighboring values and it makes sense to encourage both *block-sparsity* and *smoothness*. Such structure can be imposed using *generalized lasso*, which considers the problem

$$\underset{(\beta_0, \boldsymbol{\beta}) \in \mathbb{R} \times \mathbb{R}^p}{\text{minimize}} \left\{ \frac{1}{2} \sum_{i=1}^N (y_i - \beta_0 - \mathbf{x}_i^\top \boldsymbol{\beta})^2 + \lambda \|\mathbf{D}\boldsymbol{\beta}\|_1 \right\} \quad (6.5)$$

where $\mathbf{D} \in \mathbb{R}^{m \times p}$ is a specified *penalty matrix*. The conventional lasso regression is obtained when $\mathbf{D} = \mathbf{I}$. Note that in (6.5) we ignored the scaling constant $1/N$ from RSS criterion since this convention is more common in case of generalized lasso problems such as fused lasso explained next.

6.2.1 Fused lasso

Fused lasso (FL) penalty is defined as

$$\|\boldsymbol{\beta}\|_{\text{FL}} = \|\overline{\mathbf{D}}_p \boldsymbol{\beta}\|_1 = \sum_{j=1}^{p-1} |\beta_j - \beta_{j+1}|, \quad (6.6)$$

where $\overline{\mathbf{D}}_p$ is 1st order difference matrix, $\overline{\mathbf{D}}_p \in \mathbb{R}^{(p-1) \times p}$:

$$\overline{\mathbf{D}}_p = \begin{pmatrix} -1 & 1 & 0 & \cdots & 0 & 0 \\ 0 & -1 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & -1 & 1 \end{pmatrix}. \quad (6.7)$$

The fused lasso optimization problem is

$$\underset{(\beta_0, \boldsymbol{\beta}) \in \mathbb{R} \times \mathbb{R}^p}{\text{minimize}} \left\{ \frac{1}{2} \sum_{i=1}^N (y_i - \beta_0 - \mathbf{x}_i^\top \boldsymbol{\beta})^2 + \lambda \sum_{j=1}^{p-1} |\beta_j - \beta_{j+1}| \right\} \quad (6.8)$$

or in constrained form:

$$\underset{(\beta_0, \boldsymbol{\beta}) \in \mathbb{R} \times \mathbb{R}^p}{\text{minimize}} \frac{1}{2} \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \quad \text{s.t.} \quad \sum_{j=1}^{p-1} |\beta_j - \beta_{j+1}| \leq s$$

where $\lambda \geq 0$ (resp. $s \geq 0$) denotes the penalty (resp. constraint) parameter. Lasso penalty encourages sparsity in the coefficients, while FL penalty encourages sparsity in their differences, i.e., similarity of neighbouring coefficients as is illustrated in [Figure 6.3](#).

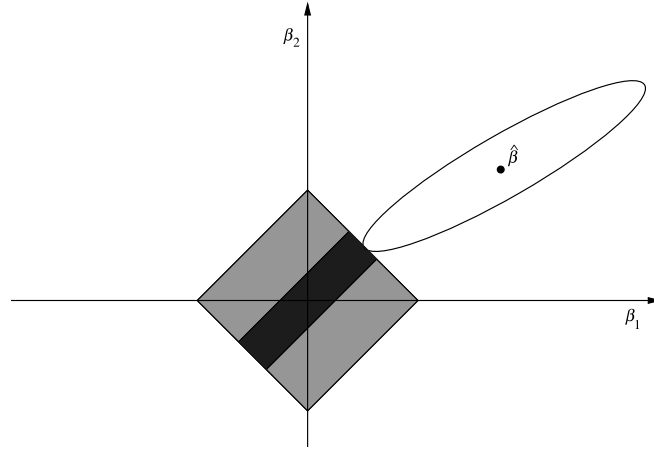


FIGURE 6.3: The geometry of FL in $p = 2$ dimension. The black shaded region denotes the region where the constraint $\|\boldsymbol{\beta}\|_{\text{FL}} = |\beta_1 - \beta_2| \leq s$ is active while the gray shaded region denotes the region where ℓ_1 -constraint $\|\boldsymbol{\beta}\|_1 = (|\beta_1| + |\beta_2|) \leq s$ is active.

Combining FL penalty with lasso penalty yields the *sparse fused lasso* (SFL) penalty:

$$P_{\text{SFL}}(\boldsymbol{\beta}; \lambda_1, \lambda_2) = \lambda_1 \|\boldsymbol{\beta}\|_1 + \lambda_2 \|\boldsymbol{\beta}\|_{\text{FL}}$$

where $\lambda_1, \lambda_2 \geq 0$ form a pair of fixed regularization parameters. The SFL optimization problem is then

$$\underset{(\beta_0, \boldsymbol{\beta}) \in \mathbb{R} \times \mathbb{R}^p}{\text{minimize}} \left\{ \frac{1}{2} \sum_{i=1}^N (y_i - \beta_0 - \mathbf{x}_i^\top \boldsymbol{\beta})^2 + \lambda_1 \|\boldsymbol{\beta}\|_1 + \lambda_2 \sum_{j=1}^{p-1} |\beta_j - \beta_{j+1}| \right\}. \quad (6.9)$$

SFL was proposed for regression by Tibshirani et al. [2005] but has longer history in image processing where it is called *total variation* (TV) penalty [Rudin et al., 1992].

Example 6.1. We generated an $N \times p$ predictor matrix \mathbf{X} of size $N = 1000$ and $p = 200$, where predictor variables have a joint multivariate Gaussian distribution with unit variance and $\rho = 0.7$ pairwise correlations. To demonstrate the usefulness of SFL regression, the true coefficient profile is piecewise constant with 50% of coefficients being exactly zero. The true coefficients β_i as a function of index i are shown in Figure 6.4 in dashed red color. The errors terms were generated as i.i.d. variables from standard Gaussian distribution (i.e., $\varepsilon_i \sim \mathcal{N}(0, 1)$). Then the outcome was generated as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

We then computed SFL solution using different values for parameter pair (λ_1, λ_2) . The results are shown in Figure 6.4. As can be noted, fused lasso ($\lambda_1 = 0$) is not able to correctly identify the non-significant (non-zero) coefficients, although its performance improves when increasing the value of 2nd penalty parameter λ_2 . On the contrary, SFL is able to correctly estimate the zero coefficients due to the used ℓ_1 penalty ($\lambda_1 \neq 0$). One can also notice that when λ_2 is fixed, and λ_1 is increased, all coefficients are shrunk towards zero. This is the familiar effect of ℓ_1 -penalty. ■

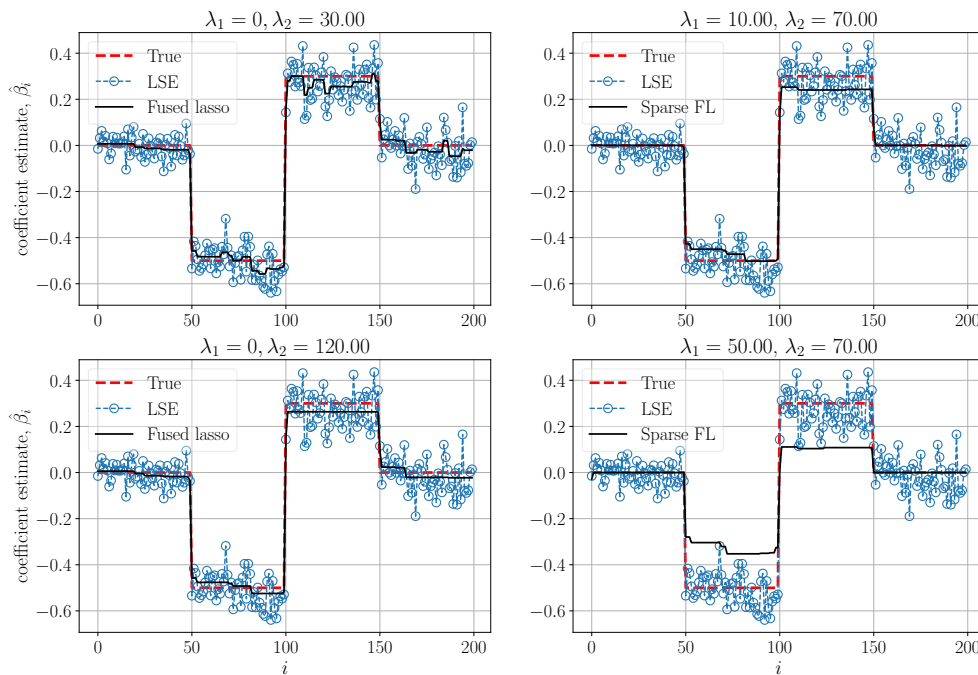


FIGURE 6.4: SFL solutions for data of Example 6.1 using different values of penalty parameter pairs (λ_1, λ_2) displayed in the figure titles. Setting: $p = 200$, $N = 1000$, and pairwise correlation between the predictors being 0.7.

Computation

Computation of FL or SFL solutions (6.8) or (6.9) is not as straightforward as EN/lasso:

- criterion is convex, but not differentiable, and coordinate descent can get stuck in the cusps.

- modification of CCD is required which involves moving pairs of parameters jointly, making it complex and computationally more involved.

One possible simpler approach is proximal gradient algorithm (PGA). The proximal operator of FL-penalty,

$$\text{prox}_{\lambda\|\cdot\|_{\text{FL}}}(\mathbf{z}) = \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{z} - \boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_{\text{FL}}, \quad (6.10)$$

has no closed-form solution, but it can be computed in linear time via a taut string method [Davies and Kovac, 2001] or dynamic programming (DP) [Johnson, 2013]. After selecting the method for evaluating the proximal operator, PGA iterations are

$$\hat{\boldsymbol{\beta}}^{(k)} = \text{prox}_{t_k \lambda \|\cdot\|_{\text{FL}}}(\hat{\boldsymbol{\beta}}^{(k-1)} + t_k \mathbf{X}^\top (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}^{(k-1)})). \quad (6.11)$$

The proximal operator of SFL penalty

$$\hat{\boldsymbol{\beta}}(\lambda_1, \lambda_2) = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{z} - \boldsymbol{\beta}\|_2^2 + \lambda_1 \|\boldsymbol{\beta}\|_1 + \lambda_2 \|\boldsymbol{\beta}\|_{\text{FL}}$$

can be computed by first evaluating the proximal map of FL penalty ($\lambda_1 = 0$), and then applying soft thresholding operator to that solution [Friedman et al., 2007]:

$$\hat{\boldsymbol{\beta}}(\lambda_1, \lambda_2) = \mathcal{S}_{\lambda_1}(\text{prox}_{\lambda_2 \|\cdot\|_{\text{FL}}}(\mathbf{z})).$$

Thus PGA for SFL is essentially the same as for FL: one simply applies soft-thresholding on top of (6.11).

Fused lasso extensions

There are several popular extensions of FL penalty. For example, FL can be generalized over a graph $\mathcal{G} = (\{1, \dots, p\}, E)$ with p nodes and edge set E by defining the penalty matrix \mathbf{D} as $|E| \times p$ matrix, whose ℓ th row is defined as

$$\mathbf{d}_\ell^\top = (0, \dots, \underset{\substack{\uparrow \\ i}}{-1}, \dots, \underset{\substack{\uparrow \\ j}}{1}, \dots, 0)$$

when (i, j) is an edge in the graph, so $(i, j) \in E$. Then

$$\|\mathbf{D}\boldsymbol{\beta}\|_1 = \sum_{(i,j) \in E} |\beta_i - \beta_j|.$$

The regression solution using FL penalty over graph is such that $\hat{\beta}_i \approx \hat{\beta}_j$ across the edges in the graph, so one encourage $\hat{\boldsymbol{\beta}}$ to be piecewise constant over the graph G .

Another extension of FL is the Fused ridge (FR) penalty:

$$\|\boldsymbol{\beta}\|_{\text{FR}}^2 = \sum_{i=1}^{p-1} (\beta_i - \beta_{i+1})^2. \quad (6.12)$$

6.2.2 Trend filtering

A special case of FL regression (with $\mathbf{X} = \mathbf{I}_{N \times N}$ and $p = N$) is *trend filtering* which is signal approximation problem that considers optimization problems of the form:

$$\underset{\boldsymbol{\beta} \in \mathbb{R}^N}{\text{minimize}} \left\{ \frac{1}{2} \sum_{i=1}^N (y_i - \beta_i)^2 + \lambda P(\boldsymbol{\beta}) \right\}, \quad (6.13)$$

where $P(\boldsymbol{\beta})$ is the penalty function and λ is the penalty parameter. Here $\{\beta_i\}_{i=1}^N$ is referred to as *signal*, so we consider a classic signal-in-noise measurement model

$$y_i = \beta_i + \varepsilon_i, \quad i = 1, \dots, N,$$

where only the corrupted measurements y_i -s are available but not the signal β_i itself.

The choice of the penalty depends on the assumed underlying signal shape. $P(\boldsymbol{\beta})$ is commonly chosen to be FL penalty $\|\cdot\|_{\text{FL}}$ when the signal β_i is piecewise constant or the SFL penalty leading to solving

$$\underset{\boldsymbol{\beta} \in \mathbb{R}^N}{\text{minimize}} \left\{ \frac{1}{2} \sum_{i=1}^N (y_i - \beta_i)^2 + \lambda_1 \|\boldsymbol{\beta}\|_1 + \lambda_2 \|\boldsymbol{\beta}\|_{\text{FL}} \right\}. \quad (6.14)$$

Fused ridge penalty $\|\cdot\|_{\text{FR}}^2$ in (6.12) on the other hand works better when signal is smoother.

It is important to notice that trend filtering is tantamount to evaluating the proximal map of the chosen penalty. For FL penalty for example the computation can be done efficiently using taut string or DP method. For general trend filtering matrix, specialized interior point methods [Kim et al., 2009] or ADMM methods [Ramdas and Tibshirani, 2016] can be used.

An example of trend filtering using SFL signal approximator is shown in Figure 6.5.

Example 6.2. We consider two cases:

- (A) signal β_i is a piecewise constant signal shown in Figure 6.5a with dashed red color.
- (B) signal is a superposition of two sine waves

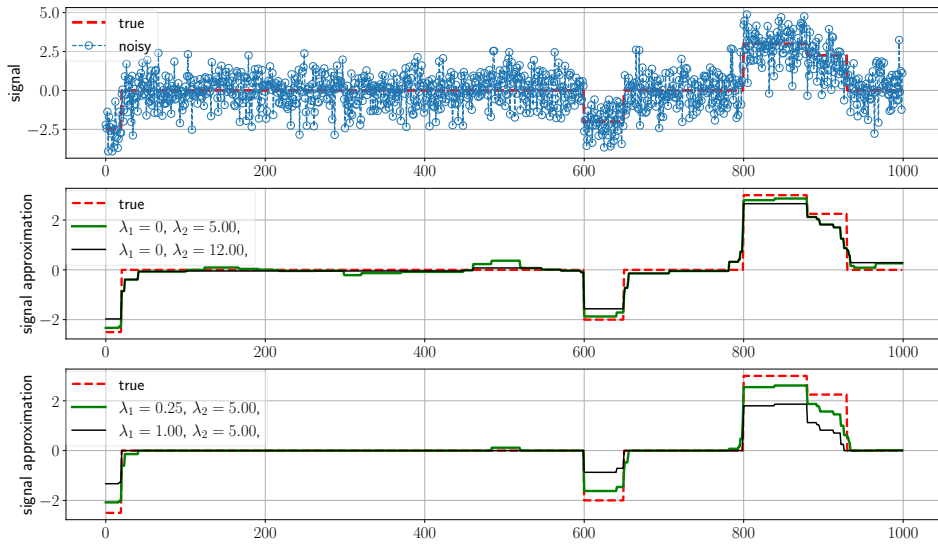
$$\beta_i = \sin((i-1)2\pi f_1) + \sin((i-1)2\pi f_2)$$

with frequencies $f_1 = 0.15$ and $f_2 = f_1/10 = 0.015$, shown in Figure 6.5b with dashed red color.

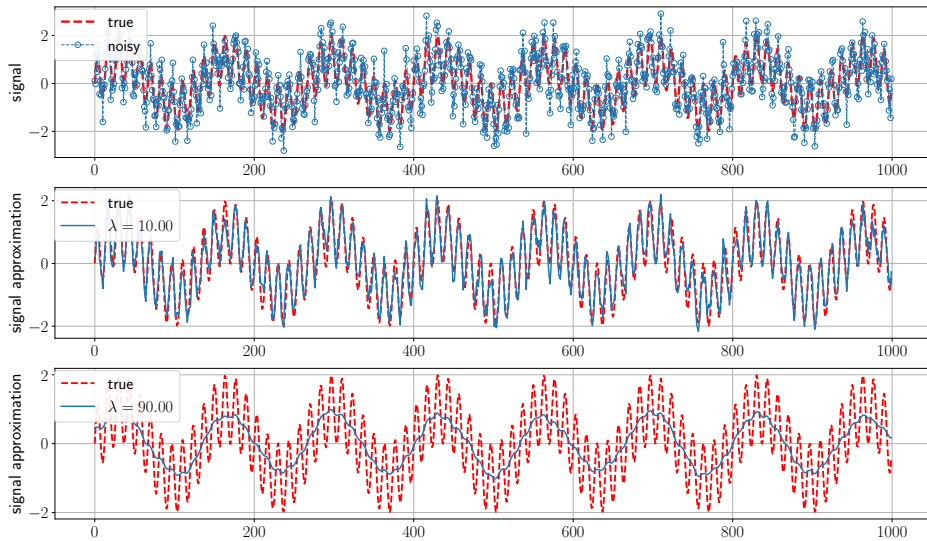
Signals are measured in additive white Gaussian noise: $\varepsilon \sim \mathcal{N}(0, 1)$ for case (A) and $\varepsilon \sim \mathcal{N}(0, 0.25)$ for case (B). Measurements y_i are then generated as

$$y_i = \beta_i + \varepsilon_i, \quad i = 1, \dots, N, \quad (6.15)$$

where the sample length is $N = 1000$. For case (A) we computed the SFL signal approximator (??) using different values options for parameter pair (λ_1, λ_2) . The results are shown in Figure 6.5a. Similar conclusions can be drawn as in the regression setup.



(A) piecewise constant signal in noise



(B) superposition of sinusoids in noise

FIGURE 6.5: (a) Results of SFL signal approximator (6.14) with different values of (λ_1, λ_2) for case (A) of Example 6.2. (b) Results of FR signal approximator with different values of λ for case (B) of Example 6.2.

FL ($\lambda_1 = 0$) performs worse than SFL ($\lambda_1 \neq 0$) in parts where true signal is exactly zero but better in parts where signal is non-zero due to shrinkage effect of lasso penalty.

For case (B) we computed the fused ridge (FR) signal approximator, i.e., solved (6.13) using penalty (6.12). Figure 6.5b displays the results using two choices of penalty parameter values. As can be noted, too large penalty ($\lambda = 90$) causes the FR approximator to adapt to low frequency component disregarding the higher frequency component. However, a better (i.e., smaller) penalty parameter value ($\lambda = 10$) is able to catch the

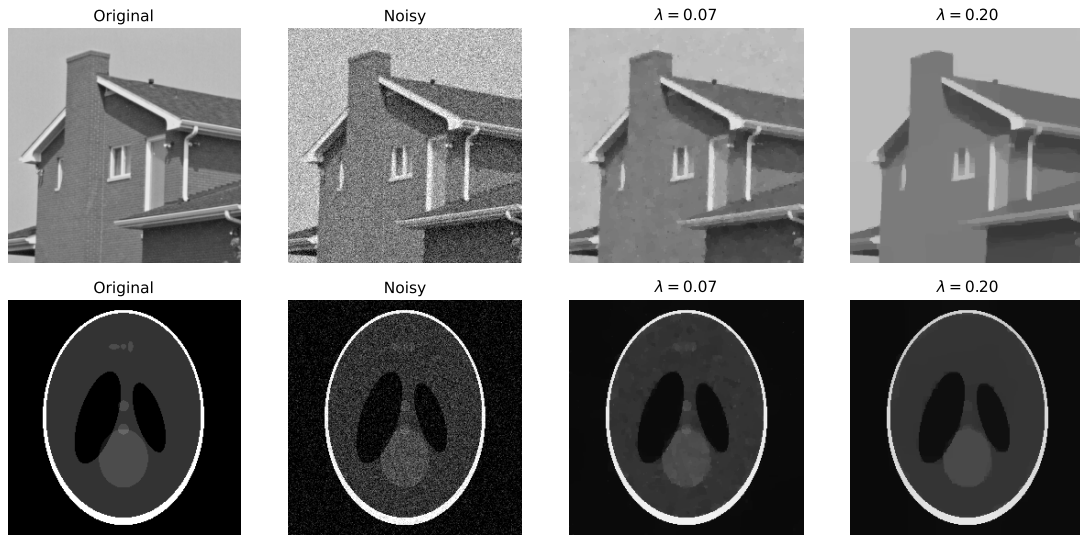


FIGURE 6.6: The original, noisy, and denoised images using total variation denoising with $\lambda = 0.07$ and $\lambda = 0.2$ for *house* and *phantom* images.

underlying signal structure very well. ■

Image denoising

As already mentioned, FL has been originally proposed in image denoising where it is called total variation (TV) denoising. In TV denoising one solves

$$\underset{\mathcal{B} \in \mathbb{R}^{N_1 \times N_2}}{\text{minimize}} \left\{ \frac{1}{2} \sum_{i=1}^N (y_{i,j} - \beta_{i,j})^2 + \lambda \sum_{i=2}^{N_1} \sum_{j=1}^{N_2} |\beta_{i,j} - \beta_{i-1,j}| + \lambda \sum_{i=1}^{N_1} \sum_{j=2}^{N_2} |\beta_{i,j} - \beta_{i,j-1}| \right\}. \quad (6.16)$$

Above $\mathbf{Y} = (y_{i,j}) \in \mathbb{R}^{N_1 \times N_2}$ is a 2D-image and $\mathcal{B} = (\beta_{i,j})$ is the denoised image, and λ is the penalty term. The idea is to enforce smoothness of neighborhood pixels both in horizontal and vertical directions of the image.

Example 6.3. We denoise two 256×256 gray scale images labelled *house* and *phantom*. Both images are scaled to scale $[0, 1]$. The latter is a standard test image ("Shepp-Logan phantom") commonly used for testing image reconstruction algorithms in medical imaging applications¹ while the former image is a standard image used for testing ability of image reconstruction algorithms in maintaining sharp edges. We added Gaussian random noise with variance $\sigma^2 = 0.01$ to the images and denoised images using (6.16) with two choices of λ . The results are displayed in Figure 6.6. As can be noted, the noise in images is clearly reduced. One can also notice that too large value of λ causes too much smoothing, making some fine details in image disappear. ■

¹https://en.wikipedia.org/wiki/Shepp-Logan_phantom

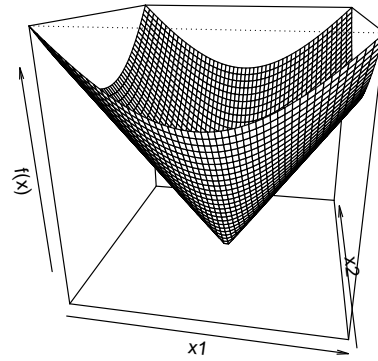
6.3 Group lasso

Group lasso is defined as the following optimization problem:

$$\underset{\beta_g \in \mathbb{R}^{p_g}}{\text{minimize}} \quad \frac{1}{2} \left\| \mathbf{y} - \sum_{g=1}^G \mathbf{X}_g \beta_g \right\|_2^2 + \lambda \sum_{g=1}^G \sqrt{p_g} \|\beta_g\|_2, \quad (6.17)$$

where $\mathbf{X}_g \in \mathbb{R}^{N \times p_g}$ is the data matrix corresponding to covariates in group g , β_g is the regression coefficients corresponding to group g , p_g is the dimensionality (number of covariates) of group g and G is the number of groups. As earlier, $\mathbf{y} \in \mathbb{R}^N$ denotes the regression response vector, N is the number of measurements, and $\lambda \geq 0$ the penalty parameter.

Notice also that in (6.17) the ℓ_2 -norm penalty is not squared. Thus it is not differentiable at zero, making it have a sharp edge at 0 as is illustrated on the right hand side figure.



This leads it to have attributes that are similar to lasso:

1. For large enough $\lambda > 0$, the entire vector β_g will be zero or all coefficients are nonzero.
2. if $p_g \equiv 1$ for all g , so we have a single covariate in each group, then the problem reduces to ordinary lasso.

Note that each group penalty is weighted according to their size, so by $\sqrt{p_g}$. However, in the original formulation of group lasso [Yuan and Lin, 2006] matrices \mathbf{X}_g were assumed to be orthogonal, so $\mathbf{X}_g^\top \mathbf{X}_g = \mathbf{I}$. However, for general matrices, Frobenius norm $\|\mathbf{X}_g\|_{\text{Fr}}$ can be used instead. Notice that $\|\mathbf{X}_g\|_{\text{Fr}} = \text{tr}(\mathbf{X}_g^\top \mathbf{X}_g) = \sqrt{p_g}$.

Some example applications where group lasso penalty is particularly useful:

1. The levels of qualitative factors are typically coded using a set of dummy variables and one would want to include or exclude this group of variables together.
2. In gene-expression arrays, genes from the same biological pathway can be highly correlated, and selecting them as a group corresponds to electing a pathway.

Computation

Subdifferential of $\|\beta\|_2$ is

$$\partial \|\beta\|_2 = \begin{cases} \beta / \|\beta\|_2 & \text{for } \beta \neq 0 \\ \{\mathbf{s} \in \mathbb{R}^p : \|\mathbf{s}\|_2 \leq 1\} & \text{for } \beta = 0 \end{cases}$$

The zero subgradient equations are given by:

$$-\mathbf{X}_g^\top \left(\mathbf{y} - \sum_{g=1}^G \mathbf{X}_g \hat{\boldsymbol{\beta}}_g \right) + \lambda \sqrt{p_g} \hat{\mathbf{s}}_g = \mathbf{0} \quad \text{for } g = 1, \dots, G,$$

where $\hat{\mathbf{s}}_g \in \mathbb{R}^{p_g}$ is an element of subdifferential of $\|\cdot\|_2$ evaluated at $\hat{\boldsymbol{\beta}}_g$.

One obvious approach for solving the zero subgradient equations is *block coordinate descent (BCD)* where one holds fixed all block vectors except one and then solve for the fixed block, repeating this for each block in turn. Since the problem is convex, and the penalty is block separable, it is guaranteed to converge to an optimal solution. For all but j^{th} block fixed, the zero subgradient equation is

$$-\mathbf{X}_j^\top (\mathbf{r}_j - \mathbf{X}_j \hat{\boldsymbol{\beta}}_j) + \lambda \sqrt{p_j} \hat{\mathbf{s}}_j = \mathbf{0}$$

where $\mathbf{r}_j = \mathbf{y} - \sum_{g \neq j}^G \mathbf{X}_g \hat{\boldsymbol{\beta}}_g$ is the j th partial residual. This equation has a simple closed-form solution when \mathbf{X}_g -s are orthonormal, given by

$$\hat{\boldsymbol{\beta}}_j = \left(1 - \frac{\lambda \sqrt{p_j}}{\|\mathbf{X}_j^\top \mathbf{r}_j\|_2} \right)_+ \mathbf{X}_j^\top \mathbf{r}_j \quad (6.18)$$

where $(t)_+ = \max\{0, t\}$ is the positive part of the function.

6.4 Discussion

The purpose of this chapter was to give a quick look at many extensions of lasso. However, many important extensions were not discussed such as Bayesian lasso [Park and Casella, 2008], adaptive lasso [Zou, 2006], lasso using nonconvex penalties such as smoothly clipped absolute deviation (SCAD) penalty [Fan and Li, 2001] or minimax concave penalty [Zhang, 2010], etc, etc.

Bibliography

- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- Leo Breiman, Jerome Friedman, Charles J Stone, and Richard A Olshen. *Classification and regression trees*. CRC press, 1984.
- Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, 55(1):119–139, 1997.
- Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001.
- Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media, 2 edition, 2009.
- Robert E Schapire. The strength of weak learnability. *Machine learning*, 5(2):197–227, 1990.
- Yoav Freund and Robert E Schapire. Experiments with a new boosting algorithm. In *Icml*, volume 96, pages 148–156. Citeseer, 1996.
- Jerome Friedman, Trevor Hastie, Robert Tibshirani, et al. Additive logistic regression: a statistical view of boosting (with discussion and a rejoinder by the authors). *The annals of statistics*, 28(2):337–407, 2000.
- Robert E Schapire and Yoram Singer. Improved boosting algorithms using confidence-rated predictions. *Machine learning*, 37(3):297–336, 1999.
- Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 785–794, 2016.
- A. E. Hoerl and R. W. Kennard. Ridge regression: biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970.
- Trevor Hastie, Robert Tibshirani, and Martin Wainwright. *Statistical Learning with Sparsity: The Lasso and Generalizations*. CRC Press, 2015.

- Paul Tseng. Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of optimization theory and applications*, 109(3):475–494, 2001.
- J. Friedman, T. Hastie, H. Höfling, and R. Tibshirani. Pathwise coordinate optimization. *Ann. Appl. Stat.*, 1(2):302–332, 2007.
- R. Tibshirani and et al. Strong rules for discarding predictors in lasso-type problems. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 74(2):245–266, 2012.
- Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *J. Royal Stat. Soc., Ser. B*, 67(2):301–320, 2005.
- M. Tibshirani, R. Saunders, S. Rosset, J. Zhu, and K. Knight. Sparsity and smoothness via the fused lasso. *J. Royal Stat. Soc., Ser. B*, 67(1):91–108, 2005.
- L. I. Rudin, S. Osher, and E. Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: Nonlinear Phenomena*, 60(1):259–268, 1992.
- P Laurie Davies and Arne Kovac. Local extremes, runs, strings and multiresolution. *The Annals of Statistics*, 29(1):1–65, 2001.
- Nicholas A Johnson. A dynamic programming algorithm for the fused lasso and l0-segmentation. *Journal of Computational and Graphical Statistics*, 22(2):246–260, 2013.
- Seung-Jean Kim, Kwangmoo Koh, Stephen Boyd, and Dmitry Gorinevsky. ℓ_1 trend filtering. *SIAM review*, 51(2):339–360, 2009.
- Aaditya Ramdas and Ryan J Tibshirani. Fast and flexible admm algorithms for trend filtering. *Journal of Computational and Graphical Statistics*, 25(3):839–858, 2016.
- M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *J. Royal Stat. Soc., Ser. B*, 68(1):49–67, 2006.
- Trevor Park and George Casella. The bayesian lasso. *Journal of the American Statistical Association*, 103(482):681–686, 2008.
- Hui Zou. The adaptive lasso and its oracle properties. *J. Amer. Stat. Assoc.*, 101(476):1418–1429, 2006.
- Jianqing Fan and Runze Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association*, 96(456):1348–1360, 2001.
- Cun-Hui Zhang. Nearly unbiased variable selection under minimax concave penalty. *The Annals of statistics*, 38(2):894–942, 2010.