

Logistic Regression Model

1) Preprocess data into h5 file.

- Loop through all images
- resize to 64×64 . \Rightarrow convert to RGB.
- assign label \rightarrow (1 = cat, 0 = not cat).
- Convert to array
- train_test_split () 80% train
20% testing.

Spent
more
time being
there.

h5 create file

2) Load data out of h5.

3) Shape: $\text{train_x} = (209, 64, 64, 3)$ $\rightarrow m = \frac{209 \times 64 \times 64 \times 3}{255}$ Pictures.
(labels): $\text{train_y} = (209, 1)$ \rightarrow 1d vector, 1 or 0
of cat.

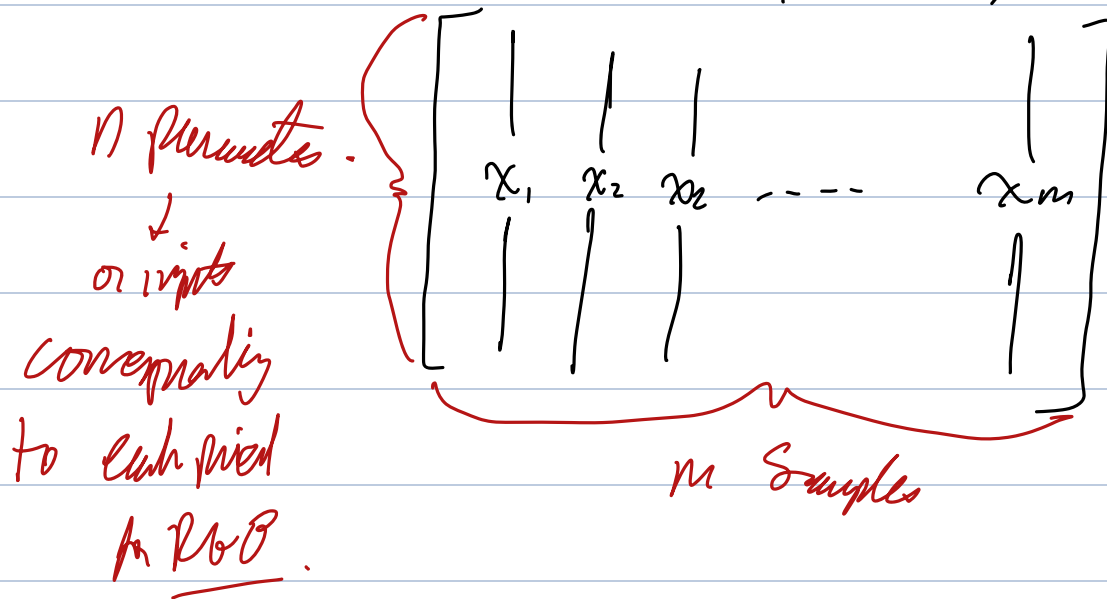
$\text{test_x} = (50, \dots)$

$\text{test_y} = (50, 1)$

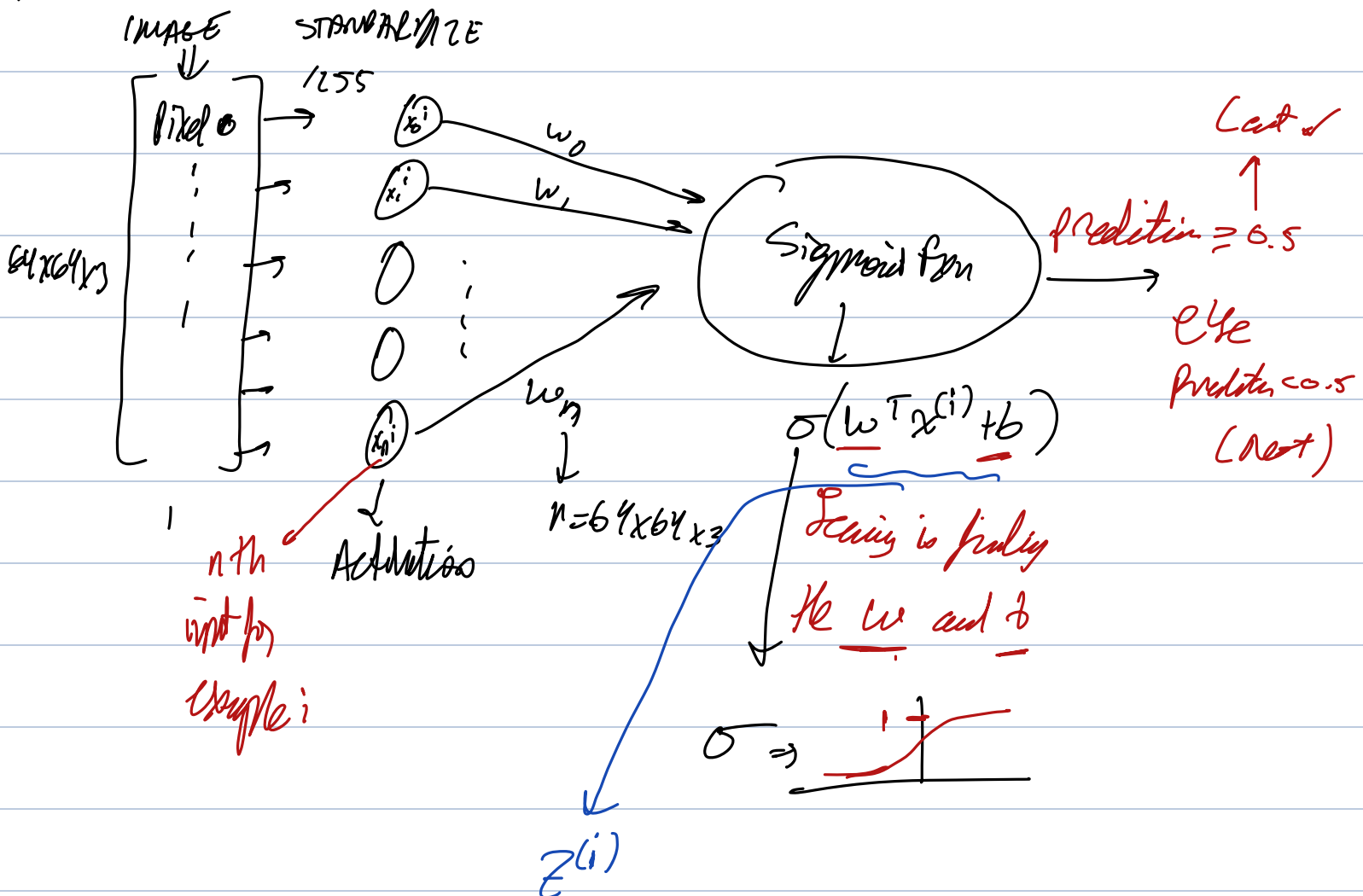
Then Flatten into a matrix:

$\text{train_x} \Rightarrow (209, 64, 64, 3)$

$$\Rightarrow (64 \times 64 \times 3, 209)$$



Local Architecture



$$z^{(i)} = w^T x^{(i)} + b.$$

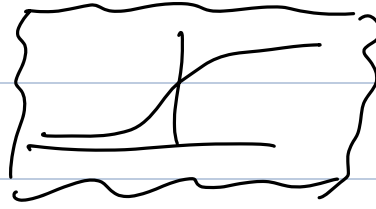
$$\hat{y} = a^{(i)} = \sigma(z^{(i)}) \rightarrow \text{prediction}$$

$$\underbrace{\mathcal{L}(a^{(i)}, y^{(i)})}_{\text{Loss Fxn}} = -y^{(i)} \log(a^{(i)}) - (1-y^{(i)}) \log(1-a^{(i)})$$

Ultimately this is a MLE problem

(MAXIMUM LIKELIHOOD ESTIMATE)

$$\text{Estimator} \Rightarrow \hat{y} = \sigma(w^T x + b).$$

$$\sigma(z) = \frac{1}{1+e^{-z}}$$


What is the likelihood fcn?

Definition: $P(\text{data} | p) = ?$. Given some unknown probability p , we assume we know p , what is the probability of getting the data, if we assume p .

Assume we have x , what is probability of y ? \hat{y}
 B/c we're using \hat{y} to estimate y .

$$\text{if } y=1 \Rightarrow P(y|x) = \hat{y} = \sigma(w^T x + b)$$

$$\text{if } y=0 \Rightarrow P(y|x) = 1 - \hat{y} = 1 - \sigma(w^T x + b)$$

But the mapping of x to y that $\in [0, 1]$, is for two intervalled neurons of our computer being 0 or 1, and that the y is the probability of a cat.

In logistic regression, assume $P(Y=1|x) = \hat{y} = \sigma(z)$
is a BERNOUlli DISTRIBUTION
success or failure.

$$\therefore P(Y=0|x) = 1 - \hat{y}.$$

\therefore The main reason is the probabilities are between 0 and 1, otherwise we could use a SOFTMAX Fun instead of sigmoid.

Now we have the two, but what is

$$P(Y=1|x) = \hat{y} \quad P(Y=0|x) = 1 - \hat{y}.$$

$$P(Y=y|x) = \hat{y}^y \cdot (1 - \hat{y})^{1-y}$$

$$\therefore \text{if } y=1 \Rightarrow \hat{y}^1 \cdot (1 - \hat{y})^{1-1} = \hat{y}$$

$$\text{if } y=0 \Rightarrow \hat{y}^0 \cdot (1 - \hat{y})^1 = 1 - \hat{y}$$

$$\therefore P(Y=y|x) = \hat{y}^y (1 - \hat{y}^n)^{1-y} \quad \text{Likelihood Fun}$$

We want to find the MLE for y .

Take log.

$$\log(P(Y=y|x)) = y \log \hat{y} + (1-y) \log (1 - \hat{y})$$

not just x , \Rightarrow w and b

$$\log(L(w, b)) = y \log y + (1-y) \log(1-y)$$

↓
equation

or likelihood?

$$\hat{y} = \sigma(z) = \sigma(w^T x + b)$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

→ log probability
of the observed data
given parameters.

So we need to maximize MLE of $\log(L(w, b)) \rightarrow$

$$\log(P(Y=y|x, w, b)) = y \log(\sigma(w^T x + b)) + (1-y) \log(1 - \sigma(w^T x + b))$$

∴ TAKE PARTIALS TO OPTIMIZE

$$\frac{\partial L(w, b)}{\partial w}$$

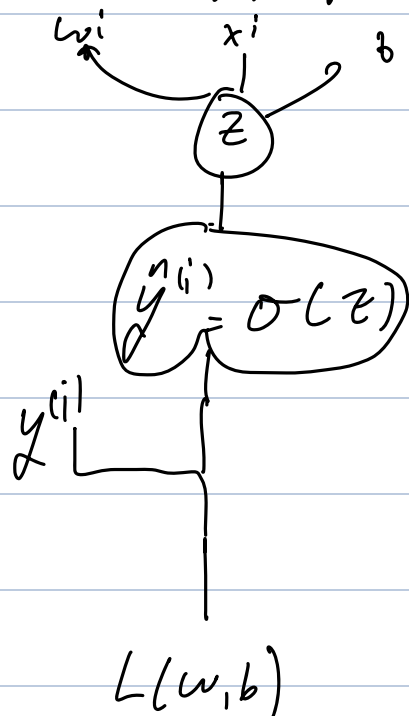
$$\text{and } \frac{\partial L(w, b)}{\partial b}$$

but there are w 's
 $w^1 \dots w^n$

So first, recall we have m examples:

$$\log(L(w, b)) = \sum_{i=1}^m y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)})$$

Now consider the chain rule and then there are



$$\frac{\partial L(w, b)}{\partial w^{(i)}} = \frac{\partial L(w^{(i)}, b)}{\partial \hat{y}^{(i)}} \cdot \frac{\partial \hat{y}^{(i)}}{\partial z^{(i)}} \cdot \frac{\partial z^{(i)}}{\partial w^{(i)}}$$

$$= \left(\sum_{i=1}^m \frac{y^{(i)}}{\hat{y}^{(i)}} + \frac{1-y^{(i)}}{1-\hat{y}^{(i)}} (-1) \right) (\sigma'(z^{(i)})) (x^{(i)})$$

$$= \left(\sum_{i=1}^M \frac{y - \hat{y} - \hat{y} + \hat{y}}{\hat{y}(1-\hat{y})} \right) (\sigma'(z^{(i)})) (x^{(i)})$$

$$= \left(\sum_{i=1}^M \frac{y - \hat{y}}{\hat{y}(1-\hat{y})} \right) (\sigma'(z^{(i)})) (1 - \sigma(z^{(i)})) (x^{(i)})$$

$$= \left(\sum_{i=1}^M \frac{y^i - \hat{y}^i}{\hat{y}^i(1-\hat{y}^i)} \right) (\hat{y}^{(i)} (1 - \hat{y}^{(i)})) (x^{(i)})$$

$$\frac{\partial}{\partial w} = \left(\sum_{i=1}^M y^{(i)} - \hat{y}^{(i)} \right) x^{(i)}$$

Similarly we can get $\frac{\partial}{\partial b} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial b} \rightarrow 1$

$$\Rightarrow \frac{\partial}{\partial b} = \left(\sum_{i=1}^M y^{(i)} - \hat{y}^{(i)} \right)$$

Step 4: Could We Just Use Gradient Ascent?

Yes, we could use gradient ascent to maximize $\log L(w, b)$ directly. However:

- Most ML frameworks (like TensorFlow, PyTorch, Scikit-Learn) are built around minimization.
- Using a minimization approach keeps it consistent with other models (like linear regression, which minimizes squared error).
- It's easier to interpret a "loss function" than a "likelihood function."

Step 5: Summary

Approach	Goal	Why We Do It
MLE (Maximize Likelihood)	Find parameters that maximize $\log L(w, b)$	Theoretical foundation for logistic regression
Gradient Descent (Minimization)	Find parameters that minimize the negative log-likelihood $J(w, b)$	Works better with standard ML optimization algorithms
Why Negative?	So we can minimize instead of maximize	Consistent with other ML models (e.g., squared error in linear regression)

Although even here we want to maximize MLE

Cost function

it makes more logical sense to end to minimize the loss over time.

$$L = \text{Loss function} = -\log(L(w, b))$$

$$\Rightarrow \frac{\partial L}{\partial w} = \sum (\hat{y}^{(i)} - y^{(i)}) x^{(i)}$$

$$\frac{\partial L}{\partial b} = \sum (\hat{y}^{(i)} - y^{(i)})$$

We want to plug the cost function into the loss function

$$J(w, b) = \frac{1}{n} \sum_{i=1}^n L(w, b)$$

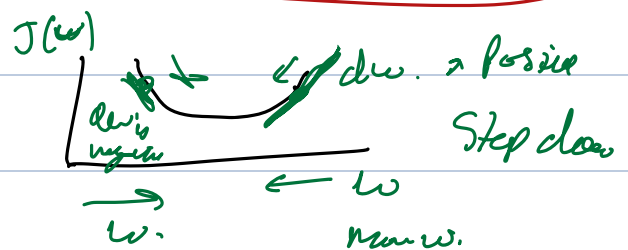
$$\frac{\partial J}{\partial w} = \frac{1}{n} \sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)}) x^{(i)}$$

$$\frac{\partial J}{\partial b} = \frac{1}{n} \sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)})$$

Gradient function is nonlinear, but just not = 0 and solve for w and b.

USE ITERATIVE OPTIMIZATION TECHNIQUES INSTEAD

GRADIENT DESCENT. Image



$$w = w - \alpha \frac{\partial J}{\partial w} \quad \text{To minimize } J(w, b) \text{ and}$$

$$b = b - \alpha \frac{\partial J}{\partial w}$$

learning rate.

maximizing likelihood

Have big of a step to take.

Difficult things like newtons method to find but is expensive

(TRANSIENT, ETC. CH2C 2)
NOTES

\therefore TRAINING (GRADIENT DESCENT)

why it works b/c \hat{y} is the derivative depends on w , and b ,
So it'll update itself as it goes.