

Linear Regression

An ice cream shop recorded its sales on various days of the year to try to determine if there is a relationship between temperature and sales. The data points are plotted below, in blue. The red line is an example of a simple linear regression obtained by minimizing the squared errors from the points to the line.

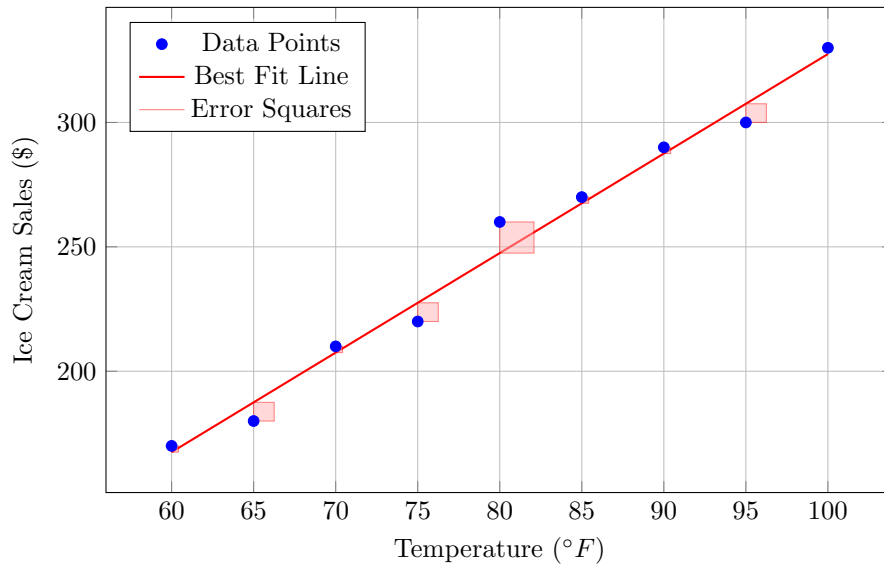


Figure 1: Temperature vs. Ice Cream Sales with Best Fit Line

1 Ordinary Least-Squares

Imagine we are given a data set $\mathcal{D} = \{(\mathbf{x}_i, y_i) \mid \mathbf{x}_i \in \mathbb{R}^d, y_i \in \mathbb{R}, i \in \{1, \dots, n\}\}$ of size n , where each data point \mathbf{x}_i is d -dimensional, to which we hope to fit a linear model. When $d = 1$, a **linear model** is a model in which each y_i is approximated by $mx_i + b$. Here, as in high school, m is the slope, while b is the y intercept.

In a regression model, Y is called the response—or dependent—variable, while the features are called the explanatory variables, or the predictors, and we are said to be “regressing Y on X .”

The feature values can be encoded in a matrix X :

$$X = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_n \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \dots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix}$$

Each \mathbf{x}_i in X is a row vector, i.e., an element of \mathbb{R}^d . When the number of features $d = 1$, the problem is called **simple regression**. Otherwise, it is called **multiple regression**.

¹This notes were compiled in conjunction with with Professors Eric Ewing and Amy Greenwald.

The response variables $\mathbf{y} \in \mathbb{R}^n$ are given by the column vector:

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

The goal is to find a set of weights $\mathbf{w} \in \mathbb{R}^d$, or parameters, or **coefficients** (e.g., m and b), so that the **estimates**,² or **predictions**,³ $X\mathbf{w}$ approximate \mathbf{y} :

$$X\mathbf{w} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \dots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} = \begin{bmatrix} w_1x_{11} + w_2x_{12} + \dots + w_dx_{1d} \\ w_2x_{21} + w_2x_{12} + \dots + w_dx_{1d} \\ \vdots \\ w_1x_{n1} + w_2x_{n2} + \dots + w_dx_{nd} \end{bmatrix} \approx \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \mathbf{y}$$

Observe that $X\mathbf{w} = \mathbf{y}$ is a system of n equations with d unknowns. A system of equations can be solved when it has the same number of equations and unknowns, i.e., when $n = d$. If $n < d$, i.e., if the matrix is wider than it is tall, then the system is **underdetermined**; it has infinitely many solutions. If $n > d$, i.e., if the matrix is taller than it is wide, which is usually the case in regression, the system is **overdetermined**, which means it has no solution. Consequently, rather than aim to solve this system of equations exactly, our goal is to solve for weights that minimize the difference between the predicted values $X\mathbf{w}$ and the observed values \mathbf{y} . As usual, we square the errors, so as to make the ensuing optimization problem smooth.

Define the **residual** of a candidate solution \mathbf{w} as the difference between \mathbf{y} and $X\mathbf{w}$, i.e.,

$$\text{residual}(\mathbf{w}) = \mathbf{y} - X\mathbf{w}$$

The **least squares objective** is to minimize the square of this residual value: $\text{loss}(\mathbf{w}) \doteq \text{residual}^2(\mathbf{w})$.

As $\text{residual}(\mathbf{w})$ is an n -dimensional column vector, we square it as follows:

$$\text{loss}(\mathbf{w}) = \text{residual}^2(\mathbf{w}) = (\mathbf{y} - X\mathbf{w})^T(\mathbf{y} - X\mathbf{w}) = \sum_{i=1}^n (\mathbf{y} - X\mathbf{w})_i^2$$

In other words, we are interested in minimizing the sum of the squared residuals.

As usual, to minimize an objective, we take its derivative and set it equal to zero.

First, let's simplify the loss:

$$\text{loss}(\mathbf{w}) \tag{1}$$

$$= (\mathbf{y} - X\mathbf{w})^T(\mathbf{y} - X\mathbf{w}) \tag{2}$$

$$= \mathbf{y}^T\mathbf{y} - \mathbf{y}^T X\mathbf{w} - (X\mathbf{w})^T\mathbf{y} + (X\mathbf{w})^T X\mathbf{w} \tag{3}$$

$$= \mathbf{y}^T\mathbf{y} - \mathbf{y}^T X\mathbf{w} - \mathbf{w}^T X^T\mathbf{y} + \mathbf{w}^T X^T X\mathbf{w} \tag{4}$$

$$= \mathbf{y}^T\mathbf{y} - 2\mathbf{w}^T X^T\mathbf{y} + \mathbf{w}^T X^T X\mathbf{w} \tag{5}$$

In Equation 3, we rely on the following fact: $(AB)^T = B^T A^T$. Equation 5 follows from the fact that $(\mathbf{y}^T X\mathbf{w})^T = \mathbf{w}^T X^T\mathbf{y}$, and that these values are both scalars, so they are equal to their transpose.

²statistical nomenclature

³machine learning nomenclature

Now let's take the derivative:

$$\nabla_{\mathbf{w}} \text{loss}(\mathbf{w}) \tag{6}$$

$$= \nabla_{\mathbf{w}} \mathbf{y}^T \mathbf{y} - 2 \nabla_{\mathbf{w}} \mathbf{w}^T X^T \mathbf{y} + \nabla_{\mathbf{w}} \mathbf{w}^T X^T X \mathbf{w} \tag{7}$$

$$= -2X^T \mathbf{y} + 2X^T X \mathbf{w} \tag{8}$$

Equation 8 relies on the following fact: $\nabla_{\mathbf{x}} \mathbf{x}^T A \mathbf{x} = 2A\mathbf{x}$.

Setting this derivative equal to zero yields a **closed-form** (i.e., analytical) solution to linear regression:

$$-2X^T \mathbf{y} + 2X^T X \mathbf{w} = 0 \tag{9}$$

$$X^T X \mathbf{w} = X^T \mathbf{y} \tag{10}$$

$$\mathbf{w} = \underbrace{(X^T X)^{-1} X^T \mathbf{y}}_{\text{pseudo-inverse}} \tag{11}$$

The matrix $X^T X$ might not be invertible. This matrix is only invertible when the columns of X are linearly independent, so that they span \mathbb{R}^d . The pseudo-inverse $(X^T X)^{-1} X^T$, however, always exists, and can be calculated via a **singular value decomposition**.

We have established that \mathbf{w} satisfies the **first-order optimality condition**, namely $\nabla_{\mathbf{w}} \text{loss}(\mathbf{w}) = 0$. To conclude that \mathbf{w} is a minimum, we also need to establish that it satisfies the second-order optimality condition. For \mathbf{w} to be a minimum, the **second-order optimality condition** requires that $\nabla_{\mathbf{w}}^2 f(\mathbf{w}) \geq 0$. Indeed, $\nabla_{\mathbf{w}}^2 \text{loss}(\mathbf{w}) = 2X^T X \geq 0$.

Alternative Derivation We describe an alternative shorter, geometric, and arguably more intuitive way, to derive the least-squares estimators. The only way the residual $\mathbf{y} - X\mathbf{w}$ can be zero is if \mathbf{y} is a linear combination of the columns of X : i.e., if \mathbf{y} lies in the column span of X . As most of the time it does not, our stated goal of minimizing the value of the residual (squared) can instead be understood as minimizing the distance from the residual to the column span of X . To minimize this distance, the residual should be **projected** onto the column span of X ; that is, it should be **orthogonal** to each column of X : for all $j \in \{1, \dots, n\}$, i.e., $(\mathbf{y} - X\mathbf{w}) \cdot X_j = 0$. Equivalently, but expressed more compactly in matrix notation, $(\mathbf{y} - X\mathbf{w})^T X = 0$. This requirement implies:

$$X^T (\mathbf{y} - X\mathbf{w}) = 0 \tag{12}$$

$$X^T X \mathbf{w} = X^T \mathbf{y} \tag{13}$$

$$\mathbf{w} = (X^T X)^{-1} X^T \mathbf{y} \tag{14}$$