

The SVM with soft constraints optimization problem:

$$\min_{\vec{w}, b} \underbrace{W^T W}_{\text{regularizer}} + C \sum_{i=1}^n \underbrace{\max(1 - \gamma_i (\vec{w}^T x + b), 0)}_{\text{loss function}}$$

- The above problem is in a form of regularizer term and loss function term
- The loss function here is "hinge-loss", and almost differentiate everywhere. Therefore, the gradient descent can be applied to find the solution.

## Empirical Risk Minimization

- A lot of ML algorithms can be written in a form of optimization problem with the objective to minimize a function  $l$  and a regularizer  $r$ .

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

- The loss function is a continuous function penalizing training error.

- The regularizer is a continuous function penalizing classifier complexity.

## Commonly Used Binary Classification Loss Functions

Different Machine Learning algorithms use different loss functions; Table 4.1 shows just a few:

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.
<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.

Table 4.1: Loss Functions With Classification  $y \in \{-1, +1\}$

Quiz: What do all these loss functions look like with respect to  $z = yh(\mathbf{x})$ ?

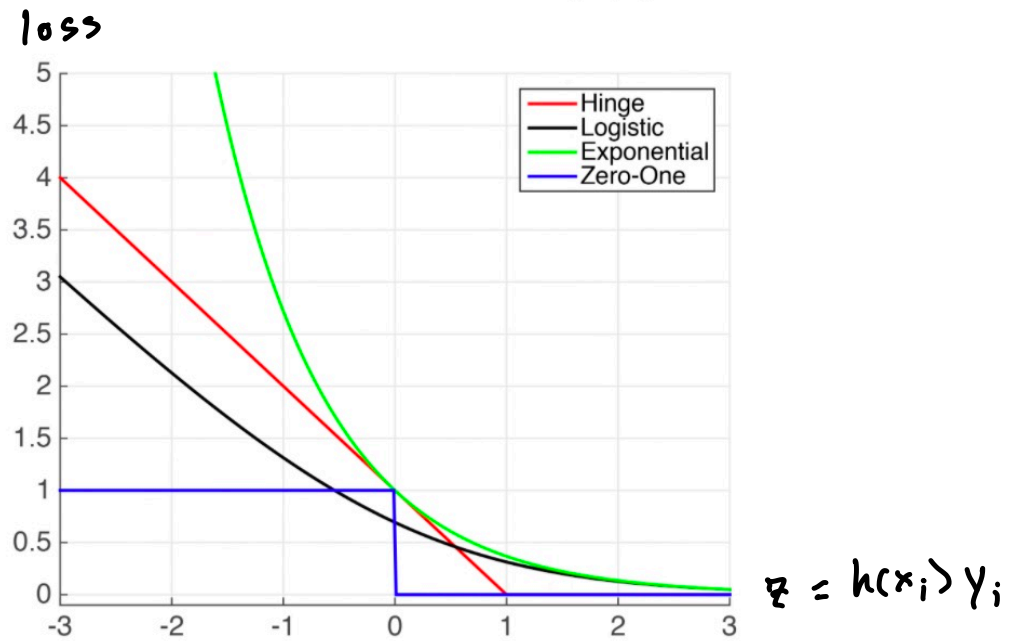


Figure 4.1: Plots of Common Classification Loss Functions - x-axis:  $h(\mathbf{x}_i)y_i$ , or "correctness" of prediction; y-axis: loss value

Some questions about the loss functions:

1. Which functions are strict upper bounds on the 0/1-loss?
2. What can you say about the hinge-loss and the log-loss as  $z \rightarrow -\infty$ ?

# Commonly Used Regression Loss Functions

Regression algorithms (where a prediction can lie anywhere on the real-number line) also have their own host of loss functions:

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>ADVANTAGE: Differentiable everywhere</li><li>DISADVANTAGE: Somewhat sensitive to outliers/noise</li><li>Also known as Ordinary Least Squares (OLS)</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>
<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>ADVANTAGE: "Best of Both Worlds" of <u>Squared</u> and <u>Absolute</u> 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></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

Table 4.2: Loss Functions With Regression, i.e.  $y \in \mathbb{R}$

Quiz: What do the loss functions in Table 4.2 look like with respect to  $z = h(\mathbf{x}_i) - y_i$ ?

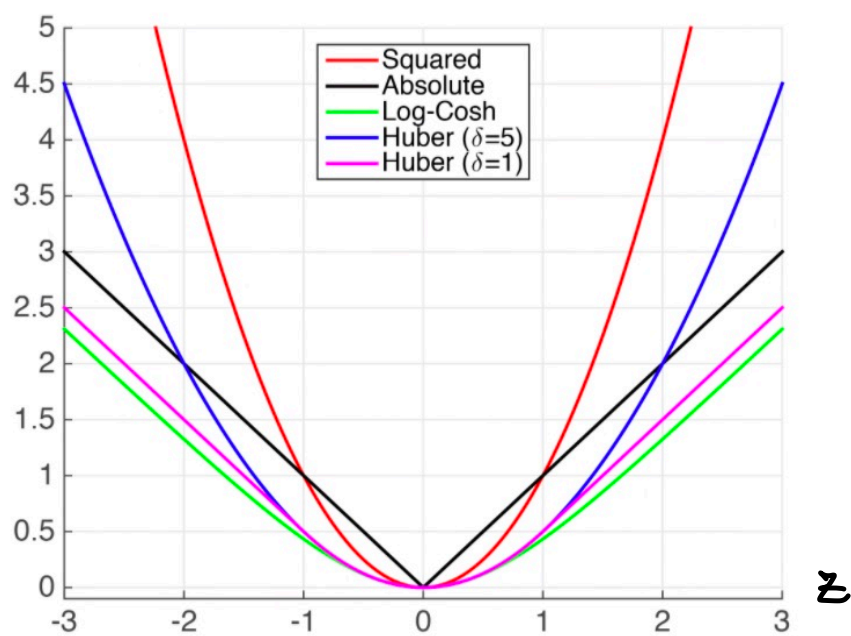


Figure 4.2: Plots of Common Regression Loss Functions - x-axis:  $h(\mathbf{x}_i) - y_i$ , or "error" of prediction; y-axis: loss value

## Regularizers

When we look at regularizers it helps to change the formulation of the optimization problem to obtain a better geometric intuition:

$$\min_{\mathbf{w}, b} \sum_{i=1}^n \ell(h_{\mathbf{w}}(\mathbf{x}), y_i) + \lambda r(\mathbf{w}) \Leftrightarrow \min_{\mathbf{w}, b} \sum_{i=1}^n \ell(h_{\mathbf{w}}(\mathbf{x}), y_i) \text{ subject to: } r(\mathbf{w}) \leq B$$

For each  $\lambda \geq 0$ , there exists  $B \geq 0$  such that the two formulations in (4.1) are equivalent, and vice versa. 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 Table 4.3.

Regularizer $r(\mathbf{w})$	Properties
<b><math>l_2</math>-Regularization</b> $r(\mathbf{w}) = \mathbf{w}^\top \mathbf{w} = \ \mathbf{w}\ _2^2$	<ul style="list-style-type: none"> <li>◦ ADVANTAGE: Strictly Convex</li> <li>◦ ADVANTAGE: 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>
<b><math>l_1</math>-Regularization</b> $r(\mathbf{w}) = \ \mathbf{w}\ _1$	<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>
<b><math>l_p</math>-Norm</b> $\ \mathbf{w}\ _p = (\sum_{i=1}^d v_i^p)^{1/p}$	<ul style="list-style-type: none"> <li>◦ (often <math>0 &lt; p \leq 1</math>)</li> <li>◦ DISADVANTAGE: Non-convex</li> <li>◦ ADVANTAGE: Very sparse solutions</li> <li>◦ Initialization dependent</li> <li>◦ DISADVANTAGE: Not differentiable</li> </ul>

Table 4.3: Types of Regularizers

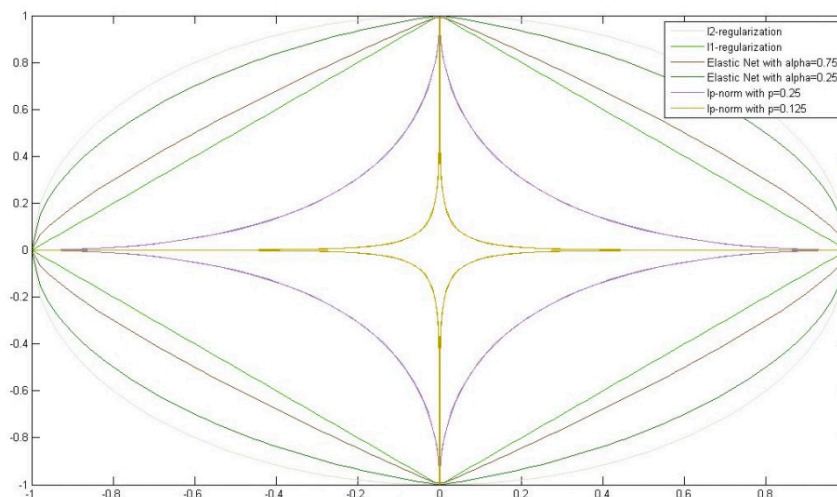


Figure 4.3: Plots of Common Regularizers



## 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. Table 4.4 provides information on their loss functions, regularizers, as well as solutions.

Loss and Regularizer	Comments
<b>Ordinary Least Squares</b> $\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$	<ul style="list-style-type: none"> <li>◦ Squared Loss</li> <li>◦ No Regularization</li> <li>◦ Closed form solution:</li> <li>◦ <math>\mathbf{w} = (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{X}\mathbf{y}^\top</math></li> <li>◦ <math>\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]</math></li> <li>◦ <math>\mathbf{y} = [y_1, \dots, y_n]</math></li> </ul>
<b>Ridge Regression</b> $\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \ \mathbf{w}\ _2^2$	<ul style="list-style-type: none"> <li>◦ Squared Loss</li> <li>◦ <math>l_2</math>-Regularization</li> <li>◦ <math>\mathbf{w} = (\mathbf{X}\mathbf{X}^\top + \lambda \mathbb{I})^{-1} \mathbf{X}\mathbf{y}^\top</math></li> </ul>
<b>Lasso</b> $\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \ \mathbf{w}\ _1$	<ul style="list-style-type: none"> <li>◦ + sparsity inducing (good for feature selection)</li> <li>◦ + Convex</li> <li>◦ - Not strictly convex (no unique solution)</li> <li>◦ - Not differentiable (at 0)</li> <li>◦ Solve with (sub)-gradient descent or <a href="#">SVEN</a></li> </ul>
<b>Elastic Net</b> $\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \alpha \ \mathbf{w}\ _1 + (1 - \alpha) \ \mathbf{w}\ _2^2$ $\alpha \in [0, 1]$	<ul style="list-style-type: none"> <li>◦ ADVANTAGE: Strictly convex (i.e. unique solution)</li> <li>◦ + sparsity inducing (good for feature selection)</li> <li>◦ + Dual of squared-loss SVM, see <a href="#">SVEN</a></li> <li>◦ DISADVANTAGE: - Non-differentiable</li> </ul>
<b>Logistic Regression</b> $\min_{\mathbf{w}, b} \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)})$	<ul style="list-style-type: none"> <li>◦ Often <math>l_1</math> or <math>l_2</math> Regularized</li> <li>◦ Solve with gradient descent.</li> <li>◦ <math>\Pr(y x) = \frac{1}{1 + e^{-y(\mathbf{w}^\top \mathbf{x} + b)}}</math></li> </ul>
<b>Linear Support Vector Machine</b> $\min_{\mathbf{w}, b} C \sum_{i=1}^n \max[1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b), 0] + \ \mathbf{w}\ _2^2$	<ul style="list-style-type: none"> <li>◦ Typically <math>l_2</math> regularized (sometimes <math>l_1</math>).</li> <li>◦ Quadratic program.</li> <li>◦ When <a href="#">kernelized</a> leads to <b>sparse</b> solutions.</li> <li>◦ Kernelized version can be solved very efficiently with specialized algorithms (e.g. <a href="#">SMO</a>)</li> </ul>

Table 4.4: Special Cases