

# 1 Overview

In general, we aim to fit a line (or hyperplane in higher dimensions) to a scattering of data.

## 1.1 Notation and Modeling

Data for regression problems goes in the form of  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$  where  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}$ .

Each input in  $\mathbf{x}$  may be a column vector of length  $N$ .

Formally, the goal of regression is the following formula:

$$\operatorname{argmin}_{b, w} \sum_{p=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i - b)^2 \quad (1)$$

where  $\mathbf{w} \in \mathbb{R}^d$  is the weight vector and  $b \in \mathbb{R}$  is the bias term.

The gradient of this cost after some chain rule:

$$\nabla g(w) = 2 \left( \sum_{p=1}^N x_p x_p^T \right) w - 2 \sum_{p=1}^N x_p y_p \quad (2)$$

Setting the gradient above to zero and solving for  $w$  gives the system of linear equations

$$\left( \sum_{p=1}^N x_p x_p^T \right) w = \sum_{p=1}^N x_p y_p \quad (3)$$

$$w^* = \left( \sum_{p=1}^N x_p x_p^T \right)^{-1} \sum_{p=1}^N x_p y_p \quad (4)$$

## 1.2 Efficacy of the Model

The efficacy of the model can be measured by the mean squared error (MSE) of the model:

$$\frac{1}{N} \sum_{p=1}^N (y_p - w^T x_p)^2 \tag{5}$$

## References

- [1] Jeremy Wattós, Reza Borhanié Aggelos K. Katsaggelos, *Machine Learning Refined: Foundations, Algorithms, and Applications*, Northwestern University.