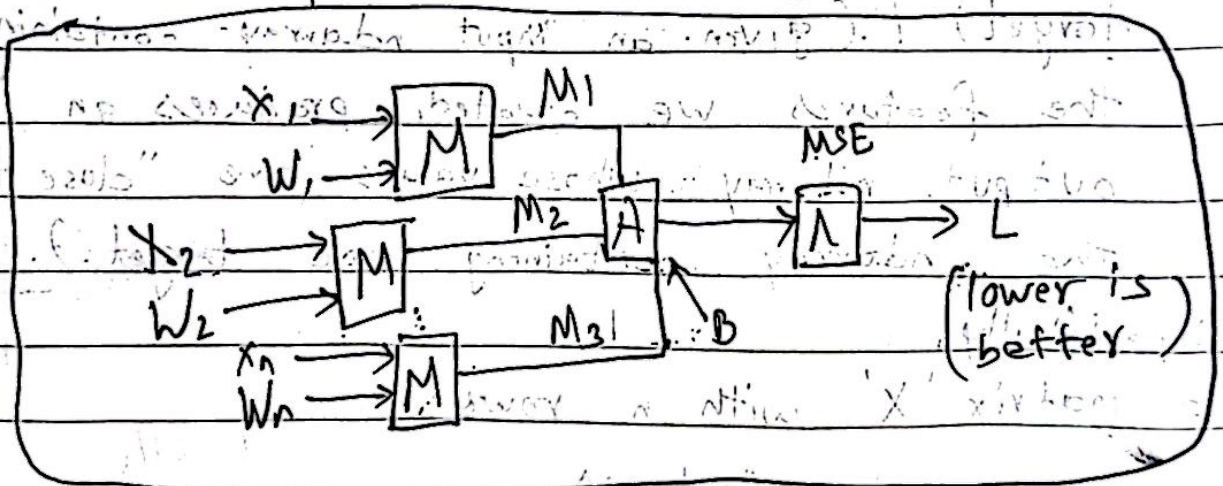


takes X batch & produce vectors of values p_i "predictions" which are "close" to the target values y_i

Linear Regression

$$y_i = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n + \epsilon$$

Numeric value of each target is a linear combination of the k features of X , plus the β_0 term to adjust the "baseline"



Suppose $X = [x_1 \ x_2 \ \dots \ x_k]$

$$W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_k \end{bmatrix}$$

$$p_i = x_i \cdot W = w_1 x_1 + w_2 x_2 + \dots + w_k x_k$$

"generating the prediction" for LR using:
the Dot Product for vectors

matrix multiplication - for matrix

For a batch of size 3.

$$X_{\text{batch}} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1k} \\ x_{21} & x_{22} & \dots & x_{2k} \\ x_{31} & x_{32} & \dots & x_{3k} \end{bmatrix}$$

$$p_{\text{batch}} = X_{\text{batch}} \times W = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1k} \\ x_{21} & x_{22} & \dots & x_{2k} \\ x_{31} & x_{32} & \dots & x_{3k} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_k \end{bmatrix}$$

$$= \begin{bmatrix} x_{11}w_1 + x_{12}w_2 + \dots + x_{1k}w_k \\ x_{21}w_1 + x_{22}w_2 + \dots + x_{2k}w_k \\ x_{31}w_1 + x_{32}w_2 + \dots + x_{3k}w_k \end{bmatrix} = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \end{bmatrix}$$

"Training" this model

takes data, combine them with parameters in some way, & produce predictions.

→ X as data, W as parameter, batch prediction using "Matrix multiplication"

whether or not these observations are good.

We use targets y_{batch} & associated with the batch of observations X_{batch} fed into the function, and we compute a single number, function of y_{batch} & p_{batch} .

"penalty". Eg: "Mean Squared Error"

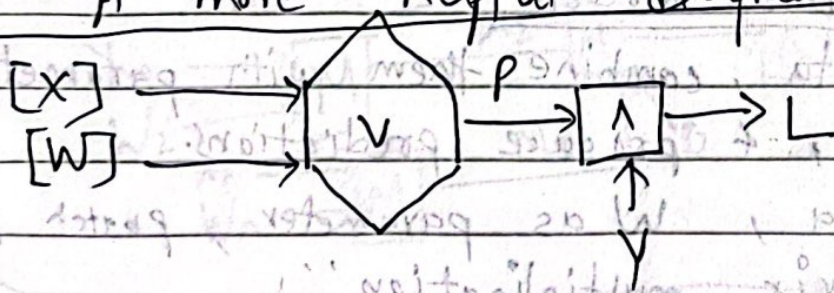
$$\text{MSE}(P_{\text{batch}}, Y_{\text{batch}}) = \text{MSE} \begin{pmatrix} p_1 \\ p_2 \\ p_3 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix}$$

$$= \frac{(y_1 - p_1)^2 + (y_2 - p_2)^2 + (y_3 - p_3)^2}{3}$$

We use all techniques from chapter 1 to compute the "gradient" of this number with respect to W .

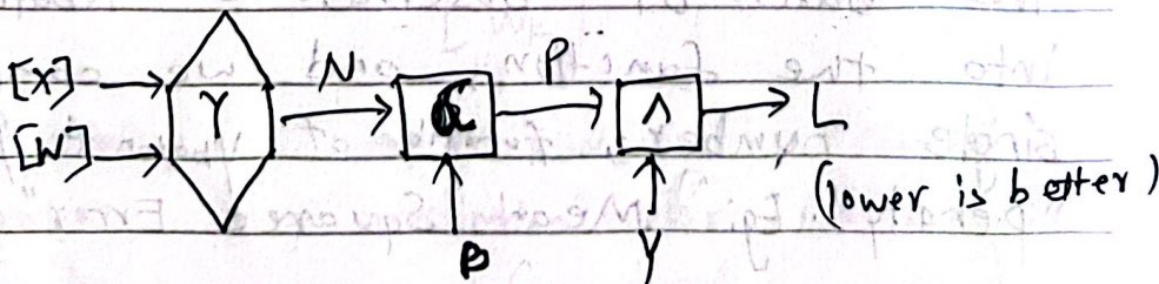
Then "we can use these derivatives to update each element of W in the direction that would cause L to decrease".
 Repeat \rightarrow "train"

LR - A more helpful diagram



$$L = \lambda(v(X, W), y)$$

Adding the intercept in



$p_{\text{batch-with-bias}} = x_i \cdot w + b$

$$= \begin{bmatrix} x_{11}w_1 + x_{12}w_2 + \dots + x_{1k}w_k + b \\ x_{21}w_1 + x_{22}w_2 + \dots + x_{2k}w_k + b \\ x_{31}w_1 + x_{32}w_2 + \dots + x_{3k}w_k + b \end{bmatrix} = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \end{bmatrix}$$

"the same number" ~~is~~ should get added as bias.
y-intercept

LR - the code

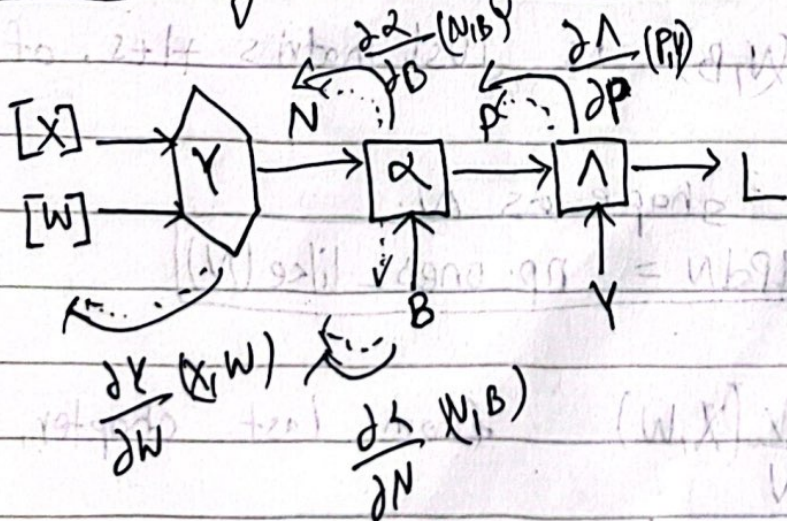
makes prediction & computes losses given
batches of observations x_{batch} & their correspo-
nding targets y_{batch} .

Training the model

"compute" $\frac{\partial L}{\partial w_i}$ for every w_i in W ,

as well as $\frac{\partial L}{\partial b}$

Calculating the gradients: A diagram



Calculating the gradients

From figure, we can see that the derivative product we ultimately want to compute is:

$$\frac{\partial \Lambda}{\partial P}(P, Y) \times \frac{\partial \alpha}{\partial N}(N, B) \times \frac{\partial V}{\partial W}(X, W)$$

First up: $\frac{\partial \Lambda}{\partial P}(P, Y)$

$\Lambda(P, Y) = (Y - P)^2$ for each element in Y if

$$\frac{\partial \Lambda}{\partial P}(P, Y) = 2(Y - P) \cdot (-1) \Rightarrow -2(Y - P)$$

code :- $dLdP = -2 * (Y - P)$

Net up: $\frac{\partial \alpha}{\partial N}(N, B)$, matrices, but α is just "addition"

With same logic as in chapter 1:

increasing any element of N by one unit will increase $P = \alpha(N, B) = N + B$ by one unit

Thus, $\frac{\partial \alpha}{\partial N}(N, B)$ is just matrix 1 's, of the same shape as N .

code :- $dPdN = \text{np.ones_like}(N)$

Finally: $\frac{\partial V}{\partial W}(X, W)$ from last chapter,

When computing derivatives of nested functions

Where one of the constituent functions is a matrix multiplication, we act as if:

$$\frac{\partial v}{\partial W}(x, W) = x^T$$

code: $dN dW = \text{np.transpose}(x, (1, 0))$

$dD dB = \text{np.ones_like}(weights['B'])$
Same as $dD dN$

Calculating the gradients: The (Full) calc
Goal: take all computed & inputted into the forward pass
of X, W, N, B, P and y

and compute:

$$\frac{\partial \mathcal{L}}{\partial W} \text{ \& \& } \frac{\partial \mathcal{L}}{\partial B}$$

'W' & 'B' as inputs in a Dict. called "weights"
rest in a Dict. Forward-info:

returns loss gradients, dict. containing
weights & Bias: 'W', 'B'

$(\text{dloss}, \text{dict}(X))$

$\frac{\partial \mathcal{L}}{\partial W} = \text{np.dot}(dN, W)$

→ iterate over

Using These gradients to train the Model.

1. Select a batch of data
2. Run the forward pass of the model
3. Run the backward pass of the model using the info computed on the forward pass.
4. Use the gradients computed on the backward pass to update the weights.

Shuffling the data ~~or~~ to ensure that it is fed through in a random order.

We run the train function for certain num of "epochs", or cycles through the entire training dataset,

parameters :- weights, bias

hyper parameters :- learning rate, epoch, batch size,

Prediction:

We simply use the weights returned earlier from the train function & write:

code:- $\text{preds} = \text{predict}(X_{\text{test}}, \text{weights})$

We just do matrix multiplication of X_{test} & weights [W] & add [weights [B]]

metrics

As RMSE is particularly common metric. Since it is on a same scale as the target. If we divide this number by mean of the target, we can get a measure of how far off a prediction is, on average, from its actual value.

We calculate the mean of y -test: 22.0776

$$\frac{\text{RMSE}}{\text{mean}} = \frac{5.0508}{22.0776} = 0.2291 \text{ on average.}$$

Analyzing the most important feature.

The reason why we do this analyze specific to Linear Regression is that we can interpret the absolute values of the coefficients as corresponding to the importance of the different features to the model; larger coefficient means the feature is more important.

Neural Networks from Scratch:

Step 1: A bunch of Linear Regressions

$$X \equiv [\text{batch_size}, \text{num_features}] \Rightarrow [\text{batch_size}, 1]$$
$$W [\text{num_features}, 1]$$

matrix multiply

To do multiple regressions, we'll simply multiply our input by a weight matrix with dimensions

$[\text{num_features}, \text{num_outputs}]$ resulting in an output of dimensions $[\text{batch_size}, \text{num_outputs}]$; now, for each observation, we have num-outputs different weighted sums of the original features.

↳ We should think of each of them as "learned features" - a combination of the original features that, once the network is trained, will represent to learn combinations of features that help it accurately predict.

Step 2: Non linear function

We will use sigmoid function to learn the non-linear pattern in the features. Also,

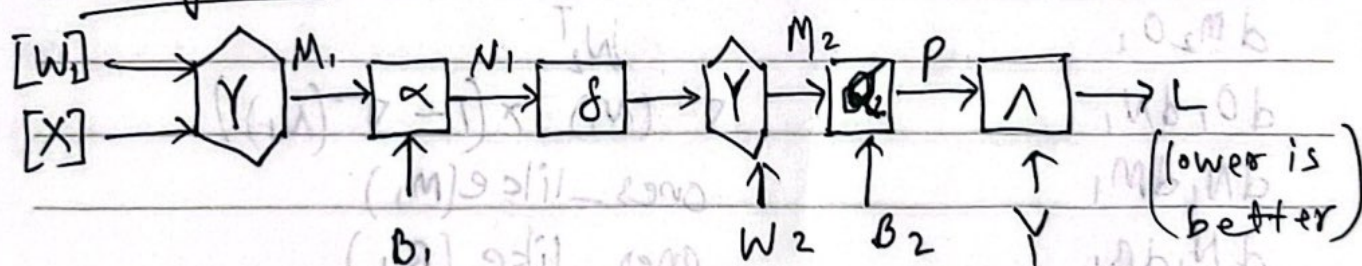
to "preserve" information about the numbers that were fed in.

$$\frac{\partial \sigma(x)}{\partial u} = \sigma(x) \times (1 - \sigma(x))$$

Step 3: Another Linear Regression

we will feed the resulting elements from step 2 through containing values betⁿ 0 and 1 into a regular linear regression and then we calculate MSE.

Diagram:



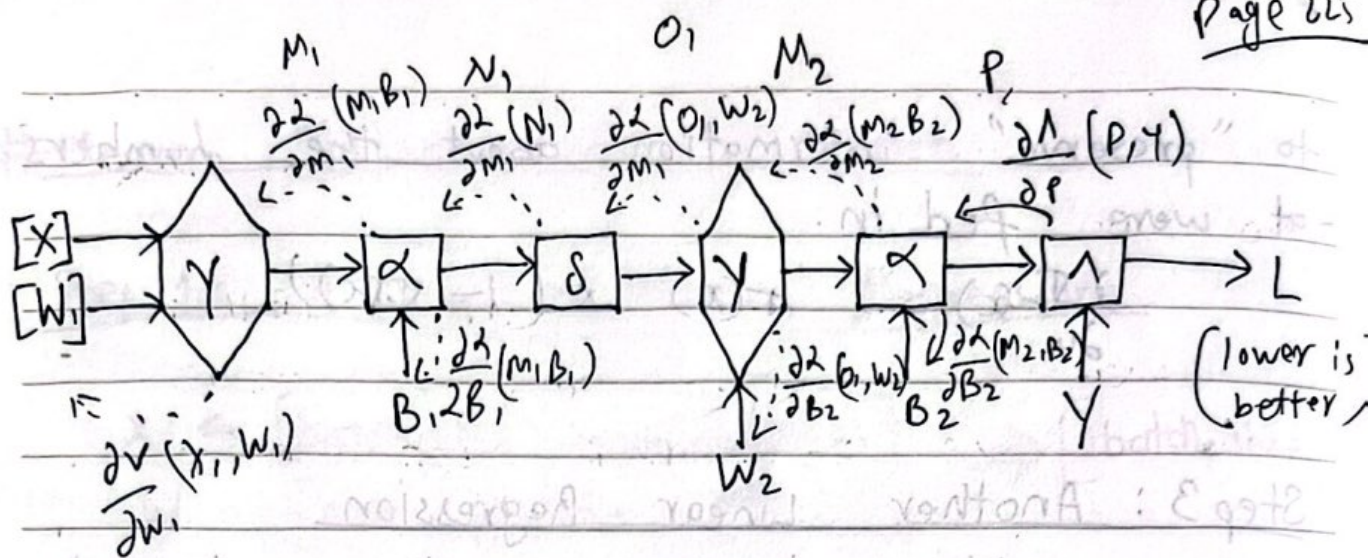
$W_1, W_2 \rightarrow$ weights matrix

$B_1, B_2 \rightarrow$ biases

$P \rightarrow$ final prediction

Neural Network - Backward Pass

- (1) Compute the derivative of each operation and evaluate it at its input.
- (2) Multiply the results together.



Derivative rep. Pseudocode formula

| | |
|-----------------|--|
| $dL dP$ | $-(Y - P)$ |
| $dP dM_2$ | ones-like (M_2) |
| $dP d\beta_2$ | ones-like (β_2) |
| $dM_2 dW_2$ | O_1^T |
| $dM_2 dO_1$ | W_2^T |
| $dO_1 dN_1$ | $\sigma(N_1) \times (1 - \sigma(N_1))$ |
| $dN_1 dM_1$ | ones-like (M_1) |
| $dN_1 d\beta_1$ | ones-like (β_1) |
| $dM_1 dW_1$ | X^T |

From chain Rule we multiply to find:

$$dL dW_2 = dL dP \times dP dM_2 \times dM_2 dW_2$$

$$dL d\beta_2 = dL dP \times dP dM_2 \times \sum dM_2 d\beta_2$$

$$dL d\beta_1 = dL dP \times dP dM_2 \times dM_2 dO_1 \times dO_1 dN_1 \times dN_1 d\beta_1$$

$$dL dW_1 = dL dP \times dP dM_2 \times dM_2 dO_1 \times dO_1 dN_1 \times dN_1 dM_1 \times dM_1 dW_1$$

$$\rightarrow dL d\beta_1 = dL dN_1 \times dN_1 \beta_1 \cdot \text{sum}(\text{axis}=0)$$

Since we added 2-bias, we need to sum along axis=0 while computing derivative of output of Fully connected layer with respect to bias

Neural Network perform better than Linear Regression

- the NN is able to learn nonlinear relationships between the features and the target
- learn the relationships between combinations of features and the target.

What it means for deep learning model to have multiple "hidden layers" is that it means it means to have multiple "hidden layers" of variables

Deep Learning Definition - A First Pass

Deep-learning models are represented by series of operations that have at least two, non-consecutive non-linear functions involved.