

Mathematics For AI

Madiba Hudson-Quansah

CONTENTS

CHAPTER 1	FOUNDATIONS	PAGE 3
1.1	Vector Notation	3
1.2	Numerical vs Analytical Solutions	4
1.3	Multivariable Calculus	4
	Common derivatives of linear algebra — 5	
1.4	Gradient Descent	5
1.5	Regression	6
1.6	Classification	6
CHAPTER 2	LINEAR MODELS	PAGE 7
2.1	Linear Regression	7
	Simple Linear Regression — 7	
2.1.1.1	Training Function	7
2.1.1.2	Loss Function	7
2.1.1.3	Optimization	7
	Multi-Linear Regression — 7	
2.1.2.1	Training Function	7
2.1.2.1.1	Parametric Models vs Non-Parametric Models	8
2.1.2.2	Loss Function	8
2.1.2.2.1	Functions with Singularities	9
2.1.2.2.2	Mean Squared Error	9
2.1.2.3	Optimization	10
2.1.2.3.1	Convex Landscapes vs Non convex Landscapes	10
2.1.2.3.2	Locating Function Minimizers	11
2.1.2.3.3	Minimizing the mean squared error function	11
2.1.2.3.4	Using the derivate approach	11
	Polynomial Regression — 11	
CHAPTER 3	DECISION TREES	PAGE 12
3.1	Decision Tree Regressors	12
CHAPTER 4	REGULARIZATION	PAGE 13
4.1	Regularization of Regression Models	13

Chapter 1

Foundations

The main structure of the machine learning part of AI is as follows:

1. Identify the problem
2. Acquire the appropriate data
3. Create a hypothesis/ learning/ prediction / training function
4. Find the numerical values of weights
5. Create an error function
6. Decide on mathematical formulas
7. Find a way to search for minimizers
8. Use the backpropagation algorithm
9. Regularize a function

1.1 Vector Notation

Definition 1.1.1: Vector

A vector is a mathematical object that has both magnitude and direction. It is represented by an arrow with a starting point and an end point. The length of the arrow represents the magnitude of the vector, and the direction of the arrow represents the direction of the vector. Denoted by \vec{v} or \mathbf{v} which in vector space \mathbb{R}^n represents:

$$\mathbf{v} = \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix}$$

The transpose of a column vector \mathbf{v} is denoted by \mathbf{v}^t and is a row vector:

$$\mathbf{v}^t = (v_1 \quad \dots \quad v_n)$$

Definition 1.1.2: Dot Product

The dot product of two vectors \mathbf{v} and \mathbf{w} is a scalar value denoted by $\mathbf{v} \cdot \mathbf{w}$ and is defined as:

$$\mathbf{v} \cdot \mathbf{w} = \sum_{i=1}^n v_i w_i$$

The dot product is the same as the product of a row vector / transposed column vector and a column vector, for

example for vectors $\mathbf{a} = \begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \\ a_5 \end{pmatrix}$ and $\mathbf{b} = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \\ b_5 \end{pmatrix}$, the dot product is:

$$\begin{aligned} \mathbf{a} \cdot \mathbf{b} &= \begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \\ a_5 \end{pmatrix} \begin{pmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \\ b_5 \end{pmatrix} \\ &= a_1 b_1 + a_2 b_2 + a_3 b_3 + a_4 b_4 + a_5 b_5 \end{aligned}$$

Or

$$\begin{aligned} \mathbf{a}^t \mathbf{b} &= (a_1 \quad a_2 \quad a_3 \quad a_4 \quad a_5) \begin{pmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \\ b_5 \end{pmatrix} \\ &= a_1 b_1 + a_2 b_2 + a_3 b_3 + a_4 b_4 + a_5 b_5 \end{aligned}$$

Moreover:

$$\|\mathbf{a}\|_2^2 = \mathbf{a}^t \mathbf{a} = a_1^2 + a_2^2 + a_3^2 + a_4^2 + a_5^2$$

1.2 Numerical vs Analytical Solutions

Definition 1.2.1: Numerical

Has to do with numbers. Are approximations of analytical solutions obtained by discretizing a continuous problem.

Definition 1.2.2: Analytical

Has to do with analysis. Are exact solutions to problems. Not all problems have analytical solutions.

1.3 Multivariable Calculus

In multivariable calculus, the gradient of a function is a vector that points in the direction of the greatest rate of increase of the function at a given point. The gradient is denoted by ∇ and is defined as:

$$\nabla F = \begin{pmatrix} \frac{\partial F}{\partial x_1} \\ \frac{\partial F}{\partial x_2} \\ \vdots \\ \frac{\partial F}{\partial x_n} \end{pmatrix}$$

1.3.1 Common derivatives of linear algebra

- When a and ω and a is a constant, the derivate of $f(\omega) = a\omega$ is

$$f'(\omega) = a$$

Therefore when \mathbf{a} and $\boldsymbol{\omega}$ and the entries of \mathbf{a} are constant, then the derivate / gradient of $f(\boldsymbol{\omega}) = \mathbf{a}^t \boldsymbol{\omega}$ is:

$$\nabla f(\boldsymbol{\omega}) = \mathbf{a}$$

And similarly the derivate of $f(\boldsymbol{\omega}) = \boldsymbol{\omega}^t \mathbf{a}$ is

$$\nabla f(\boldsymbol{\omega}) = \mathbf{a}$$

- When s is constant and ω is scalars, then the derivate of the quadratic function $f(\omega) = s\omega^2$ is :

$$f'(\omega) = 2s\omega$$

Therefore when S is a symmetric matrix with constant entries, then the derivative of the function $f(\boldsymbol{\omega}) = \boldsymbol{\omega}^t S \boldsymbol{\omega}$ is:

$$\nabla f(\boldsymbol{\omega}) = 2S\boldsymbol{\omega}$$

Because $\boldsymbol{\omega}^t \boldsymbol{\omega}$, if $\boldsymbol{\omega} = \begin{bmatrix} \omega_0 \\ \omega_1 \\ \omega_2 \end{bmatrix}$, then $\boldsymbol{\omega}^t = [\omega_0 \quad \omega_1 \quad \omega_2]$ is defined as:

$$\boldsymbol{\omega}^t \boldsymbol{\omega} = [\omega_0 \quad \omega_1 \quad \omega_2] \begin{bmatrix} \omega_0 \\ \omega_1 \\ \omega_2 \end{bmatrix}$$

1.4 Gradient Descent

Definition 1.4.1: Gradient

The gradient of a function is a vector that points in the direction of the greatest rate of increase of the function at a given point.

Definition 1.4.2: Gradient Descent

An optimization algorithm used to minimize some function by iteratively moving in the direction of steepest descent as defined by the negative of the gradient.

This algorithm is used in the optimization of training functions using a loss function. Gradient descent is carried out as follows:

1. Choose a random starting point within the search space
2. Calculate the gradient of the function at that point
3. Move in the direction of the negative of the gradient by subtracting the gradient from the current point
4. Repeat steps 2 and 3 until the gradient is close to zero

As denoted by the following formula:

$$\boldsymbol{\omega}_{n+1} = \boldsymbol{\omega}_n - \alpha \nabla F(\boldsymbol{\omega}_n)$$

Where $\boldsymbol{\omega}_{n+1}$ is the new point, $\boldsymbol{\omega}_n$ is the current point, α is the learning rate, and $\nabla F(\boldsymbol{\omega}_n)$ is the gradient of the function at the current point.

Definition 1.4.3: Learning Rate

A hyperparameter that controls how much the weights are updated during training. A large learning rate means the weights are updated by a large amount, while a small learning rate means the weights are updated by a small amount.

1.5 Regression

Definition 1.5.1: Regression

A statistical method used to determine the strength and character of the relationship between one dependent variable and one or more independent variables.

1.6 Classification

Definition 1.6.1: Classification

A supervised learning approach that categorizes input data into one of a number of classes.

Chapter 2

Linear Models

2.1 Linear Regression

A linear regression model is a linear approach to modelling the relationship between a dependent variable and one or more independent variables. In using a linear model we make the assumption that the dependent variable / predicted value depends linearly on the independent variables / features.

2.1.1 Simple Linear Regression

In the case of simple linear regression, there is only one independent variable. This relationship can be expressed as:

$$y = \alpha + \beta x$$

Where y is the predicted value / dependent variable, α is the intercept, β is the slope, and x is the feature / independent variable.

2.1.1.1 Training Function

2.1.1.2 Loss Function

2.1.1.3 Optimization

2.1.2 Multi-Linear Regression

In the case of multi-linear regression, there are multiple independent variables. This relationship can be expressed as:

$$y = f(x_1, x_2, x_3, \dots, x_n)$$

Where y is the predicted value, $x_1, x_2, x_3, \dots, x_n$ are the features, and f is the function that maps the features to the predicted value.

2.1.2.1 Training Function

As with all linear models the assumption of linear dependence is made. This means that the predicted value is a linear combination of its features plus a bias term ω_0 , giving the following training function:

$$y = \omega_0 + \omega_1 x_1 + \omega_2 x_2 + \dots + \omega_n x_n$$

Where y is the predicted value, ω_0 is the bias term, $\omega_1, \omega_2, \dots, \omega_n$ are the weights, and x_1, x_2, \dots, x_n are the features. In order to find the appropriate weights ω_n , they have to be learnt from the training data. This process is called training the model. Ergo a trained model is model that has decided on the appropriate weights to use.

Each weight ω is a scalar value multiple by a corresponding feature x , and represent the importance of a specific feature in determining the predicted value. This means the larger the weight the more important the feature is in determining

the predicted value, and vice versa. *Dead features* are features that have a weight of zero and do not affect the predicted value. This means for each row in the total amount of rows m the prediction of this model is:

$$\begin{aligned} y_{\text{predict}}^1 &= \omega_0 + \omega_1 x_1^1 + \omega_2 x_2^1 + \dots + \omega_n x_n^1 \\ y_{\text{predict}}^2 &= \omega_0 + \omega_1 x_1^2 + \omega_2 x_2^2 + \dots + \omega_n x_n^2 \\ &\vdots \\ y_{\text{predict}}^m &= \omega_0 + \omega_1 x_1^m + \omega_2 x_2^m + \dots + \omega_n x_n^m \end{aligned}$$

Which can be expressed as:

$$\begin{pmatrix} y_{\text{predict}}^1 \\ y_{\text{predict}}^2 \\ \vdots \\ y_{\text{predict}}^m \end{pmatrix} = \omega_0 \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} + \omega_1 \begin{pmatrix} x_1^1 \\ x_1^2 \\ \vdots \\ x_1^m \end{pmatrix} + \dots + \omega_n \begin{pmatrix} x_n^1 \\ x_n^2 \\ \vdots \\ x_n^m \end{pmatrix}$$

$$\begin{pmatrix} y_{\text{predict}}^1 \\ y_{\text{predict}}^2 \\ \vdots \\ y_{\text{predict}}^m \end{pmatrix} = \begin{pmatrix} 1 & x_1^1 & x_2^1 & \dots & x_n^1 \\ 1 & x_1^2 & x_2^2 & \dots & x_n^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^m & x_2^m & \dots & x_n^m \end{pmatrix} \begin{pmatrix} \omega_0 \\ \omega_1 \\ \omega_2 \\ \vdots \\ \omega_n \end{pmatrix}$$

This can be then expressed as

$$\mathbf{y}_{\text{predict}} = \mathbf{X} \boldsymbol{\omega}$$

Where $\mathbf{y}_{\text{predict}}$ is the vector of predictions \mathbf{X} is the training subset, and $\boldsymbol{\omega}$ is the vector of unknown weights.

2.1.2.1.1 Parametric Models vs Non-Parametric Models

Definition 2.1.1: Parametric Model

A model that has parameters that are fixed in number, regardless of the size of the training data. For example:

$$y = \omega_0 + \omega_1 x_1 + \omega_2 x_2 + \dots + \omega_n x_n$$

The number of parameters (ω) is fixed at $n + 1$. This means that the formula of training function is fixed ahead of training the model, and the model is trained to find the appropriate values of the parameters.

Definition 2.1.2: Non-Parametric Model

A model that does not specify the formula for the training function with it's parameters before training the model. This results in the number of parameters being dependent on the size of the training data, as such models adapt to the data and determines the required number of parameters during training.

Both parametric and non-parametric models have other parameters that are not weights called hyperparameters, that need to be tuned during training.

Definition 2.1.3: Hyperparameters

Parameters that are not weights, and are not learnt during training. They are set before training the model, and determine the behaviour of the model.

2.1.2.2 Loss Function

To determine the weights of the model using the training data we need to define and optimize a loss function.

Definition 2.1.4: Loss Function

Determines the error of the model, by measuring the difference between the predicted values generated by the model and the actual values in the training data.

For example assuming we trained a model with n features, giving us the weights $\omega_0, \omega_1, \dots, \omega_n$, then the i th predicted values using the i th row of the training data would be:

$$y_{\text{predict}}^i = \omega_0 + \omega_1 x_1^i + \omega_2 x_2^i + \dots + \omega_n x_n^i$$

However for the i th row of the training data the actual value is y_{true}^i . In determining the error between the predicted value and the actual value we can use the following approaches.

Absolute Value Distance - $|y_{\text{predict}} - y_{\text{true}}|$ derived from $|x|$

Squared Distance - $|y_{\text{predict}} - y_{\text{true}}|^2$ derived from $|x|^2$

In examining both functions $|x|$ and $|x|^2$ we notice that the squared distance function is smoother than the absolute value distance function. This means that the squared distance function is differentiable at all points, while the absolute value distance function is not differentiable at $x = 0$, this is called a *singularity*.

Other than the difference in *regularity* of both functions another consideration must be made, *if a number is large, then its square is even larger*. This means that the squared distance function is more sensitive to outliers in the data than the absolute value distance function.

Definition 2.1.5: Regularity

The property of a function that determines how smooth it is. A function is regular if it is differentiable at all points.

2.1.2.2.1 Functions with Singularities

Definition 2.1.6: Singularity

A point at which a function is not differentiable.

Generally graphs of differential functions do not have sharp corners, cusps, or vertical tangents. This is because at that point you can draw two different tangents lines to the graph of the function at that point giving two different slopes at that point. This discontinuity in the slope of the tangent creates a problem for methods that rely on evaluating the derivative of the function such as the gradient descent algorithm.

2.1.2.2.2 Mean Squared Error

For linear regression models the most common loss function is the mean squared error function.

Definition 2.1.7: Mean Squared Error

The average of the squared differences between the predicted values and the actual values in the training data. Defined as:

$$\text{Mean Squared Error} = \frac{1}{m} \sum_{i=1}^m |y_{\text{predict}}^i - y_{\text{true}}^i|^2$$

Where m is the number of rows. And in linear algebra notation:

$$\text{Mean Squared Error} = \frac{1}{m} (\vec{y}_{\text{predict}} - \vec{y}_{\text{true}})^t (\vec{y}_{\text{predict}} - \vec{y}_{\text{true}}) = \frac{1}{m} \|\vec{y}_{\text{predict}} - \vec{y}_{\text{true}}\|_{l^2}^2$$

Where l^2 denotes the l^2 norm of the vector which by definition is the $\sqrt{\text{sum of squares of its components}}$.
And using $\mathbf{y}_{\text{predict}} = X\boldsymbol{\omega}$

$$\text{Mean Squared Error} = \frac{1}{m} (X\boldsymbol{\omega} - \mathbf{y}_{\text{true}})^t (X\boldsymbol{\omega} - \mathbf{y}_{\text{true}}) = \frac{1}{m} \|X\boldsymbol{\omega} - \mathbf{y}_{\text{true}}\|_{l^2}^2$$

2.1.2.3 Optimization**Definition 2.1.8: Optimization**

Finding the best / optimal / maximal / minimal / extreme solution.

Given our linear training function:

$$y = \omega_0 + \omega_1 x_1 + \omega_2 x_2 + \dots + \omega_n x_n$$

The goal of the optimization step is to find the values of $\omega_0, \omega_1, \omega_2, \dots, \omega_n$, that make our training function best fit the training data, where an optimal fit is characterized by the loss function. Therefore we need the weights that minimize the value of our loss function, denoted by:

$$\min_{\boldsymbol{\omega}} \text{loss function}$$

For our loss function, the mean squared error is:

$$\min_{\boldsymbol{\omega}} \frac{1}{m} \|X\boldsymbol{\omega} - \mathbf{y}_{\text{true}}\|_{l^2}^2$$

2.1.2.3.1 Convex Landscapes vs Non convex Landscapes**Definition 2.1.9: Convex Function**

Let f be a function that is differentiable over an interval I . The function f is convex if the first derivative of the function f' is increasing over I , which means the second derivative of the function, f'' is greater than zero over the interval I , $f'' > 0$

Definition 2.1.10: Non Convex / Concave Function

Let f be a function that is differentiable over an interval I . The function f is concave if f' is decreasing over I , which means f'' is less than zero over the interval I , $f'' < 0$

Most functions and equations we will deal with are non-linear, but can be linearised at certain points of interest. This is done by drawing a tangent line to the point in 2D space or a tangent hyperplane in greater dimensions. To do this the derivative of the function must be found with respect to all of its variables to determine the slope of the approximating flat space.

One important type of function is a function that is the maximum of two or more convex functions, which are always convex. Linear functions are flat so they are both convex and concave. This is important because some functions are

defined as the maxima of linear functions, which are not guaranteed to be linear, but are guaranteed to be convex. That is even if we lose linearity when taking the maximum of linear functions we still have convexity.

The Rectified Linear Unit function (ReLU), is an example of a non-linear function defined as the maximum of two linear functions:

$$\text{ReLU}(x) = \max(0, x)$$

There is one more relationship between linearity and convexity. If we have a non-linear convex function, the maximum of all the linear functions that stay below the function is exactly equal to it.

Overall the landscape of a convex function is good for minimization problems as there is no chance at getting stuck at a local minima as any local minimum is also a global minimum for a convex function. However a non convex function could have peaks, valleys and saddle points, which runs the risk of getting stuck at a local minima, and never finding the global minima.

2.1.2.3.2 Locating Function Minimizers

There are generally two approaches to locating minimizers/maximizers of functions with the trade-off usually being between:

- Calculating only one derivative and converging to the minimum slowly. These are called *gradient* methods, where the gradient is one derivate of a function of several variables as is our case with $\omega_0, \dots, \omega_n$.
- Calculating two derivatives and converging to the minimum faster. These are called *Newton's* methods, and the *Hessian* or an approximation of the Hessian is used in these methods

Definition 2.1.11: Hessian Matrix

A square matrix consisting of second derivatives of a function

One key idea from calculus remains fundamental: minimizers / maximizers occur at critical points. So to locate these minimizers / maximizers we must search through both the boundary points (points where the function is undefined) and interior critical points (where the first derivative of the function is equal to 0). This can be done using one of two approaches

Approach 1

1. Find the derivate of the function
2. Equate the derivative to 0
3. Solve for the ω 's that make the derivative 0.

This approach works when the derivate is linear and the ω s can easily be found but is not ideal when the derivate is non linear.

Approach 2 Follow the gradient direction to descend towards the minimum or ascend towards the maximum, using gradient descent 1.4.

2.1.2.3.3 Minimizing the mean squared error function

Let $L(\omega)$ be

$$\frac{1}{m} (X\omega - \mathbf{y}_{\text{true}})^t (X\omega - \mathbf{y}_{\text{true}})$$

2.1.2.3.4 Using the derivate approach

2.1.3 Polynomial Regression

Chapter 3

Decision Trees

3.1 Decision Tree Regressors

Chapter 4

Regularization

4.1 Regularization of Regression Models

Chapter 5

Evaluation

5.1 Evaluating Regression Models