

Mathematics behind Supervised Machine Learning Algorithms

Narsimha Chilkuri

Contents

1 K-Nearest Neighbours	2
2 Linear Regression	3
2.1 Statistical Inference for Linear Regression	5
2.2 Modifications to Linear Regression	5
2.3 Regularization for Linear Regression	6
3 Logistic Regression	7
4 Linear Discriminant Analysis	7
4.1 Quadratic Discriminant Analysis	8
5 Back-propagation for binary Neural Networks	8

Introduction

As of now, there are three major paradigms in machine learning. They are supervised learning, unsupervised learning and reinforcement learning. I find it hard to say anything here that applies to all three types of learning, other than the fact that none of them make use of explicit programming to solve problems but instead rely on data, whether instructive (in-case of supervised), evaluative (in-case of reinforcement) or @@@ (in-case of unsupervised), to solve problems of various kinds. Let us discuss each of these briefly down below.

Supervised Learning: Consider the following two conditions:

1. We have access to a dataset, D , consisting of pairs (x, y) :

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\},$$

where $x \in X \subset \mathbb{R}^d$ and $y \in Y \subset \mathbb{R}^l$.

2. We have reasons to believe that there exists a pattern to this data i.e., there is function f such that

$$\begin{aligned} Y &= f(X) + \text{noise}, \\ &= f(X) + \epsilon, \text{ where we assume } \epsilon \sim N(0, \sigma^2). \end{aligned}$$

Given the above two conditions, the goal of supervised learning is to obtain an *estimate* of this function f using nothing but the dataset D !

Unsupervised Learning:

Reinforcement Learning: The main idea of reinforcement learning is to learn by interacting with the environment to achieve a goal. More specifically, we deal with a Markov Decision Process (MDP) where the transition probabilities and rewards are unknown before-hand and the goal is to learn a policy.

In supervised machine learning, a distinction is usually made between the types of problems, dividing them into two classes known as classification and regression. I believe that, although this distinction is useful and sensible, it is more fundamental to draw a line between two classes of solution methods, ones that learn features from raw data by themselves and ones that do not. If I may be a bit brash, I would say that the ones that do not learn features are fairly boring and are reminiscent of basic data fitting problems that one encounters in something like experimental physics. The solution methods that learn features are fascinating and I believe they do more justice to the word “learning”.

In supervised machine learning, problems fall into one of two classes, classification and regression. The difference between them being the nature of the output we are trying to predict. If the output is discrete, then its a classification problem; if the output is continuous, then its a regression problem. We start by looking at regression.

1 K-Nearest Neighbours

If I were thinking about the data-driven approach to solving problems all by myself, I believe I would have discovered the K-Nearest Neighbours (KNN) approach first. This approach is not only intuitive, but also uses elementary mathematics. The crux of this approach can be stated as follows: a point has similar properties as the points lying in its proximity. Mathematically, we can write the previous statement as:

$$f(\vec{x}_*) = \frac{1}{K} \sum_{i=1}^K f(\vec{x}_i),$$

i.e., we predict the output of a point \vec{x}_* by averaging the outputs of K nearest points.

I have written down eq (1), but I have made no mention of what K should be. Is there a formula relating K to the number of training examples, dimension of

the problem, distribution of the training and testing data etc.? Unfortunately, as is presently the case with most of the algorithms in machine learning, there is no standard answer for values of parameters such as K that work well in all scenarios. It has to be figured out by trial and error.

Another thing that I have left out is the meaning of “near” (in nearest K neighbours) or the choice of the distance function. We use a generalized form of distance called the Minkowski metric defined as:

$$Distance(\vec{x}, \vec{y}; p) = \left(\sum_i |x_i - y_i|^p \right)^{\frac{1}{p}},$$

where x_i and y_i are components of the vectors \vec{x} and \vec{y} . In Fig. ?? we see all the points that are equidistant from the origin for various values of p . Note that we obtain the familiar Euclidean distance if we plug $p = 2$ in the above equation. Before we discuss more about this, let us talk a bit about how our intuition breaks down in “high” dimensions.

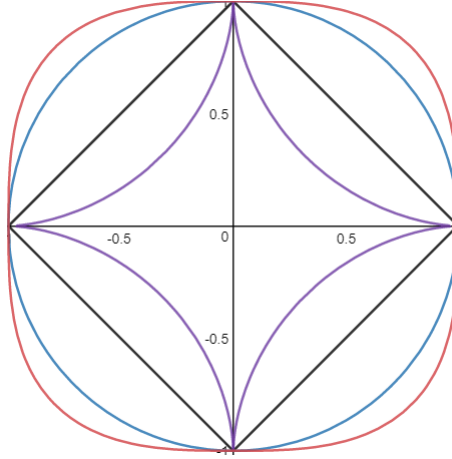


Figure 1:

2 Linear Regression

Linear regression is a relatively old technique, first published by Legendre in 1805 and then by Gauss in 1809. They developed the method to make sense of astronomical data.

Consider a dataset $\{\vec{x}_i, y_i\}$, where \vec{x}_i is an n -dimensional input vector of real values and y_i is the corresponding real valued output. Denoting the predicted output for a given input \vec{x}_i by \hat{y}_i , we can define the error term, which, when minimized outputs the quantity of interest.

$$E = \sum_i (\hat{y}_i - y_i)^2.$$

Here, we make an assumption which explains the term ‘linear’ in linear-regression. We assume that the *true* function f that relates Y and X is a linear function of X :

$$Y = f(X) + \epsilon,$$

$$Y = \vec{w} \cdot \vec{X} + b + \epsilon, \text{ this is the linear assumption,}$$

and our job is to estimate the function f , or, in other words, estimate the parameters \vec{w} and b . We denote the estimates using a hat, $\hat{\vec{w}}$ and \hat{b} . Once we have these estimates of parameters, we can compute predictions as follows:

$$\hat{y}_i = \hat{\vec{w}} \cdot \vec{x}_i + b.$$

Plugging in the above assumption into eq (4), we obtain the following:

$$E = \sum_i ((\hat{\vec{w}} \cdot \vec{x}_i + b) - y_i)^2.$$

The above equation can be written in a slightly different form, where the intercept term (also called bias) b is absorbed into the coefficients vector (also called weight vector). This is possible if we define another vector $\hat{\vec{w}}' = [\hat{b}, \hat{w}_1, \hat{w}_2, \dots]$ and $\vec{x}'_i = [1, x_{1i}, x_{2i}, \dots]$. We can now write:

$$E = \sum_i (\hat{\vec{w}}' \cdot \vec{x}'_i - y_i)^2.$$

It is a lot cleaner to solve the above minimization problem if we write the above equation in-terms of matrices and vectors. When this is done, we will be left with the following equation:

$$E = ||X\hat{\vec{W}} - Y||^2,$$

where the matrix X is made up by stacking the input vectors \vec{x}'_i in rows, \vec{Y} is the vector containing all the output values y_i and $\hat{\vec{W}} = \hat{\vec{w}}'$.

Minimizing the Error

We now minimize the error term with respect to the weights $\hat{\vec{W}}$. Differentiating the error term with respect to $\hat{\vec{W}}$ and setting it to zero, we obtain:

$$\begin{aligned} 2X^T(X\hat{\vec{W}} - Y) &= 0, \\ \rightarrow X^T X \hat{\vec{W}} - X^T Y &= 0. \end{aligned}$$

If X has more rows than columns, then it is very likely that the columns are independent and this implies that $X^T X$ is also invertible. Assuming X has independent columns, we can multiply the above equation by the inverse of $X^T X$ to end up with:

$$\hat{\vec{W}} = (X^T X)^{-1} X^T Y.$$

We thus have a nice closed form solution for the weigh vector. This is a sight to be cherished as something like this is hard to come by in machine-learning.

2.1 Statistical Inference for Linear Regression

One of the main advantages with linear regression is that it lends itself very well to statical analysis, thus permitting us to construct confidence-intervals and to perform hypothesis testing. Before we dive into this, let us look at the problem of confidence intervals and such for a far simpler case.

From introductory statistics, consider the classic problem of estimating the mean height of a population; we do not know anything about the population, but we do have access to a sample of size N . Let us call the sample mean $\hat{\mu}$ and sample standard deviation $\hat{\sigma}$. Given this information, the tools of statistical-inference help us construct a confidence interval as:

$$\left[\hat{\mu} - t_{n-1}^* \frac{\hat{\sigma}}{\sqrt{n}}, \hat{\mu} + t_{n-1}^* \frac{\hat{\sigma}}{\sqrt{n}} \right].$$

In a similar fashion, we can construct confidence intervals for $\hat{\vec{W}}$ as shown below:

$$\left[\hat{w}_i - t_{n-p}^* SE(\hat{w}_i), \hat{w}_i + t_{n-p}^* SE(\hat{w}_i) \right].$$

It can be shown that the estimates come from a t_{n-p} distribution with mean and variance derived below:

$$\begin{aligned} E[\hat{\vec{W}}] &= E[(X^T X)^{-1} X^T Y] \\ &= E[(X^T X)^{-1} X^T (X \vec{W} + \epsilon)] \\ &= E[(X^T X)^{-1} X^T X \vec{W} + (X^T X)^{-1} X^T \epsilon] \\ &= 0 \\ var(\vec{W}) &= var((X^T X)^{-1} X^T Y) \\ &= var((X^T X)^{-1} X^T (X \vec{W} + \epsilon)) \\ &= var((X^T X)^{-1} X^T X \vec{W} + (X^T X)^{-1} X^T \epsilon) \end{aligned}$$

Using $var(A\vec{B} + \vec{a}) = A var(\vec{B}) A^T$ (where \vec{a} is a constant), the above equation can be written as:

$$\begin{aligned} var(\vec{W}) &= (X^T X)^{-1} X^T var(\epsilon) ((X^T X)^{-1} X^T)^T \\ &= (X^T X)^{-1} X^T \sigma^2 \mathbf{I} ((X^T X)^{-1} X^T)^T \\ &= (X^T X)^{-1} X^T \sigma^2 \mathbf{I} X (X^T X)^{-1} \\ &\rightarrow var(\vec{W}) = \sigma^2 (X^T X)^{-1} \end{aligned}$$

Thus we obtain the co-variance matrix for the estimated parameters \vec{W} .

2.2 Modifications to Linear Regression

Taylor and Fourier Bases

In-order to deal with non-linear data, we will have to relax the assumption that we made in the earlier section i.e., instead of assuming that the predicted

output is a linear function of the input, we consider the output to be any general function of the input. We can demonstrate this in one-dimension as follows:

$$\text{linear} : \hat{y}_i = wx_i + b$$

$$\text{general} : \hat{y}_i = b + w_1x_i + w_2x_i^2 + w_3x_i^3 + \dots$$

$$\text{general} : \hat{y}_i = b + w_1 \cos(x_i) + w_2 \sin(x_i) + w_3 \sin(2x_i) + w_4 \cos(2x_i) + \dots$$

We can recognize the above two general functions as the Taylor and Fourier expansions of a function. Let us pick the Fourier basis to fit non-linear data as the polynomial one is dealt with in many books. Although by the looks of it, the non-linear problem looks quite different from the linear one, we can convert this problem into pretty much the exact same one as before.

$$\begin{aligned} x_i &\rightarrow [1, \cos(x_i), \sin(x_i), \sin(2x_i), \cos(2x_i), \dots], \\ X &\rightarrow \begin{bmatrix} 1 & \cos(x_1) & \sin(x_1) & \sin(2x_1) & \cos(2x_1) \dots \\ 1 & \cos(x_2) & \sin(x_2) & \sin(2x_2) & \cos(2x_2) \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}, \\ \vec{W} &\rightarrow [b, w_1, w_2, w_3, \dots]. \end{aligned}$$

Using the above definitions, we can again compute the weight vector using the same equation as before:

$$\vec{W} = (X^T X)^{-1} X^T Y.$$

2.3 Regularization for Linear Regression

Ridge Regularization

Compared to linear regression, here we minimize a slightly modified error term; to the residual-squared-sum, we add another term as shown below.

$$E_{\text{ridge}} = \|X\hat{\vec{W}} - \vec{Y}\|^2 + \lambda \|\hat{\vec{W}}\|_2^2.$$

Differentiating with respect to $\hat{\vec{W}}$ and setting it to zero, we obtain:

$$\begin{aligned} 2X^T(X\hat{\vec{W}} - \vec{Y}) + 2\lambda\hat{\vec{W}} &= 0, \\ \rightarrow X^T X\hat{\vec{W}} - X^T \vec{Y} + \lambda\hat{\vec{W}} &= 0, \\ \rightarrow (X^T X + \lambda\mathbf{I})\hat{\vec{W}} - X^T \vec{Y} &= 0, \\ \rightarrow \hat{\vec{W}} &= (X^T X + \lambda\mathbf{I})^{-1} X^T \vec{Y}. \end{aligned}$$

Lasso Regularization

$$E_{\text{lasso}} = \|X\hat{\vec{W}} - \vec{Y}\|^2 + \lambda \|\hat{\vec{W}}\|_1.$$

For lasso, there is no closed form solution for the general X . We use coordinate descent to arrive at the parameters that minimize the loss function.

3 Logistic Regression

Consider a multi-class dataset $\{x^{(i)}, y^{(i)}\}$, where y can take the values $0, 1, 2, \dots, K-1$. Say we have some $y^{(k)} = m$, then we can easily convert this to a K -dimensional vector with the m th element set to 1 and the rest set to zero. If we do this for all y values and if we consider the probability function to be a K -dimensional vector as well,

$$\Pr(Y = y^{(i)} | X = x^{(i)}) = p(x^{(i)}) = \frac{\exp(Wx^{(i)} + b)}{\|\exp(Wx^{(i)} + b)\|},$$

then we can represent the log-likelihood function as follows:

$$\mathcal{L} = \sum_i^m \sum_j^K \log(p_j(x^{(i)})) y_j^{(i)}.$$

4 Linear Discriminant Analysis

The Bayes' theorem states that

$$\Pr(Y = k | X = x) = \frac{\Pr(Y = k) \Pr(X = x | Y = k)}{\sum_{i=1}^K \Pr(Y = i) \Pr(X = x | Y = i)}.$$

Using the notation $\Pr(Y = k | X = x) = p_k(x)$, $\Pr(Y = k) = \pi_k$ and $\Pr(X = x | Y = k) = f_k(x)$, we can write the above theorem as:

$$p_k(x) = \frac{\pi_k f_k(x)}{\sum_{i=1}^K \pi_i f_i(x)}.$$

Linear Discriminant Analysis (LDA) makes use of the Bayes' theorem and a few assumptions on the probability density function, $f_k(x)$, in-order to classify points. The assumptions on the density function are: (1) it is gaussian with mean μ_k and variance Σ_k^2 (2) $\Sigma_1 = \Sigma_2 = \dots = \Sigma_K = \Sigma$.

For a given set of predictors x , LDA classification is done as follows:

$$\hat{y} = \operatorname{argmax}_k p_k(x).$$

Since the denominator term in $p_k(x)$ is simply a scaling factor, we can instead do the classification as:

$$\hat{y} = \operatorname{argmax}_k \pi_k f_k(x).$$

We simplify the above expression further by taking the log of the function that is being maximized.

$$\begin{aligned}
\log(\pi_k f_k(x)) &= \log(\pi_k) + \log(f_k(x)), \\
&= \log(\pi_k) + \log\left(\frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1}(x - \mu_k)\right)\right), \\
&= \log(\pi_k) + \log\left(\frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}}\right) + \log\left(\exp\left(-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1}(x - \mu_k)\right)\right), \\
&\text{(ignoring the term common to all densities)} \\
&= \log(\pi_k) + \log\left(\exp\left(-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1}(x - \mu_k)\right)\right), \\
&= \log(\pi_k) - \frac{1}{2}(x - \mu_k)^T \Sigma^{-1}(x - \mu_k), \\
&= \log(\pi_k) - \frac{1}{2} [x^T \Sigma^{-1} x - x^T \Sigma^{-1} \mu_k - \mu_k^T \Sigma^{-1} x + \mu_k^T \Sigma^{-1} \mu_k], \\
&\text{(again ignoring the term common to all densities)} \\
&= \log(\pi_k) - \frac{1}{2} [-x^T \Sigma^{-1} \mu_k - \mu_k^T \Sigma^{-1} x + \mu_k^T \Sigma^{-1} \mu_k], \\
&\text{(using } x^T \Sigma^{-1} \mu_k = \mu_k^T \Sigma^{-1} x) \\
&= \log(\pi_k) + x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k.
\end{aligned}$$

Hence, we can write:

$$\hat{y} = \operatorname{argmax}_k \log(\pi_k) + x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k.$$

4.1 Quadratic Discriminant Analysis

Quadratic discriminant analysis is exactly like LDA but without the second assumption that the density distributions have equal variances. Therefore, without the assumption that $\Sigma_1 = \Sigma_2 = \dots = \Sigma_K = \Sigma$, we can no longer ignore the terms we ignored in the above derivation. Thus we classify using the following equation:

$$\begin{aligned}
\hat{y} &= \operatorname{argmax}_k \log(\pi_k) + x^T \Sigma_{\mathbf{k}}^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma_{\mathbf{k}}^{-1} \mu_k - \frac{1}{2} x^T \Sigma_{\mathbf{k}}^{-1} x + \log\left(\frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma_{\mathbf{k}}|^{\frac{1}{2}}}\right), \\
&= \operatorname{argmax}_k \log(\pi_k) + x^T \Sigma_{\mathbf{k}}^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma_{\mathbf{k}}^{-1} \mu_k - \frac{1}{2} x^T \Sigma_{\mathbf{k}}^{-1} x - \frac{1}{2} \log(|\Sigma_{\mathbf{k}}|).
\end{aligned}$$

5 Back-propagation for binary Neural Networks

Let us consider a general neural network. No matter what the structure of the network is, given we are dealing with binary classification, the last layer will

contain only one node and the activation function for the last layer will be the sigmoid function. Let us do the computation for derivatives of the parameters W and b in detail for the last two layers of this general network i.e, we compute $W^{[L]}, b^{[L]}$ and $W^{[L-1]}, b^{[L-1]}$, where L denotes the final or output layer. One we have these, it is not too hard to see the jump to the general result.

$$\mathcal{L} = - \sum_{i=1}^N y^{(i)} \log(a^{[L](i)}) + (1 - y^{(i)}) \log(1 - a^{[L](i)})$$

1 Differentiating with respect to $a^{[L](j)}$, we get:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial a^{[L](j)}} &= - \sum_{i=1}^N \frac{\partial}{\partial a^{[L](j)}} \left[y^{(i)} \log(a^{[L](i)}) + (1 - y^{(i)}) \log(1 - a^{[L](i)}) \right] \\ &= - \frac{\partial}{\partial a^{[L](j)}} \left[y^{(j)} \log(a^{[L](j)}) + (1 - y^{(j)}) \log(1 - a^{[L](j)}) \right] \\ &= - \left[\frac{y^{(j)}}{a^{[L](j)}} - \frac{(1 - y^{(j)})}{1 - a^{[L](j)}} \right] \end{aligned}$$

Vectorized form of the above equation that lets us compute that above derivative for all j can be written as follows (we assume we are dealing with numpy arrays):

$$\frac{\partial \mathcal{L}}{\partial A^{[L]}} = - \left[\frac{Y}{A^{[L]}} - \frac{(1 - Y)}{1 - A^{[L]}} \right]$$

2 Now lets compute the derivative of loss with respect to z .

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial z^{[L](j)}} &= \frac{\partial \mathcal{L}}{\partial a^{[L](j)}} \frac{\partial a^{[L](j)}}{\partial z^{[L](j)}} \\ &= \frac{\partial \mathcal{L}}{\partial a^{[L](j)}} \frac{\partial g(z^{[L](j)})}{\partial z^{[L](j)}}, \text{ where } g \text{ is the activation function} \end{aligned}$$

In case of binary classification, $g(z) = \sigma(z)$, where σ is the sigmoid function. Therefore:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial z^{[L](j)}} &= \frac{\partial \mathcal{L}}{\partial a^{[L](j)}} \frac{\partial \sigma(z^{[L](j)})}{\partial z^{[L](j)}} \\ &= \frac{\partial \mathcal{L}}{\partial a^{[L](j)}} \frac{\partial}{\partial z^{[L](j)}} \left[\frac{1}{1 + \exp(-z^{[L](j)})} \right] \\ &= \frac{\partial \mathcal{L}}{\partial a^{[L](j)}} \left(\frac{1}{1 + \exp(-z^{[L](j)})} \right) \left(1 - \frac{1}{1 + \exp(-z^{[L](j)})} \right) \\ &= \frac{\partial \mathcal{L}}{\partial a^{[L](j)}} a^{[L](j)} (1 - a^{[L](j)}) \end{aligned}$$

Again, we can write the above equation in vectorized form as:

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial Z^{[L]}} &= \frac{\partial \mathcal{L}}{\partial A^{[L]}} * A^{[L]}(1 - A^{[L]}), \\ &= - \left[\frac{Y}{A^{[L]}} - \frac{(1 - Y)}{a^{[L]}(j)} \right] * A^{[L]}(1 - A^{[L]}),\end{aligned}$$

where ‘*’ indicates element-wise multiplication.

3

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial W_j^{[L]}} &= \sum_{i=1}^N \frac{\partial \mathcal{L}}{\partial a^{[L]}(i)} \frac{\partial a^{[L]}(i)}{\partial z^{[L]}(i)} \frac{\partial z^{[L]}(i)}{\partial W_j^{[L]}} \\ &= \sum_{i=1}^N \frac{\partial \mathcal{L}}{\partial a^{[L]}(i)} \frac{\partial a^{[L]}(i)}{\partial z^{[L]}(i)} \frac{\partial}{\partial W_j^{[L]}} \left(\sum_{l=1}^{n_{L-1}} a_i^{[L-1](j)} W_l^{[L]} + b_j^{[L]} \right) \\ &= \sum_{i=1}^N \frac{\partial \mathcal{L}}{\partial a^{[L]}(i)} \frac{\partial a^{[L]}(i)}{\partial z^{[L]}(i)} a_j^{[L-1](i)}\end{aligned}$$

Therefore, the vectorized form is:

$$\frac{\partial \mathcal{L}}{\partial W^{[L]}} = \frac{\partial \mathcal{L}}{\partial Z^{[L]}} \cdot A^{[L-1]T}.$$

4

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial b_j^{[L]}} &= \sum_{i=1}^N \frac{\partial \mathcal{L}}{\partial a^{[L]}(i)} \frac{\partial a^{[L]}(i)}{\partial z^{[L]}(i)} \frac{\partial z^{[L]}(i)}{\partial b_j^{[L]}} \\ &= \sum_{i=1}^N \frac{\partial \mathcal{L}}{\partial a^{[L]}(i)} \frac{\partial a^{[L]}(i)}{\partial z^{[L]}(i)} \frac{\partial}{\partial b_j^{[L]}} \left(\sum_{l=1}^{n_{L-1}} a_i^{[L-1](j)} W_l^{[L]} + b_j^{[L]} \right) \\ &= \sum_{i=1}^N \frac{\partial \mathcal{L}}{\partial a^{[L]}(i)} \frac{\partial a^{[L]}(i)}{\partial z^{[L]}(i)} (1)\end{aligned}$$

Therefore, the vectorized form is:

$$\frac{\partial \mathcal{L}}{\partial W^{[L]}} = \frac{\partial \mathcal{L}}{\partial Z^{[L]}} \cdot \mathbf{1},$$

where $\mathbf{1}$ is a row vector of ones of dimension $N \times 1$.

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial a_j^{[L-1]}} &= \frac{\partial \mathcal{L}}{\partial a^{[L](i)}} \frac{\partial a^{[L](i)}}{\partial z^{[L](i)}} \frac{\partial z^{[L](i)}}{\partial a_j^{[L-1]}} \\
&= \frac{\partial \mathcal{L}}{\partial a^{[L](i)}} \frac{\partial a^{[L](i)}}{\partial z^{[L](i)}} \frac{\partial}{\partial a_j^{[L-1]}} \left(\sum_{l=1}^{n_{L-1}} a_l^{[L-1](j)} W_l^{[L]} + b_j^{[L]} \right) \\
&= \frac{\partial \mathcal{L}}{\partial a^{[L](i)}} \frac{\partial a^{[L](i)}}{\partial z^{[L](i)}} W_j^{[L]}
\end{aligned}$$

Therefore, the vectorized form is:

$$\frac{\partial \mathcal{L}}{\partial A^{[L-1]}} = \frac{\partial \mathcal{L}}{\partial Z^{[L]}} * W^{[L]}.$$

Equations we derived

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial A^{[L]}} &= - \left[\frac{Y}{A^{[L]}} - \frac{(1-Y)}{A^{[L]}} \right] \\
\frac{\partial \mathcal{L}}{\partial Z^{[L]}} &= \frac{\partial \mathcal{L}}{\partial A^{[L]}} * A^{[L]} (1 - A^{[L]}) \\
\frac{\partial \mathcal{L}}{\partial W^{[L]}} &= \frac{\partial \mathcal{L}}{\partial Z^{[L]}} \cdot A^{[L-1]T} \\
\frac{\partial \mathcal{L}}{\partial W^{[L]}} &= \frac{\partial \mathcal{L}}{\partial Z^{[L]}} \cdot \mathbf{1} \\
\frac{\partial \mathcal{L}}{\partial A^{[L-1]}} &= \frac{\partial \mathcal{L}}{\partial Z^{[L]}} * W^{[L]}
\end{aligned}$$

General Equations

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial A^{[l]}} &= - \left[\frac{Y}{A^{[l]}} - \frac{(1-Y)}{A^{[l]}} \right] \\
\frac{\partial \mathcal{L}}{\partial Z^{[l]}} &= \frac{\partial \mathcal{L}}{\partial A^{[l]}} * g'(Z^{[l]}) \\
\frac{\partial \mathcal{L}}{\partial W^{[l]}} &= \frac{\partial \mathcal{L}}{\partial Z^{[l]}} \cdot A^{[l-1]T} \\
\frac{\partial \mathcal{L}}{\partial W^{[l]}} &= \frac{\partial \mathcal{L}}{\partial Z^{[l]}} \cdot \mathbf{1} \\
\frac{\partial \mathcal{L}}{\partial A^{[l-1]}} &= W^{[l]T} \cdot \frac{\partial \mathcal{L}}{\partial Z^{[l]}}
\end{aligned}$$