

# Empirical Risk Minimization (ERM)

Kun He (何琨)

Data Mining and Machine Learning Lab  
(John Hopcroft Lab)  
Huazhong University of Science & Technology

*brooklet60@hust.edu.cn*

2022年5月



# Table of Contents

## 1 Definition

- Recap
- Background
- Empirical Risk Minimization

## 2 Binary Classification Loss Functions

## 3 Regression Loss Functions

## 4 Regularizers

## 5 Famous Special Cases

# Table of Contents

## 1 Definition

- Recap
- Background
- Empirical Risk Minimization

## 2 Binary Classification Loss Functions

## 3 Regression Loss Functions

## 4 Regularizers

## 5 Famous Special Cases

Remember the unconstrained SVM Formulation

$$\min_{\mathbf{w}} \underbrace{C \sum_{i=1}^n \max[1 - y_i(\underbrace{\mathbf{w}^\top \mathbf{x}_i + b}_{h(\mathbf{x}_i)}, 0]}_{\text{Hinge-Loss}} + \underbrace{\|\mathbf{w}\|_z^2}_{\ell_2\text{-Regularizer}}$$

The hinge loss is the SVM's error function of choice, whereas the  $\ell_2$ -regularizer reflects the complexity of the solution, and penalizes complex solutions. This is an example of empirical risk minimization with a loss function  $\ell$  and a regularizer  $r$ ,

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \underbrace{l(h_{\mathbf{w}}(\mathbf{x}_i), y_i)}_{\text{Loss}} + \underbrace{\lambda r(\mathbf{w})}_{\text{Regularizer}},$$

where the loss function is a continuous function which penalizes training error, and the regularizer is a continuous function which penalizes classifier complexity. Here, we define  $\lambda$  as  $\frac{1}{C}$  from the previous lecture.

# Background

- Consider the following situation, a general setting of many supervised learning problems.
- We have two spaces of objects  $X$  and  $Y$  and would like to learn a function  $h : X \rightarrow Y$  which outputs an object  $y \in Y$ , given  $x \in X$ .
- To do so, we have at our disposal a training set of  $n$  examples  $(x_1, y_1), \dots, (x_n, y_n)$  where  $x_i \in X$  is an input and  $y_i \in Y$  is the corresponding response that we wish to get from  $h(x_i)$ .
- To put it more formally, we assume that there is a **joint probability distribution**  $p(x, y)$  over  $X$  and  $Y$ , and that the training set consists of  $n$  instances  $(x_1, y_1), \dots, (x_n, y_n)$  drawn i.i.d. from  $p(x, y)$ .
- Note that the assumption of a joint probability distribution allows us to model uncertainty in predictions because  $y$  is not a deterministic function of  $x$ , but rather a random variable with conditional distribution  $p(x|y)$  for a fixed  $x$ .

## Background

We also assume that we are given a **non-negative real-valued loss function**  $L(\hat{y}, y)$  which measures how different the prediction  $\hat{y}$  of a hypothesis is from the true outcome  $y$ . The risk associated with hypothesis  $h(x)$  is then defined as the expectation of the loss function:

$$R(h) = E[L(h(x), y)] = \int L(h(x), y) dP(x, y)$$

A loss function commonly used in theory is the 0-1 loss function:  $L(\hat{y}, y) = I(\hat{y} \neq y)$ , where  $I(\cdot)$  is the indicator notation.

The ultimate goal of a learning algorithm is to find a hypothesis  $h^*$  among a fixed class of functions  $H$  for which the risk  $R(h)$  is minimal:

$$h^* = \underset{h \in H}{\operatorname{argmin}} R(h).$$

# Empirical Risk Minimization

In general, the risk  $R(h)$  cannot be computed because the distribution  $P(x, y)$  is unknown to the learning algorithm (this situation is referred to as agnostic learning). However, we can compute an approximation, called empirical risk, by averaging the loss function on the training set:

$$R_{emp}(h) = \frac{1}{n} \sum_{i=1}^n L(h(x_i), y_i).$$

The empirical risk minimization principle states that the learning algorithm should choose a hypothesis  $\hat{h}$  which minimizes the empirical risk:

$$h^* = \underset{h \in H}{\operatorname{argmin}} R(h).$$

Thus **Empirical Risk Minimization (ERM)** is a principle in statistical learning theory which defines a family of learning algorithms and is used to give theoretical bounds on their performance.

# Table of Contents

- 1 Definition
  - Recap
  - Background
  - Empirical Risk Minimization
- 2 Binary Classification Loss Functions
- 3 Regression Loss Functions
- 4 Regularizers
- 5 Famous Special Cases

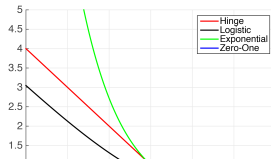


# Binary Classification Loss Functions

Different Machine Learning algorithms use different loss functions

Table 1: Binary Classification Loss Functions,  $y \in \{-1, +1\}$

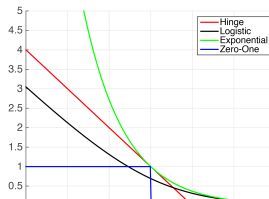
Loss $\ell(h_{\mathbf{w}}(\mathbf{x}_i, y_i))$	Usage	Comments
<b>Hinge-Loss:</b> $\max[1 - h_{\mathbf{w}}(\mathbf{x}_i)y_i, 0]^p$	<ul style="list-style-type: none"><li>• Standard SVM (<math>p = 1</math>)</li><li>• (Differentiable) Squared Hingeless SVM (<math>p = 2</math>)</li></ul>	When used for Standard SVM, the loss function denotes the size of the margin between linear separator and its closest points in either class. Only differentiable everywhere with $p = 2$ .
<b>Log-Loss:</b> $\log(1 + e^{-h_{\mathbf{w}}(\mathbf{x}_i)y_i})$	Logistic Regression	One of the most popular loss functions in Machine Learning, since its outputs are well-calibrated probabilities.



# Binary Classification Loss Functions

Table 1: Binary Classification Loss Functions,  $y \in \{-1, +1\}$

Loss $\ell(h_{\mathbf{w}}(\mathbf{x}_i, y_i))$	Usage	Comments
<b>Exponential Loss:</b> $e^{-h_{\mathbf{w}}(\mathbf{x}_i)y_i}$	AdaBoost	This function is very aggressive. The loss of a mis-prediction increases <i>exponentially</i> with the value of $-h_{\mathbf{w}}(\mathbf{x}_i)y_i$ . This can lead to nice convergence results, for example in the case of Adaboost, but it can also cause problems with noisy data.
<b>Zero-One Loss:</b> $\delta(\text{sign}(h_{\mathbf{w}}(\mathbf{x}_i)) \neq y_i)$	Actual Classification Loss	Non-continuous and thus impractical to optimize.



# Binary Classification Loss Functions

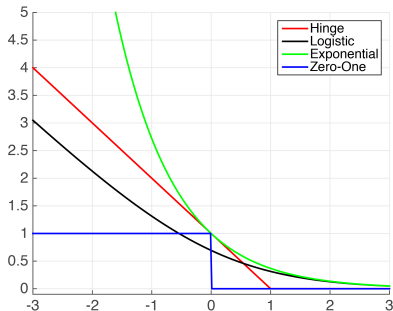


Figure 1: Plots of Common Classification Loss Functions. x-axis:  $h(\mathbf{x}_i)y_i$ , or "correctness" of prediction; y-axis: loss value

**Quiz:** Which functions are strict upper bounds on the 0/1-loss?

**Quiz:** What can you say about the hinge-loss and the log-loss as  $z \rightarrow -\infty$ ?

# Table of Contents

## 1 Definition

- Recap
- Background
- Empirical Risk Minimization

## 2 Binary Classification Loss Functions

## 3 Regression Loss Functions

## 4 Regularizers

## 5 Famous Special Cases

# Regression Loss Functions

Loss for regression algorithms (a prediction can lie anywhere on the real-number line):

Table 2: Regression Loss Functions,  $y \in \mathbb{R}$

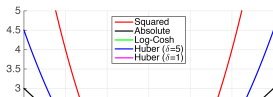
Loss $\ell(h_{\mathbf{w}}(\mathbf{x}_i, y_i))$	Comments
<b>Squared Loss:</b> $(h(\mathbf{x}_i) - y_i)^2$	<ul style="list-style-type: none"><li>• Most popular regression loss function</li><li>• Estimates <u>Mean</u> Label</li><li>• Also known as Ordinary Least Squares (OLS)</li><li>• DISADVANTAGE: Somewhat sensitive to outliers/noise</li></ul>
<b>Absolute Loss:</b> $ h(\mathbf{x}_i) - y_i $	<ul style="list-style-type: none"><li>• Also a very popular loss function</li><li>• Estimates <u>Median</u> Label</li><li>• ADVANTAGE: Less sensitive to noise</li><li>• DISADVANTAGE: Not differentiable at 0</li></ul>



# Regression Loss Functions

Table 2: Regression Loss Functions,  $y \in \mathbb{R}$

Loss $\ell(h_{\mathbf{w}}(\mathbf{x}_i, y_i))$	Comments
<b>Huber Loss:</b> <ul style="list-style-type: none"><li><math>\frac{1}{2} (h(\mathbf{x}_i) - y_i)^2</math> if <math> h(\mathbf{x}_i) - y_i  &lt; \delta</math>,</li><li>otherwise <math>\delta( h(\mathbf{x}_i) - y_i  - \frac{\delta}{2})</math></li></ul>	<ul style="list-style-type: none"><li>Also known as Smooth Absolute Loss</li><li>Once-differentiable</li><li>Takes on behavior of Squared-Loss when loss is small, and Absolute Loss when loss is large.</li><li>ADVANTAGE: "Best of Both Worlds" of <u>Squared</u> and <u>Absolute</u> Loss</li></ul>
<b>Log-Cosh Loss:</b> $\log(\cosh(h(\mathbf{x}_i) - y_i))$ , $\cosh(x) = \frac{e^x + e^{-x}}{2}$	ADVANTAGE: Similar to Huber Loss, but twice differentiable everywhere



# Regression Loss Functions

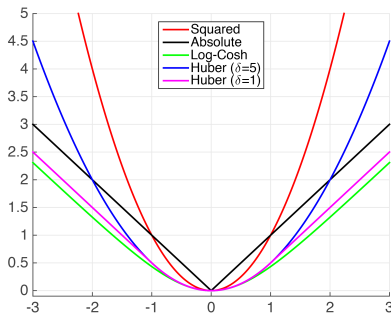


Figure 2: Plots of Common Regression Loss Functions - x-axis:  $h(\mathbf{x}_i)y_i$  , or "error" of prediction; y-axis: loss value

**Quiz:** What do the loss functions in Table 2 look like with respect to  $z = h(\mathbf{x}_i) - y_i$  ?

# Table of Contents

## 1 Definition

- Recap
- Background
- Empirical Risk Minimization

## 2 Binary Classification Loss Functions

## 3 Regression Loss Functions

## 4 Regularizers

## 5 Famous Special Cases



- In mathematics, statistics, and computer science, particularly in machine learning problems, regularization is the process of adding information in order to solve an ill-posed problem or to prevent overfitting.
- When we look at regularizers it helps to change the formulation of the optimization problem to obtain a better geometric intuition.
- In previous sections,  $l_2$ -regularizer has been introduced as the component in SVM that reflects the complexity of solutions.

Besides the  $l_2$ -regularizer, other types of useful regularizers and their properties are listed in the following Table.

Table 3: Types of Regularizers

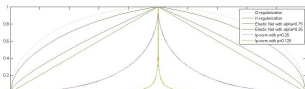
Regularizer $r(\mathbf{w})$	Properties
$l_2$ -Regularization $r(\mathbf{w}) = \mathbf{w}^\top \mathbf{w} = \ \mathbf{w}\ _2^2$	<ul style="list-style-type: none"><li>• ADVANTAGE: Strictly Convex; Differentiable</li><li>• DISADVANTAGE: Uses weights on all features, i.e. relies on all features to some degree (ideally we would like to avoid this) - these are known as <u>Dense Solutions</u>.</li></ul>
$l_1$ -Regularization $r(\mathbf{w}) = \ \mathbf{w}\ _1$	<ul style="list-style-type: none"><li>• Convex (but not strictly)</li><li>• DISADVANTAGE: Not differentiable at 0 (the point which minimization is intended to bring us to)</li><li>• Effect: <u>Sparse</u> (i.e. not <u>Dense</u>) Solutions</li></ul>



# Regularizers

Table 3: Types of Regularizers

Regularizer $r(\mathbf{w})$	Properties
$l_p$ -Norm $\ \mathbf{w}\ _p = \left(\sum_{i=1}^d v_i^p\right)^{1/p}$	<ul style="list-style-type: none"><li>• (often <math>0 &lt; p \leq 1</math>)</li><li>• Initialization dependent</li><li>• ADVANTAGE: Very sparse solutions</li><li>• DISADVANTAGE: Non-convex; Not differentiable</li></ul>
<b>Elastic Net:</b> $\alpha\ \mathbf{w}\ _1 + (1 - \alpha)\ \mathbf{w}\ _2^2$ $\alpha \in [0, 1]$	<ul style="list-style-type: none"><li>• ADVANTAGE: Strictly convex (i.e. unique solution)</li><li>• ADVANTAGE: Sparsity inducing (good for feature selection); Dual of squared-loss SVM, see <b>SVEN</b></li><li>• DISADVANTAGE: Non-differentiable</li></ul>



# Regularizers

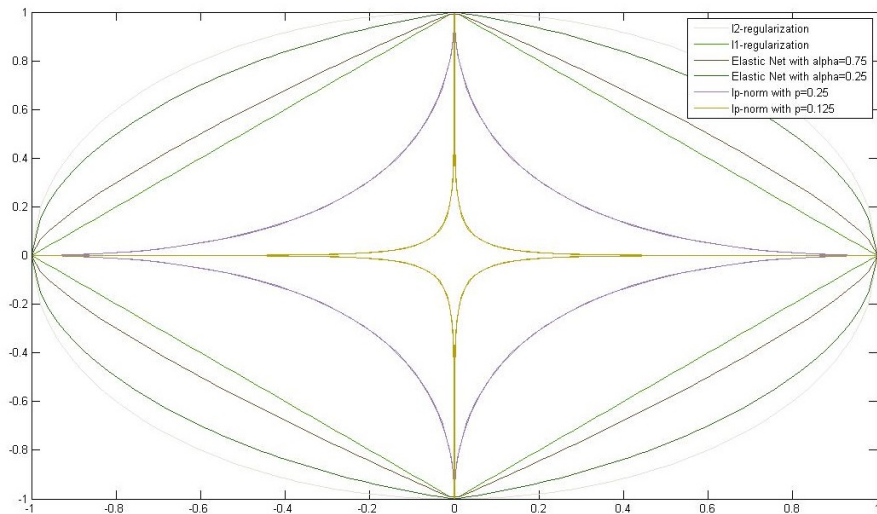


Figure 3: Plots of Common Regularizers

# Table of Contents

## 1 Definition

- Recap
- Background
- Empirical Risk Minimization

## 2 Binary Classification Loss Functions

## 3 Regression Loss Functions

## 4 Regularizers

## 5 Famous Special Cases

## Famous Special Cases

This section includes several special cases that deal with risk minimization, such as Ordinary Least Squares, Ridge Regression, Lasso, and Logistic Regression. The following blocks provide information on their loss functions, regularizers, as well as solutions.

### Ordinary Least Squares

In statistics, ordinary least squares (OLS) is a type of linear least squares method for estimating the unknown parameters in a linear regression model.

**goal:**

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$

**Comments:**

- Squared Loss
- No Regularization
- Closed form solution:  $\mathbf{w} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y}^\top$ ,  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ ,  $\mathbf{y} = [y_1, \dots, y_n]$
- There is an interesting connection between OLS and the first principal component of PCA (Principal Component Analysis). PCA also minimizes square loss, but looks at perpendicular loss (the horizontal distance between each point and the regression

## Ridge Regression

Tikhonov regularization, is the most commonly used method of regularization of ill-posed problems. In statistics, the method is known as **ridge regression**.

**goal:**

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_2^2$$

**Comments:**

- Squared Loss
- $l_2$ -Regularization
- $\mathbf{w} = (\mathbf{X}\mathbf{X}^\top + \lambda\mathbf{I})^{-1}\mathbf{X}\mathbf{y}^\top$
- Fast if data isn't too high dimensional
- Just 1 line of Julia / Python

## Lasso

**LASSO** is a regression analysis method that performs both variable selection and regularization in order to enhance the prediction accuracy and interpretability of the statistical model it produces.

### Goal:

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_1$$

### Comments:

- ADVANTAGE sparsity inducing (good for feature selection); Convex
- DISADVANTAGE Not strictly convex (no unique solution); Not differentiable (at 0)
- Solve with (sub)-gradient descent or **SVEN**



## Elastic Net

### Goal

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \alpha \|\mathbf{w}\|_1 + (1 - \alpha) \|\mathbf{w}\|_2^2$$
$$\alpha \in [0, 1)$$

### Comments:

- ADVANTAGE: Strictly convex (i.e. unique solution)
- +Sparsity inducing (good for feature selection)
- +Dual of squared-loss SVM, see **SVEN**
- DISADVANTAGE - Non-differentiable
- Solve with (sub)-gradient descent or **SVEN**

## Logistic Regression

In regression analysis, **logistic regression** is estimating the parameters of a logistic model; it is a form of binomial regression.

**goal:**

$$\min_{\mathbf{w}, b} \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)})$$

**Comments:**

- Often  $l_1$  or  $l_2$  regularized
- Solve with gradient descent.
- $\Pr(y|x) = \frac{1}{1 + e^{-y(\mathbf{w}^\top \mathbf{x} + b)}}$

## Linear Support Vector Machine

goal:

$$\min_{\mathbf{w}, b} C \sum_{i=1}^n \max[1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b), 0] + \|\mathbf{w}\|_2^2$$

Comments:

- Typically  $l_2$  regularized (sometimes  $l_1$ ).
- Quadratic program.
- When kernelized leads to **sparse** solutions.
- Kernelized version can be solved very efficiently with specialized algorithms (e.g. **SMO**)

# The End