

# Derivations for Linear Regression

Sayantana Auddy

This document contains derivations for linear regression with and without regularization. All the necessary linear algebra proofs have been explained in the Appendix. Refer to the links in the footnotes for more information about matrix calculus.

## 1 Linear Regression<sup>1</sup>

We have a *design matrix*  $\Phi \in \mathbb{R}^{n \times m}$  (e.g. constructed using  $n$  data points and for a  $m$ -dimensional polynomial). We want to find a weight vector  $\mathbf{w} \in \mathbb{R}^m$  such that the following is true where  $\mathbf{t}$  is the vector of labels.

$$\begin{aligned} & \Phi \mathbf{w} \approx \mathbf{t} \\ \Rightarrow & \begin{bmatrix} \Phi_{11} & \Phi_{12} & \cdots & \Phi_{1m} \\ \Phi_{21} & \Phi_{22} & \cdots & \Phi_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \Phi_{n1} & \Phi_{n2} & \cdots & \Phi_{nm} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \end{bmatrix} \approx \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{bmatrix} \\ \Rightarrow \Phi \mathbf{w} = & \begin{bmatrix} w_1 \Phi_{11} + w_2 \Phi_{12} + \cdots + w_m \Phi_{1m} \\ w_1 \Phi_{21} + w_2 \Phi_{22} + \cdots + w_m \Phi_{2m} \\ \vdots \\ w_1 \Phi_{n1} + w_2 \Phi_{n2} + \cdots + w_m \Phi_{nm} \end{bmatrix} \approx \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{bmatrix} \end{aligned}$$

Thus, we need to minimize the differences between our predicted label vector ( $\Phi \mathbf{w}$ ) and our target label vector  $\mathbf{t}$ . The vector containing the difference for each data point can be computed as

$$\Phi \mathbf{w} - \mathbf{t} = \begin{bmatrix} (\Phi w)_1 - t_1 \\ (\Phi w)_2 - t_2 \\ \vdots \\ (\Phi w)_n - t_n \end{bmatrix}$$

The squared error is a scalar number which gives the cumulative error over all the  $n$  data points. This can be obtained by adding the squared errors for each prediction. In other words,

$$\begin{aligned} E(\mathbf{w}) &= \frac{1}{2} \left( ((\Phi w)_1 - t_1)^2 + ((\Phi w)_2 - t_2)^2 + \cdots + ((\Phi w)_n - t_n)^2 \right) \left( \frac{1}{2} \text{ is used for mathematical convenience} \right) \\ &= \frac{1}{2} \left[ \begin{matrix} ((\Phi w)_1 - t_1) & ((\Phi w)_2 - t_2) & \cdots & ((\Phi w)_n - t_n) \end{matrix} \right] \begin{bmatrix} ((\Phi w)_1 - t_1) \\ ((\Phi w)_2 - t_2) \\ \vdots \\ ((\Phi w)_n - t_n) \end{bmatrix} \\ &= \frac{1}{2} (\Phi \mathbf{w} - \mathbf{t})^T (\Phi \mathbf{w} - \mathbf{t}) \end{aligned}$$

So, now we need to find  $\mathbf{w}$  such that  $E(\mathbf{w})$  is minimized. This is done by computing the gradient of  $E(\mathbf{w})$  with respect to  $\mathbf{w}$  and setting it to 0. Thus,

$$\mathbf{w} = \arg \min_{\mathbf{w}} E(\mathbf{w})$$

---

<sup>1</sup>This is based on [http://www.haija.org/derivation\\_lin\\_regression.pdf](http://www.haija.org/derivation_lin_regression.pdf), which contains a more condensed form of this proof.

$$\begin{aligned}
\Rightarrow \nabla_{\mathbf{w}} E(\mathbf{w}) &= \nabla_{\mathbf{w}} \left( \frac{1}{2} (\Phi \mathbf{w} - \mathbf{t})^T (\Phi \mathbf{w} - \mathbf{t}) \right) = 0 \\
&= \frac{1}{2} \nabla_{\mathbf{w}} \left( (\Phi \mathbf{w} - \mathbf{t})^T (\Phi \mathbf{w} - \mathbf{t}) \right) = 0 \\
&= \frac{1}{2} \nabla_{\mathbf{w}} \left( (\mathbf{w}^T \Phi^T - \mathbf{t}^T) (\Phi \mathbf{w} - \mathbf{t}) \right) = 0 \\
&\text{(since } (A - B)^T = A^T - B^T \text{ and } (AB)^T = B^T A^T) \\
&= \frac{1}{2} \nabla_{\mathbf{w}} \left( \mathbf{w}^T \Phi^T \Phi \mathbf{w} - \mathbf{t}^T \Phi \mathbf{w} - \mathbf{w}^T \Phi^T \mathbf{t} + \mathbf{t}^T \mathbf{t} \right) = 0 \\
&= \frac{1}{2} \nabla_{\mathbf{w}} \left( \mathbf{w}^T \Phi^T \Phi \mathbf{w} - \mathbf{t}^T \Phi \mathbf{w} - \mathbf{t}^T \Phi \mathbf{w} + \mathbf{t}^T \mathbf{t} \right) = 0 \\
&\text{(since } \mathbf{t}^T \Phi \mathbf{w} = \mathbf{w}^T \Phi^T \mathbf{t} = \text{same scalar number)} \\
&= \frac{1}{2} \nabla_{\mathbf{w}} \left( \mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{t}^T \Phi \mathbf{w} + \mathbf{t}^T \mathbf{t} \right) = 0 \\
&= \frac{1}{2} \nabla_{\mathbf{w}} \left( \mathbf{w}^T (\Phi^T \Phi) \mathbf{w} \right) - \nabla_{\mathbf{w}} \left( \mathbf{t}^T \Phi \mathbf{w} \right) + \frac{1}{2} \nabla_{\mathbf{w}} \left( \mathbf{t}^T \mathbf{t} \right) = 0 \\
&= \frac{1}{2} 2(\Phi^T \Phi) \mathbf{w} - \left( \mathbf{t}^T \Phi \right)^T + 0 = 0 \\
&\text{(for the first part note that } (\Phi^T \Phi) \text{ is a symmetric matrix (see section A.3), and using this fact along with the rule in section A.2 gives us } \nabla_{\mathbf{w}} \mathbf{w}^T (\Phi^T \Phi) \mathbf{w} = 2(\Phi^T \Phi) \mathbf{w}; \text{ for the second part see section A.1; the third part does not depend on } \mathbf{w} \text{ and so its gradient w.r.t. } \mathbf{w} \text{ is 0)} \\
&= (\Phi^T \Phi) \mathbf{w} - \Phi^T \mathbf{t} = 0 \\
&\Rightarrow (\Phi^T \Phi) \mathbf{w} = \Phi^T \mathbf{t} \\
&\Rightarrow \mathbf{w} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \text{ (multiplying both sides by } (\Phi^T \Phi)^{-1} \text{)}
\end{aligned}$$

## 2 Linear Regression with Regularization

In linear regression without regularization, our error function was defined as

$$E(\mathbf{w}) = \frac{1}{2} (\Phi \mathbf{w} - \mathbf{t})^T (\Phi \mathbf{w} - \mathbf{t})$$

In order to minimize overfitting, we want to regularize the linear regression by introducing a penalty on the weights (such that low values of weights are preferred). We do this by modifying the error function as given below

$$E(\mathbf{w}) = \frac{1}{2} (\Phi \mathbf{w} - \mathbf{t})^T (\Phi \mathbf{w} - \mathbf{t}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

(note that the sum of squared weights is given by  $\mathbf{w}^T \mathbf{w} = w_1^2 + w_1^2 + \dots + w_m^2$  and  $\lambda$  is the regularization coefficient)

Then, similar to section 1, we compute the gradient of the error and set it to 0.

$$\begin{aligned}
\mathbf{w} &= \arg \min_{\mathbf{w}} E(\mathbf{w}) \\
\Rightarrow \nabla_{\mathbf{w}} E(\mathbf{w}) &= \nabla_{\mathbf{w}} \left( \frac{1}{2} (\Phi \mathbf{w} - \mathbf{t})^T (\Phi \mathbf{w} - \mathbf{t}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \right) = 0 \\
&= \nabla_{\mathbf{w}} \left( \frac{1}{2} (\Phi \mathbf{w} - \mathbf{t})^T (\Phi \mathbf{w} - \mathbf{t}) \right) + \nabla_{\mathbf{w}} \left( \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \right) = 0 \\
&= \left( (\Phi^T \Phi) \mathbf{w} - \Phi^T \mathbf{t} \right) + \frac{\lambda}{2} \nabla_{\mathbf{w}} \left( \mathbf{w}^T \mathbf{w} \right) = 0 \\
&\text{(the first part comes from the derivation in section 1)} \\
&= \left( (\Phi^T \Phi) \mathbf{w} - \Phi^T \mathbf{t} \right) + \lambda \mathbf{w} = 0 \\
&\text{(since } \nabla_{\mathbf{w}} \mathbf{w}^T \mathbf{w} = 2\mathbf{w} \text{); see section A.4)}
\end{aligned}$$

$$\begin{aligned} &\Rightarrow ((\Phi^T \Phi) + \lambda \mathbf{I}) \mathbf{w} = \Phi^T \mathbf{t} \\ &\Rightarrow \mathbf{w} = ((\Phi^T \Phi) + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{t} \end{aligned}$$

## A Appendix: Useful Matrix Calculus rules<sup>2</sup>

### A.1 Rule 1

$$\begin{aligned} \nabla_{\mathbf{x}} \mathbf{b}^T \mathbf{x} &= \mathbf{b} \\ \text{where } \mathbf{x} \text{ and } \mathbf{b} \text{ are vectors of size } n \text{ (} \mathbf{x} \in \mathbb{R}^n \text{ and } \mathbf{b} \in \mathbb{R}^n \text{)} \end{aligned} \tag{1}$$

*Proof.*

$$\begin{aligned} \text{Let } f(\mathbf{x}) &= \mathbf{b}^T \mathbf{x} = \sum_{i=1}^n b_i x_i \\ \text{So } \frac{\partial f(\mathbf{x})}{\partial x_k} &= \frac{\partial}{\partial x_k} \sum_{i=1}^n b_i x_i = b_k \\ \text{Therefore } \nabla_{\mathbf{x}} \mathbf{b}^T \mathbf{x} &= \begin{bmatrix} \frac{\partial}{\partial x_1} \sum_{i=1}^n b_i x_i = b_1 \\ \frac{\partial}{\partial x_2} \sum_{i=1}^n b_i x_i = b_2 \\ \vdots \\ \frac{\partial}{\partial x_n} \sum_{i=1}^n b_i x_i = b_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} = \mathbf{b} \end{aligned}$$

■

### A.2 Rule 2

$$\begin{aligned} \nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} &= 2\mathbf{A} \mathbf{x} \\ \text{where } \mathbf{x} \text{ is a vector of size } n \text{ (} \mathbf{x} \in \mathbb{R}^n \text{)} \text{ and } \mathbf{A} \text{ is a symmetric matrix of size } n \times n \text{ (} \mathbf{A} \in \mathbb{R}^{n \times n} \text{)} \end{aligned} \tag{2}$$

*Proof.*

$$\begin{aligned} \text{Let } f(\mathbf{x}) &= \mathbf{x}^T \mathbf{A} \mathbf{x} \\ &= [x_1 \quad x_2 \quad \cdots \quad x_n] \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1n} \\ A_{21} & A_{22} & \cdots & A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{n1} & A_{n2} & \cdots & A_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \\ &= [(x_1 A_{11} + x_2 A_{21} + \cdots + x_n A_{n1}) \quad \cdots \quad (x_1 A_{1n} + x_2 A_{2n} + \cdots + x_n A_{nn})] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \\ &= x_1(x_1 A_{11} + x_2 A_{21} + \cdots + x_n A_{n1}) + \cdots + x_n(x_1 A_{1n} + x_2 A_{2n} + \cdots + x_n A_{nn}) \\ &= x_1 \left( \sum_{i=1}^n x_i A_{i1} \right) + \cdots + x_n \left( \sum_{i=1}^n x_i A_{in} \right) \end{aligned}$$

<sup>2</sup>Some of the proofs are based on <http://cs229.stanford.edu/section/cs229-linalg.pdf> but have been simplified further by showing the matrix elements where ever possible. Useful sources for matrix calculus are <https://atmos.washington.edu/~dennis/MatrixCalculus.pdf> and of course [https://en.wikipedia.org/wiki/Matrix\\_calculus](https://en.wikipedia.org/wiki/Matrix_calculus).

$$\begin{aligned}
&= \sum_{j=1}^n \sum_{i=1}^n x_j x_i A_{ij} \\
&= \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j \\
\frac{\partial f(\mathbf{x})}{\partial x_k} &= \frac{\partial}{\partial x_k} \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j \\
&= \frac{\partial}{\partial x_k} \left[ \sum_{i \neq k} \sum_{j \neq k} A_{ij} x_i x_j + \sum_{i \neq k} A_{ik} x_i x_k + \sum_{j \neq k} A_{kj} x_k x_j + A_{kk} x_k x_k \right]
\end{aligned}$$

Separating  $f(\mathbf{x})$  into 4 parts:  $(i \neq k, j \neq k), (i \neq k, j = k), (i = k, j \neq k), (i = k, j = k)$

$$\begin{aligned}
&= \sum_{i \neq k} A_{ik} x_i + \sum_{j \neq k} A_{kj} x_j + 2A_{kk} x_k \\
&= \sum_{i \neq k} A_{ik} x_i + A_{kk} x_k + \sum_{j \neq k} A_{kj} x_j + A_{kk} x_k \\
&= \sum_{i=1}^n A_{ik} x_i + \sum_{j=1}^n A_{kj} x_j \\
&= 2 \sum_{i=1}^n A_{ik} x_i \quad (\text{since } A \text{ is symmetric, } A_{ik} = A_{kj})
\end{aligned}$$

$$\begin{aligned}
\nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} &= \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_k} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix} = \begin{bmatrix} 2 \sum_{i=1}^n A_{i1} x_i \\ 2 \sum_{i=1}^n A_{i2} x_i \\ \vdots \\ 2 \sum_{i=1}^n A_{ik} x_i \\ \vdots \\ 2 \sum_{i=1}^n A_{in} x_i \end{bmatrix} = 2 \begin{bmatrix} A_{11} x_1 + A_{21} x_2 + \cdots + A_{n1} x_n \\ A_{12} x_1 + A_{22} x_2 + \cdots + A_{n2} x_n \\ \vdots \\ A_{1k} x_1 + A_{2k} x_2 + \cdots + A_{nk} x_n \\ \vdots \\ A_{1n} x_1 + A_{2n} x_2 + \cdots + A_{nn} x_n \end{bmatrix} \\
&= 2 \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1n} \\ A_{21} & A_{22} & \cdots & A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{n1} & A_{n2} & \cdots & A_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \\
&= 2 \begin{bmatrix} A_{11} & A_{21} & \cdots & A_{n1} \\ A_{12} & A_{22} & \cdots & A_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ A_{1n} & A_{2n} & \cdots & A_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad (\text{since } A \text{ is symmetric, } A^T = A) \\
&= 2\mathbf{A}\mathbf{x}
\end{aligned}$$

■

### A.3 Rule 3

For any matrix  $\Phi \in \mathbb{R}^{m \times n}$ , the matrix  $\Phi^T \Phi$  is a symmetric matrix.

*Proof.* For any matrix  $\mathbf{A} \in \mathbb{R}^{n \times m}$  and  $\mathbf{B} \in \mathbb{R}^{m \times n}$ ,  $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$

$$\begin{aligned}(\Phi^T \Phi)^T &= \Phi^T (\Phi^T)^T \text{ (taking } \mathbf{A} = \Phi^T \text{ and } \mathbf{B} = \Phi) \\ &= \Phi^T \Phi \text{ (since } (\Phi^T)^T = \Phi)\end{aligned}$$

Since the matrix  $(\Phi^T \Phi)$  is equal to its transpose, it is symmetric. ■

### A.4 Rule 4

$$\nabla_{\mathbf{w}} (\mathbf{w}^T \mathbf{w}) = 2\mathbf{w} \text{ (where } \mathbf{w} \in \mathbb{R}^m)$$

*Proof.*

$$\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_k \\ \vdots \\ w_m \end{bmatrix}$$

$$\mathbf{w}^T \mathbf{w} = [w_1 \quad w_2 \quad \cdots \quad w_k \quad \cdots \quad w_m] \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_k \\ \vdots \\ w_m \end{bmatrix} = (w_1^2 + w_2^2 + \cdots + w_k^2 + \cdots + w_m^2)$$

$$\frac{\partial}{\partial w_k} \mathbf{w}^T \mathbf{w} = 2w_k$$

$$\nabla_{\mathbf{w}} (\mathbf{w}^T \mathbf{w}) = \begin{bmatrix} 2w_1 \\ 2w_2 \\ \vdots \\ 2w_k \\ \vdots \\ 2w_m \end{bmatrix} = 2\mathbf{w}$$
■