# Descending into ML: Training and Loss

**Estimated Time:** 6 minutes

**Training** a model simply means learning (determining) good values for all the weights and the bias from labeled examples. In supervised learning, a machine learning algorithm builds a model by examining many examples and attempting to find a model that minimizes loss; this process is called **empirical risk minimization**.

Loss is the penalty for a bad prediction. That is, **loss** is a number indicating how bad the model's prediction was on a single example. If the model's prediction is perfect, the loss is zero; otherwise, the loss is greater. The goal of training a model is to find a set of weights and biases that have *low* loss, on average, across all examples. For example, Figure 3 shows a high loss model on the left and a low loss model on the right. Note the following about the figure:

- The arrows represent loss.
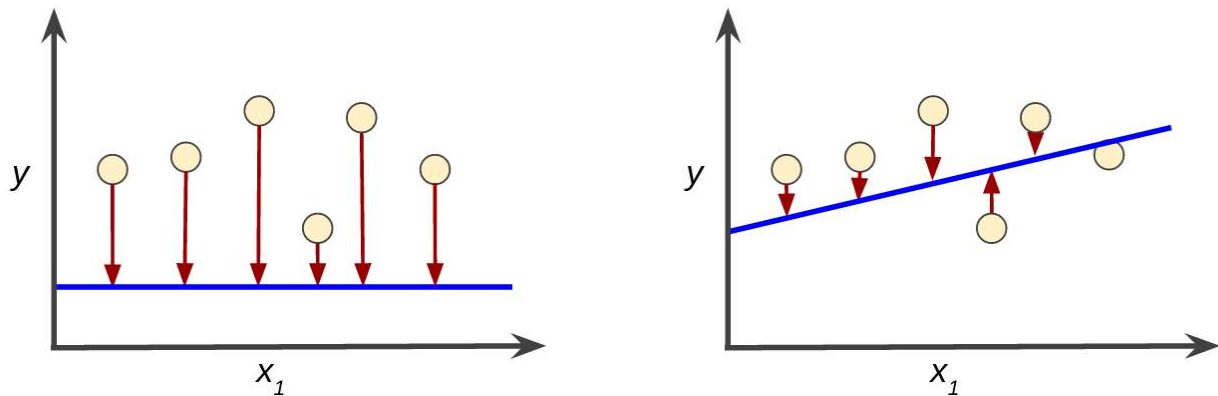
- The blue lines represent predictions.



**Figure 3. High loss in the left model; low loss in the right model.**

Notice that the arrows in the left plot are much longer than their counterparts in the right plot. Clearly, the line in the right plot is a much better predictive model than the line in the left plot.

You might be wondering whether you could create a mathematical function—a loss function —that would aggregate the individual losses in a meaningful fashion.

# Squared loss: a popular loss function

The linear regression models we'll examine here use a loss function called **squared loss** (also known as **L$_2$ loss**). The squared loss for a single example is as follows:

```
= the square of the difference between the label and the prediction
= (observation - prediction(x))²
= (y - y')²
```

**Mean square error** (**MSE**) is the average squared loss per example over the whole dataset. To calculate MSE, sum up all the squared losses for individual examples and then divide by the number of examples:

$$MSE = \frac{1}{N} \sum_{(x,y)\in D} (y - prediction(x))^2$$

where:

- $(x, y)$ is an example in which

  - $x$ is the set of features (for example, chirps/minute, age, gender) that the model uses to make predictions.

  - $y$ is the example's label (for example, temperature).

- $prediction(x)$ is a function of the weights and bias in combination with the set of features $x$.

- $D$ is a data set containing many labeled examples, which are $(x, y)$ pairs.

- $N$ is the number of examples in $D$.

Although MSE is commonly-used in machine learning, it is neither the only practical loss function nor the best loss function for all circumstances.

**Key Terms**

- empirical risk minimization (/machine-learning/glossary#ERM)

- mean squared error (/machine-learning/glossary#MSE)

- loss (/machine-learning/glossary#loss)

- squared loss (/machine-learning/glossary#squared_loss)