

~~SUPERVISED MACHINE LEARNING: REGRESSION & CLASSIFICATION~~

~~CLASSIFICATION~~

~~EDUCATIONAL : ENTRE. TO MACHINE LEARNING~~

~~KK1~~ ~~DATA BASES~~

~~WHAT IS MACHINE LEARNING~~

SUPERVISED vs. UNSUPERVISED MACHINE LEARNING

- Machine learning is a field of study that gives computers the ability to learn without being explicitly programmed.

- Common ML algorithms →
 - (i) Supervised learning

[i] Unsupervised learning [ii] Recommender systems [iii] Reinforcement learning

(A) SUPERVISED LEARNING:

Refers to algorithm that takes $x \rightarrow y$ or input to output mapping. Key characteristic = in Supervised learning, we give the learning ^{algorithm} examples to learn from that include the right answers i.e. correct label, y for a given ⁺ input x .

- It's by using those examples and learning its relationships, that the algorithm is able to give a reasonably accurate prediction of the output given only the input.

- Regression, a type of supervised learning, refers to predicting a continuous number from infinitely many possible outputs.
- Classification algos, another type of supervised learning. This is different from regression because we're trying to predict only a small number of output while regression try to predict a continuous number.
→ output class or output category = output.

Classification algorithms predict categories and they don't have to be numbers
→ that needs how to fit a boundary line to the data in case of 2 or more inputs.

(B) UNSUPERVISED LEARNING

~~Topic~~ In the supervised learning, here we're given data that is not associated with any output labels y . It involves finding something interesting in unlabeled data.

Algorithm types : (i) Clustering used in google news [words ~~in~~ in heading to rate top news]

(ii) Data only comes w/ inputs & no output labels y . The algorithm has to find structure in the data.

(iii) Clustering involves grouping similar data points together

(iv) Anomaly detection = involves finding unusual data points useful for fraud detection

(v) Dimensionality reduction = compressing data from big to small w/out losing too much information.

WK 1 REGRESSION MODEL

Cost Function: Takes the predictor \hat{y} & subtracts the target y from it

the squares the difference. $\frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 = J_{(w,b)}$

m = number of training examples

Squared error cost function.

$$J_{(w,b)} = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

$$\hat{y}^{(i)} = f_{w,b}(x^{(i)})$$

$$f_{w,b}(x^{(i)}) = w x^{(i)} + b$$

In linear regression, we use the

cost function to find the value of w & b plotting all possible values of

w & b that minimize $J_{(w,b)}$

Goal of Linear Regression.

WK 1 TRAIN THE MODEL w/ GRADIENT DESCENT

A computational way of finding the exact w & b that result in the min. $J_{(w,b)}$

Gradient descent can work w/ several parameters to minimize the cost function

↳ Interesting property of Gradient Descent: It contains local minima i.e.

different loss points ^{w/ only one very} that can be reached based on your starting point.

Gradient descent is done in such away that whenever you are, you take the steepest baby step and keep on taking the steepest baby steps around you until you reach a local minimum.

Gradient Descent Algorithm

$$\textcircled{1} \quad w = w - \alpha \left(\frac{\partial}{\partial w} J(w, b) \right)$$

derivative of $J(w, b)$ tells us the direction & the step

$$\textcircled{2} \quad b = b - \alpha \left(\frac{\partial}{\partial b} J(w, b) \right)$$

learning rate → learning rate, α , controls how big of a step you take downhill.

large α = large steps, small α = baby steps.

- In the gradient descent algorithm, you repeat the update step $\textcircled{1} \& \textcircled{2}$ until the algorithm converges. Convergence at a local minimum where the parameters $w \& b$ no longer change much with each step taken.

- Note: Both $w \& b$ in $\textcircled{1} \& \textcircled{2}$ are updated simultaneously

Example: 1. $\text{tmp_w} = w - \alpha \frac{\partial}{\partial w} J(w, b)$

2. $\text{tmp_b} = b - \alpha \frac{\partial}{\partial b} J(w, b)$

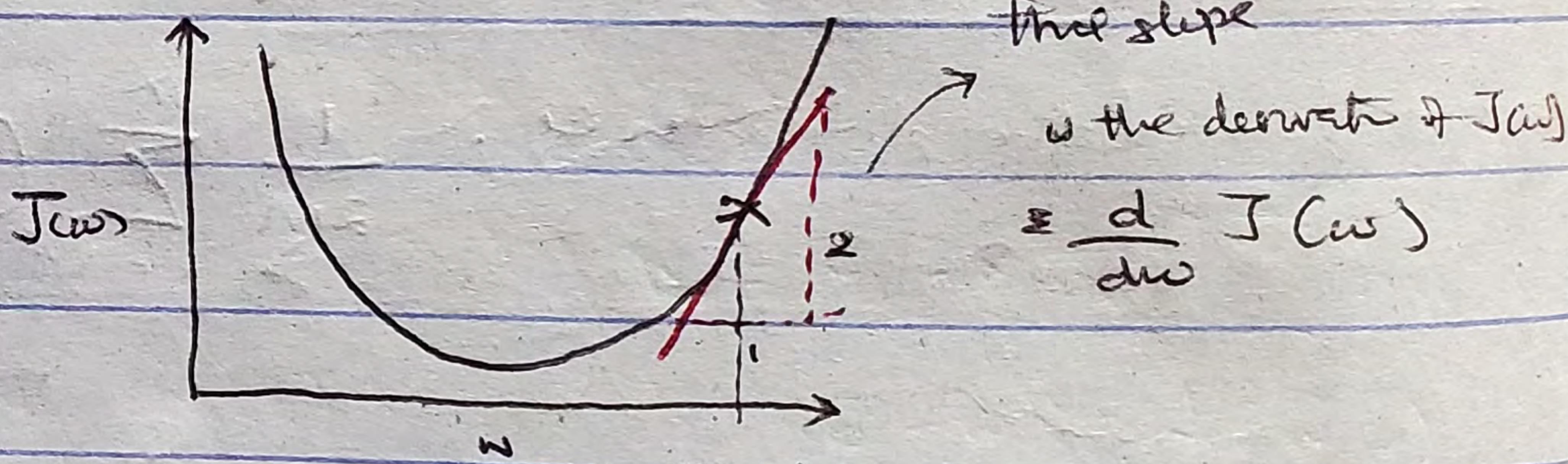
3. $w, b = \text{tmp_w}, \text{tmp_b}$

w = weight
 b = bias

Gradient Descent Intuition in Derivative

- Looking at $J(w)$:

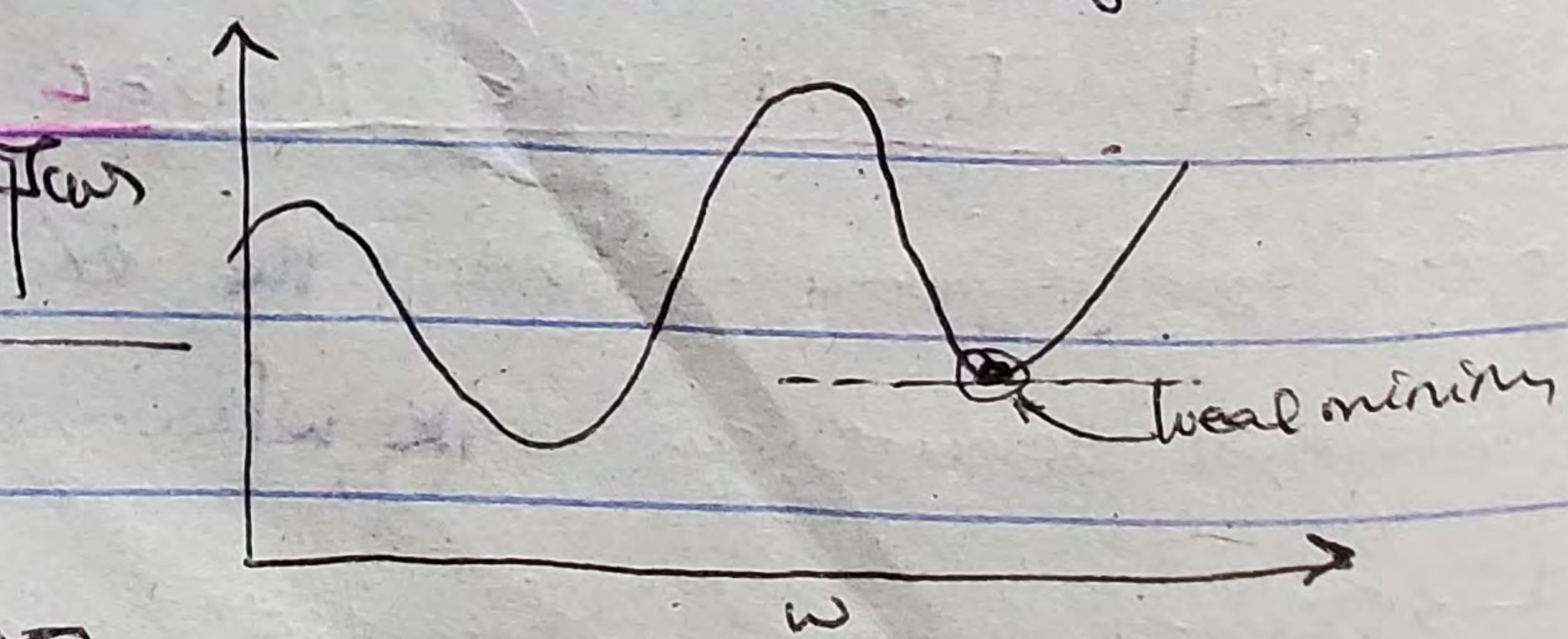
setty intercept $b = 0$



LEARNING RATE: α has a huge impact on the efficiency of gradient descent

- If α is too small; gradient descent may be slow. If α is too large; gradient descent may overshoot & fail to reach the minimum i.e. it may fail to converge instead it'll diverge.

What if you've already reached a local minimum?



Note: This is not a squared error cost function

and it has 2 local minima.

$$w = w - \alpha \frac{d}{dw} J(w)$$

The slope @ the local minimum is equal to 0

∴ $w = w - \alpha(0)$

$w = w$... no update

- As we approach a local minimum;
 - The derivative becomes smaller
 - \therefore The update steps become smaller
 - This means that we can reach minimum w/o decreasing learning rate α

GRADIENT DESCENT FOR LINEAR REGRESSION

- for our linear regression model we have; $f(x)$

$$f_{w,b}(x) = w_0 + b$$

our cost function, $J(w, b) = \frac{1}{2m} \sum_{i=1}^m [f_{w,b}(x^{(i)}) - y^{(i)}]^2$

- the gradient descent algorithm

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b) ; \text{ where } \frac{\partial}{\partial w} J(w, b) = \frac{1}{m} \sum_{i=1}^m [f_{w,b}(x^{(i)}) - y^{(i)}] x^{(i)}$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b) ; \text{ where } \frac{\partial}{\partial b} J(w, b) = \frac{1}{m} \sum_{i=1}^m [f_{w,b}(x^{(i)}) - y^{(i)}]$$

- The squared error cost function always has \geq global minimum & will never have $<$ multiple local minima. In other words, it is a convex function.

Convex Function: A bowl shaped function that can't have any local minimum other than the single global minimum.

- The gradient descent used = batch gradient descent. "Batch" because each step of the gradient descent uses all the training examples. $\Rightarrow m = \text{length}$

WEEK 2 MULTIPLE LINEAR REGRESSION

A: Multiple Features (Variables)

	Size-feet ² (x_1)	No. of bedrooms (x_2)	No. of floors (x_3)	Age in years (x_4)	Price(\$) in \$1000's
1	2104	5	1	45	460
2	1416	3	2	40	232
3	1534	3	2	30	315
4	852	2	1	36	178
...

$x_j^{(i)}$ = j^{th} feature; $n = \text{no. of features}$; $\vec{x}^{(i)}$ = feature of i^{th} training example

$x_j^{(i)}$ = value of feature j in i^{th} training example | e.g. $\vec{x}^{(2)} = [1416, 3, 2, 40]$

$$x_3^{(2)} = 2$$

$$j = 1 \dots 4$$

$$n = 4$$

New model; $f_{w,b}(x) = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b$

$$\rightarrow f_{\vec{w},b}(\vec{x}) = w_1x_1 + w_2x_2 + \dots + w_nx_n + b \quad \vec{w} = [w_1, w_2, \dots, w_n]$$

$$\Rightarrow f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

- Note this is not multivariate regression

B: VECTORIZATION \rightarrow Helps makes code shorter + more efficient

\hookrightarrow Make use of Numpy functions.

When implementing a learning algorithm

C: GRADIENT DESCENT for MULTIPLE LINEAR REGRESSION

Previous notation

Parameters w_1, \dots, w_n

b

Model $f_{\vec{w},b}(\vec{x}) = w_1x_1 + \dots + w_nx_n + b$

Cost function $J(w_1, \dots, w_n, b)$

Gradient descent repeat {

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(w_1, \dots, w_n, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w_1, \dots, w_n, b)$$

}

Vector notation

vector of length n

$$\vec{w} = [w_1, \dots, w_n]$$

b

$$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

$$J(\vec{w}, b)$$

repeat {

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$$

| Normal equation \rightarrow An alternative to gradient descent.

| Only for linear regression & solves for w, b without iterations.

Disadvantages:

- (i) Doesn't generalize to other learning algorithms
- (ii) Slow when no. of features n large (> 10000)

| What you need to know:

- Normal eqn method may be used w/ ML libraries to implement linear regression & not by yourself.
- Gradient descent is a better way

Wk2 Gradient Descent in Practice

A. FEATURE SCALING: When the the possible range of values of a feature is large, it's more likely that a good model would learn to choose a relatively small parameter value. Likewise when the possible values of a feature are small, the reasonable value for the parameters will be large.

How then does this relate to Gradient Descent?

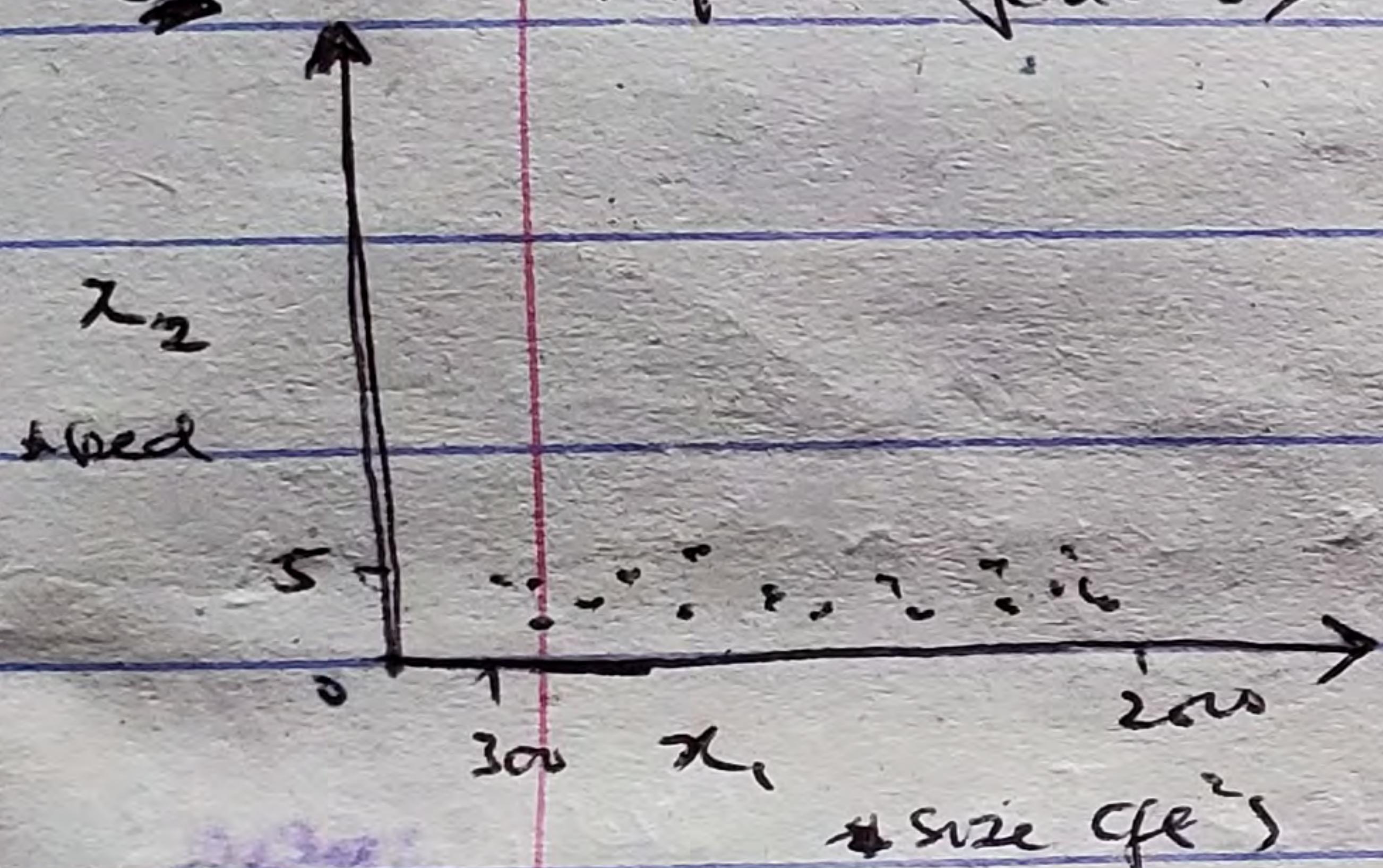
E.g.: price = $w_1x_1 + w_2x_2 + b$: x_1 : size (feet²) ; x_2 : # of bedrooms
 Range = 300 - 2000 Range = 0 - 5

possible

- size & parameters for a house w/ $x_1 = 2000$, $x_2 = 5$ & price = \$500k

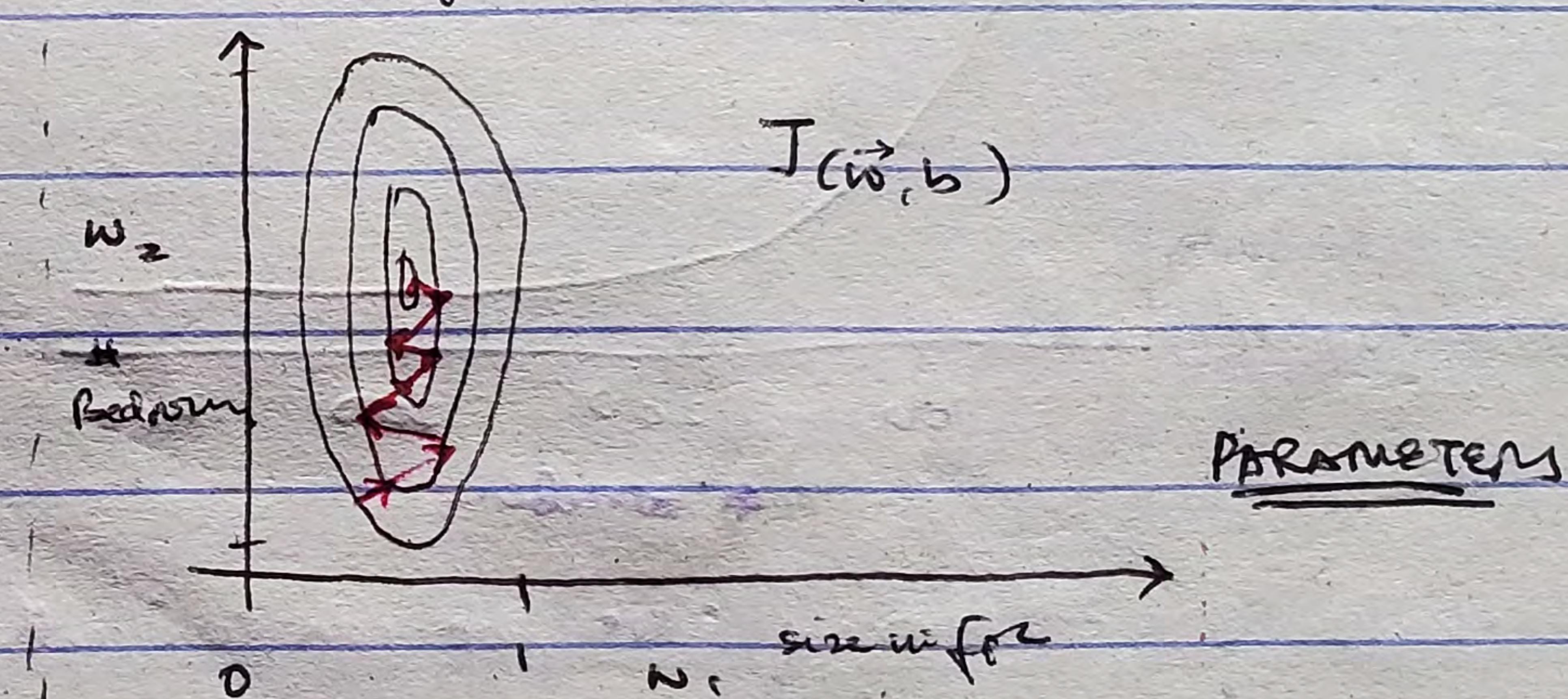
$$w_1 = 0.1; w_2 = 50, b = 50 \rightarrow \text{price} = 0.1(2000) + 50(5) + 50 = \$500k$$

\rightarrow <scatterplot & features>



Features

contour plot
cost function J(w, b)

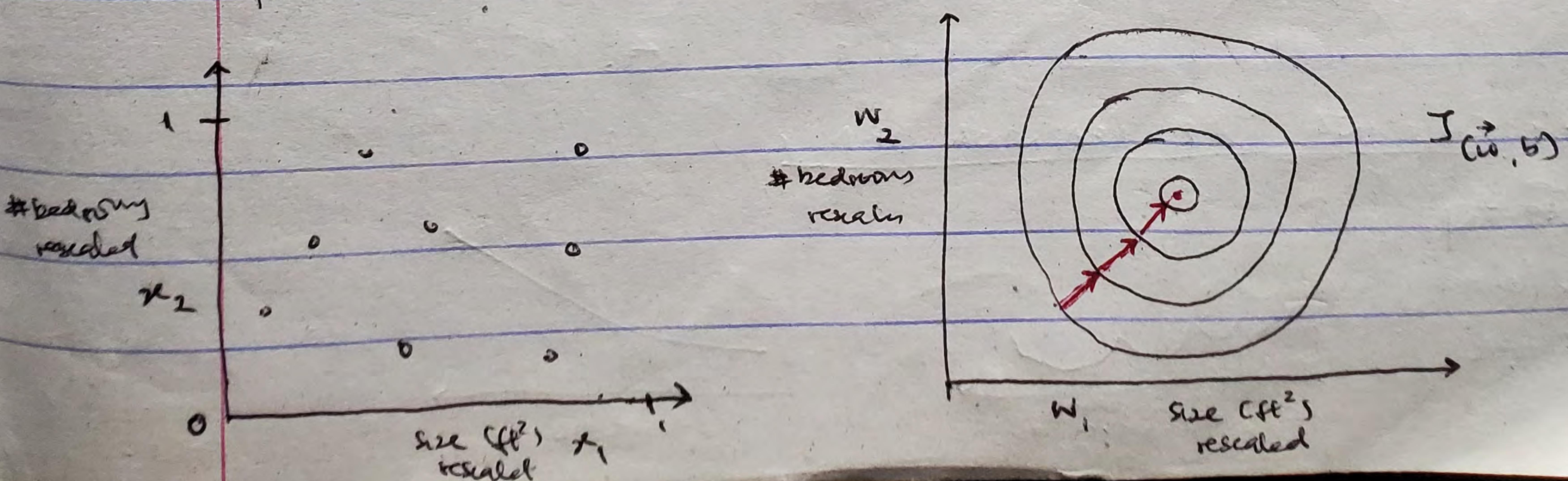


Horizontal axis has a much narrower range, while the vertical axis has a larger range. This is a very small change to w_1 can have a very large impact on the estimated price & thus a large impact on the cost J . The opp. w/ w_2

\rightarrow Because the contours are so tall & skinny, gradient descent might end up

\rightarrow bouncing back & forth for a long time before ~~reaching~~^{settling} on the global minimum. In order to avoid this, the features need to be scaled i.e. transform the data so that x_2 range from ~0 to ~1

the, the features need to be scaled i.e. transform the data so that x_2 range from ~0 to ~1



Summary: Rescaling features speed up gradient descent significantly

Scales matters

- Mean normalization $\rightarrow 300 \leq x_i \leq 2000$

$$x_i = \frac{x_i - \bar{x}_i}{\text{max} - \text{min}} \quad \begin{matrix} \nearrow \text{mean } \bar{x}_i \\ \searrow \text{max-min} \end{matrix}$$

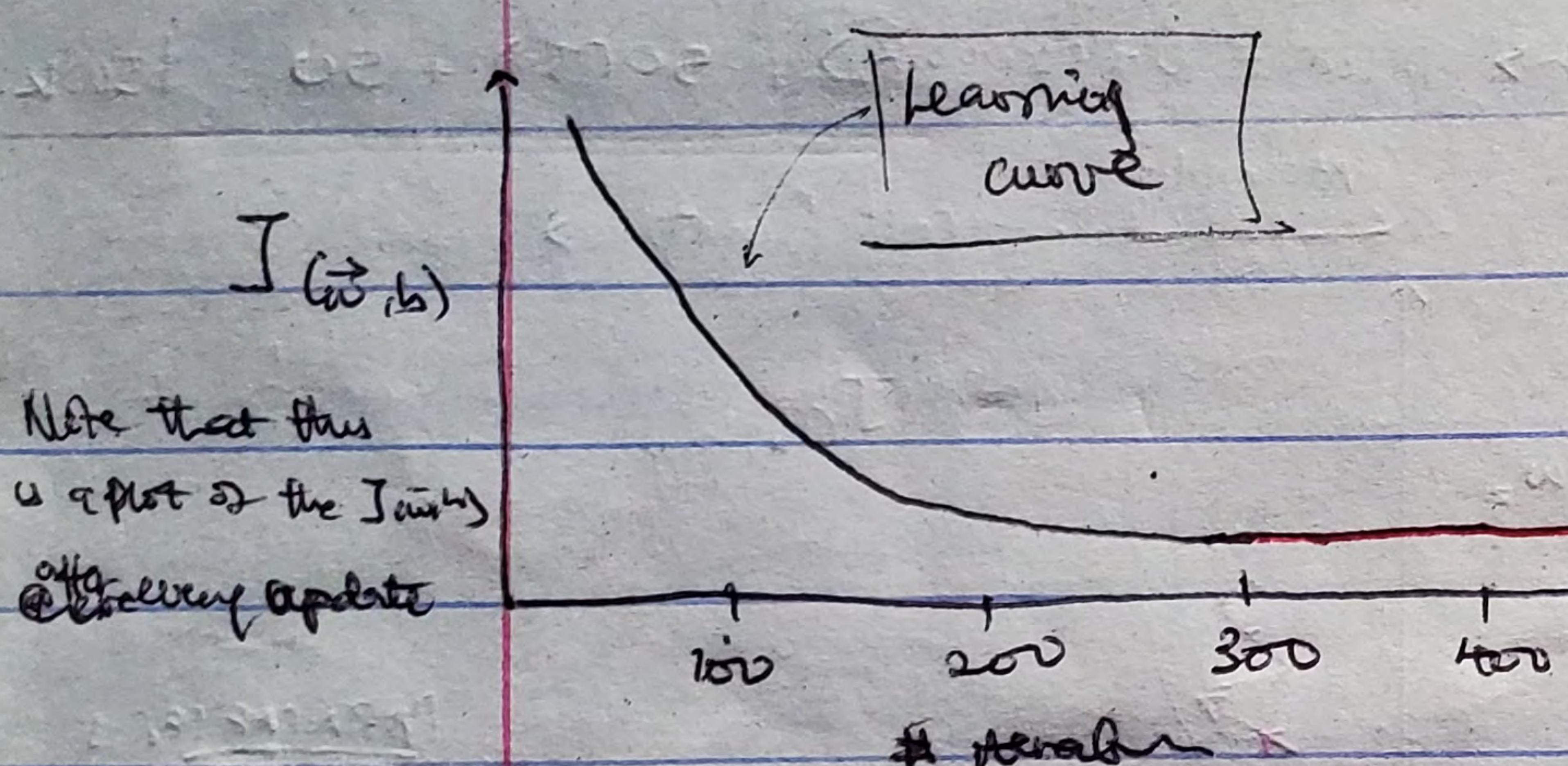
More come $\rightarrow Z$ -score normalization = $\frac{x_i - \mu_i}{\sigma_i}$

Ans for Feature scaling :

$-1 \leq x_i \leq 1$ $-3 \leq x_i \leq 3$ $-0.3 \leq x_i \leq 0.3$	} acceptable range for each feature
--	--

B. MONITORING GRADIENT DESCENT FOR CONVERGENCE

The objective of gradient descent is to minimize the cost: $\min_{\vec{w}, b} J(\vec{w}, b)$



The graph helps us visualize how $J(\vec{w}, b)$ changes after every single iterat

- If it ever increases, it either means that α is too large or there's a bug in the cost

- Looking at the learning curve, we can confidently say that $J(\vec{w}, b)$ likely converged by 400 iteration because the curve is no longer decreasing.

- Note that this # iteration needed to converge varies

Although not as reliable as the learning curve, we can use the automatic convergence

Although not a test. We set ϵ (epsilon) to be a number: 10^{-3} or 0.001. We say that

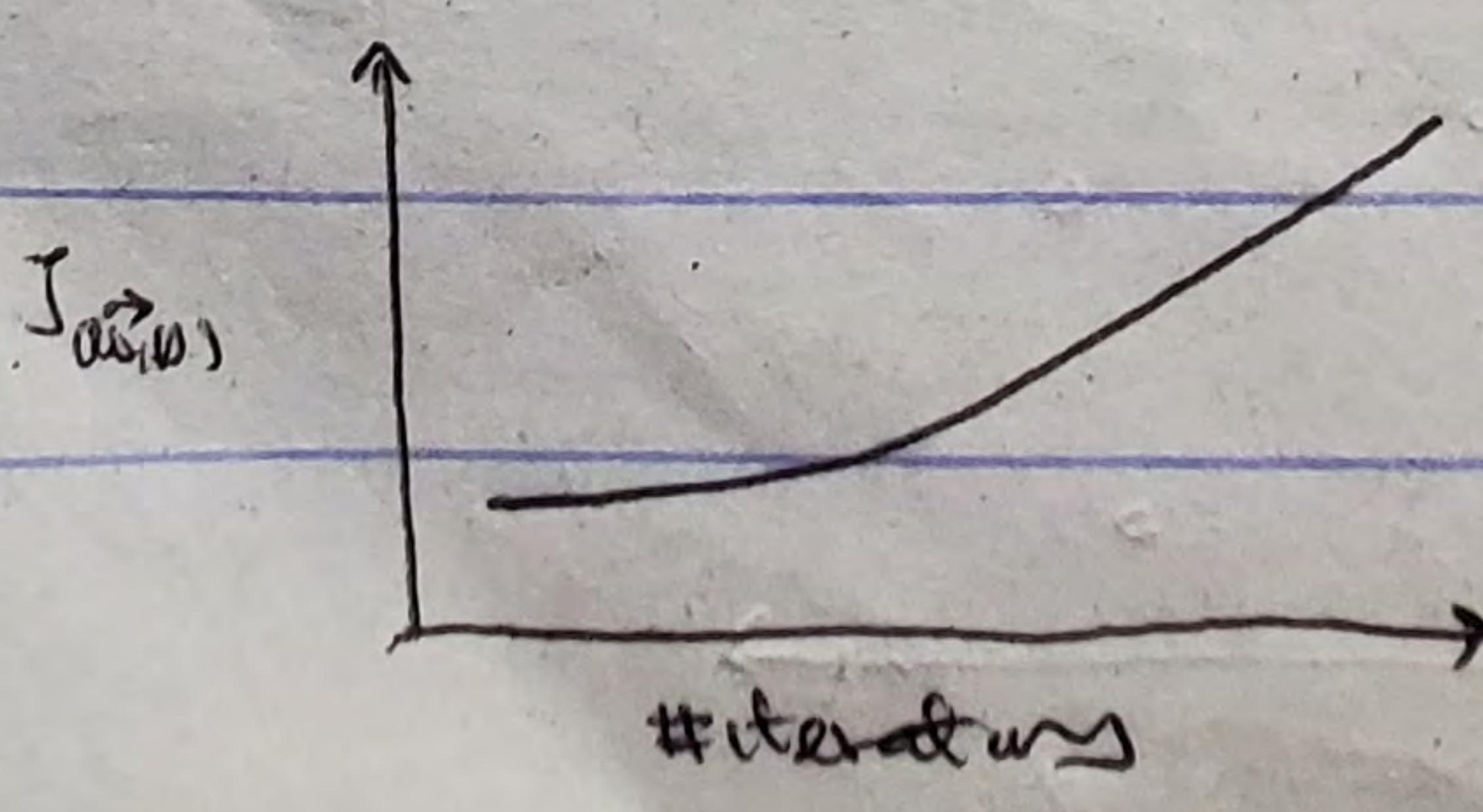
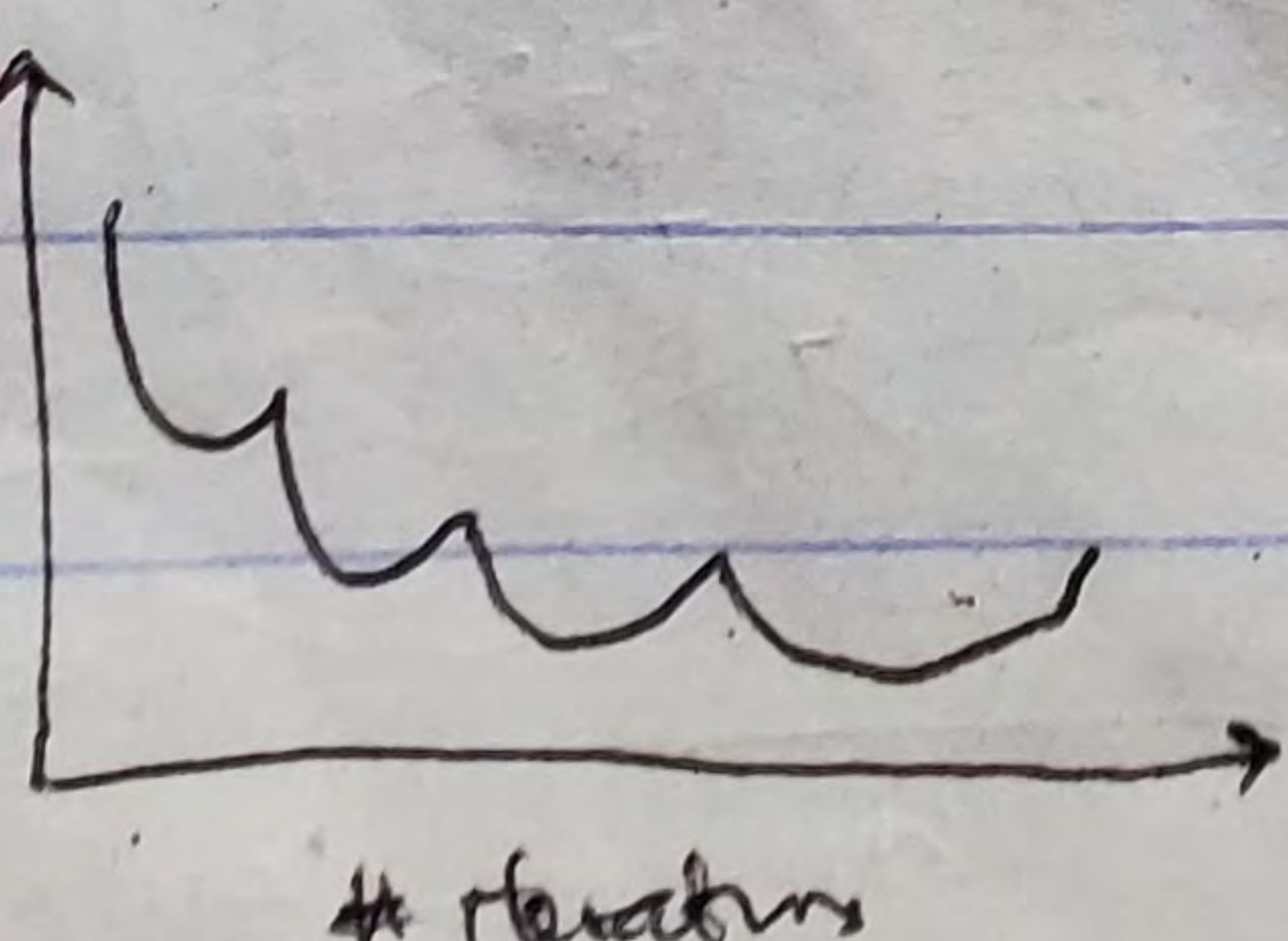
if $J(\vec{w}, b)$ decreases by $\leq \epsilon$ in one iteration, we can declare convergence

CONVERGENCE = We've found parameters \vec{w}, b that give us the global minima

where our cost, $J(\vec{w}, b)$, is least.

C. CHOOSING THE LEARNING RATE

- Faulty grade learning curves



most likely a long

- Both can be solved by choosing a smaller learning rate.

No Debugging | → When in error, set α to a very low value & see if the learning curve decreases. If it doesn't, then it's a bug.

Tip

$\rightarrow \times 3$

Use values for $\alpha = [0.001, 0.01, 0.1, 1]; [0.003, 0.03, 0.3, 3]$

D: FEATURE ENGINEERING

→ This involves using your knowledge or intuition about the problem to design new features, usually by transforming or combining original features.

Example: Suppose you have y ^{a target} of price of a house with features, x_1 : frontage of the house (ft) & x_2 : depth of the house i.e. basically length & width. You might decide that area would be a useful feature and have $x_3 = x_1 x_2$ (ft²) → Your function would be; $f(x) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$

E: POLYNOMIAL REGRESSION

Wk 3 CLASSIFICATION WITH LOGISTIC REGRESSION

• Binary classification $\in 2$ classes or categories.

Example: Question

Is the email spam?

Is the transaction fraudulent?

Is the tumor malignant?

Answer "y"

No Yes

No Yes

No Yes

False True

$0 \text{ or } 1$ for model
↑
the class

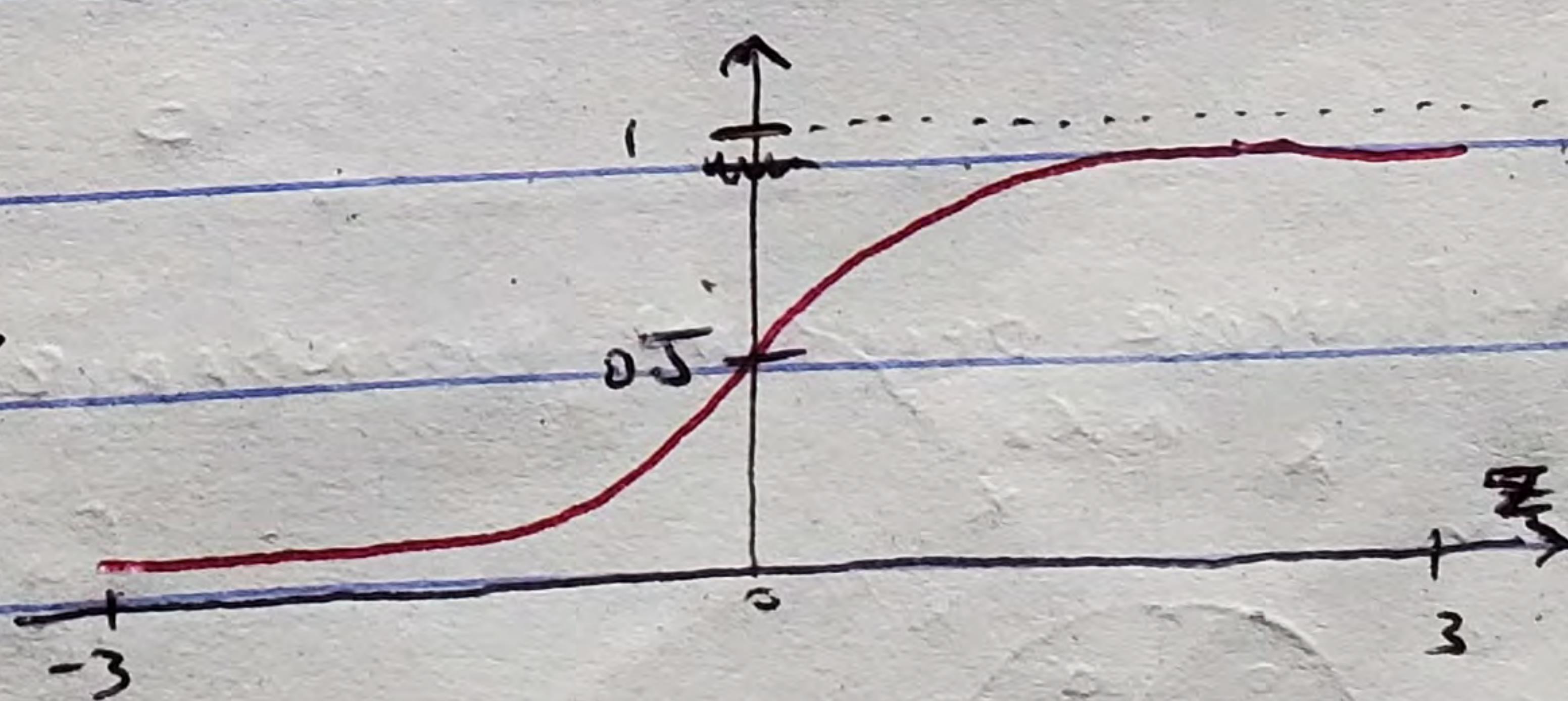
How do you build a classification model?

F: LOGISTIC REGRESSION

Sigmoid or logistic function \rightarrow

Outputs values between 0 & 1

$$g(z) = \frac{1}{1+e^{-z}}, 0 < g(z) < 1$$



$$\textcircled{a} \quad f_{w,b}(z) \quad | z = w \cdot x + b$$

$$\textcircled{b} \quad f_{w,b}(z) = \frac{1}{1+e^{-(w \cdot x + b)}}$$

$$\textcircled{c} \quad g(z) = \frac{1}{1+e^{-z}}$$

INTERPREATION OF LOGISTIC REGRESSION OUTPUT

$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-\vec{w} \cdot \vec{x} + b}}$$

} think of the output as the probability that the class or "y" is 1 given a certain input, x

Examp: $x = \text{"tumor size"}; y = \begin{cases} 0 & (\text{not malignant}) \\ 1 & (\text{malignant}) \end{cases}$

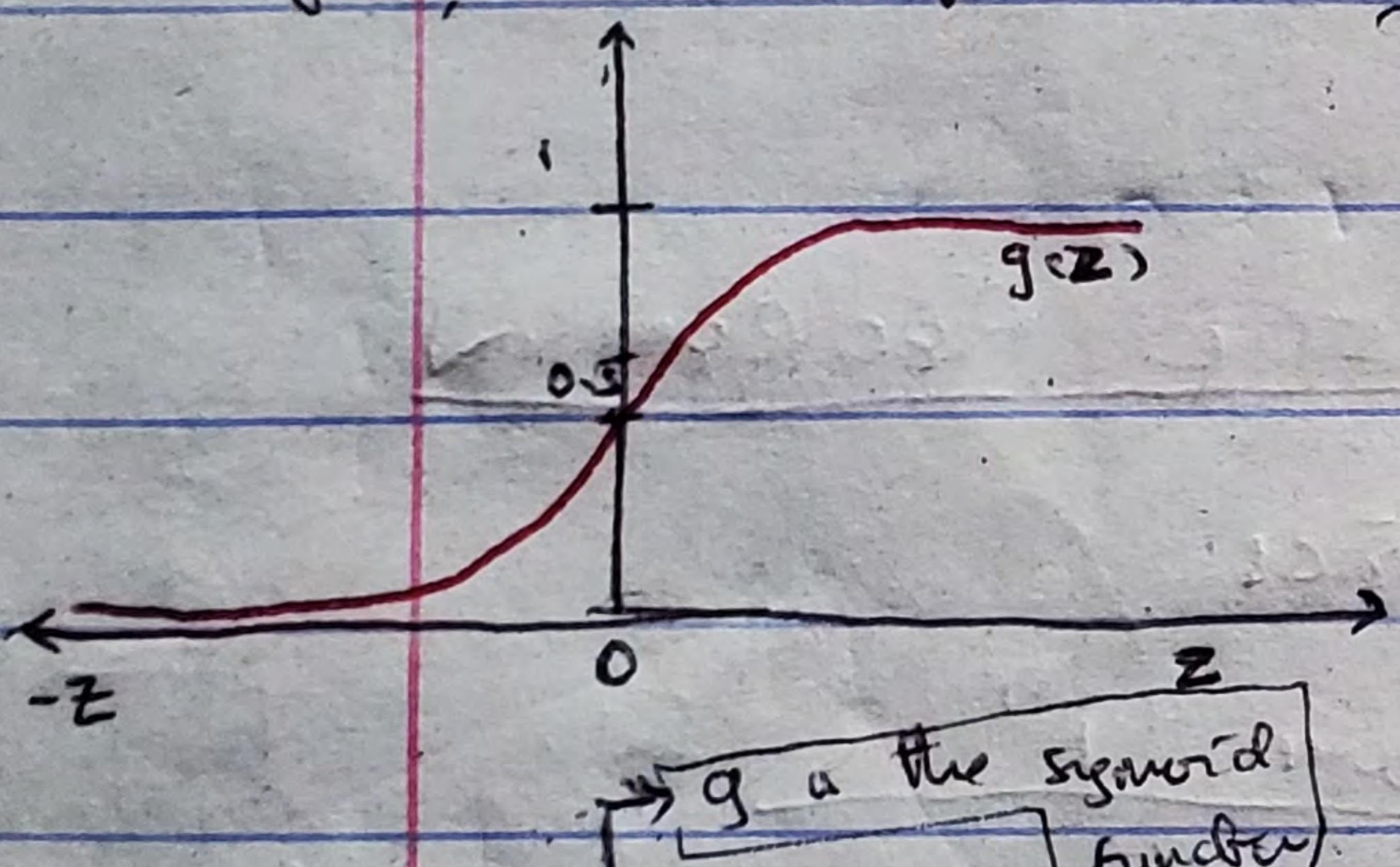
If $f_{\vec{w}, b}(\vec{x}) = 0.7$, then the model thinks that the patient has a 70% chance of their tumor turning out to be malignant.

Note when reading research papers, $\rightarrow f_{\vec{w}, b}(\vec{x}) = P(y=1 | \vec{x}; \vec{w}, b)$; semicolon denotes what the regression params are.

You might see $\boxed{\text{the mean}} \rightarrow f_{\vec{w}, b}(\vec{x}) = P(y=1 | \vec{x}; \vec{w}, b)$; the probability that $y=1$, given the input \vec{x} with parameters \vec{w}, b

B: DECISION BOUNDARY: Choosing what threshold defines whether

For example, take the sigmoid plot $\hat{y} = 1 \text{ or } 0$ based on $f_{\vec{w}, b}(\vec{x})$



Suppose $\circ f_{\vec{w}, b}(\vec{x}) \geq 0.5$, let $\hat{y} = 1$

then when $\circ \hat{y} = 1 \rightarrow g(z) \geq 0.5$

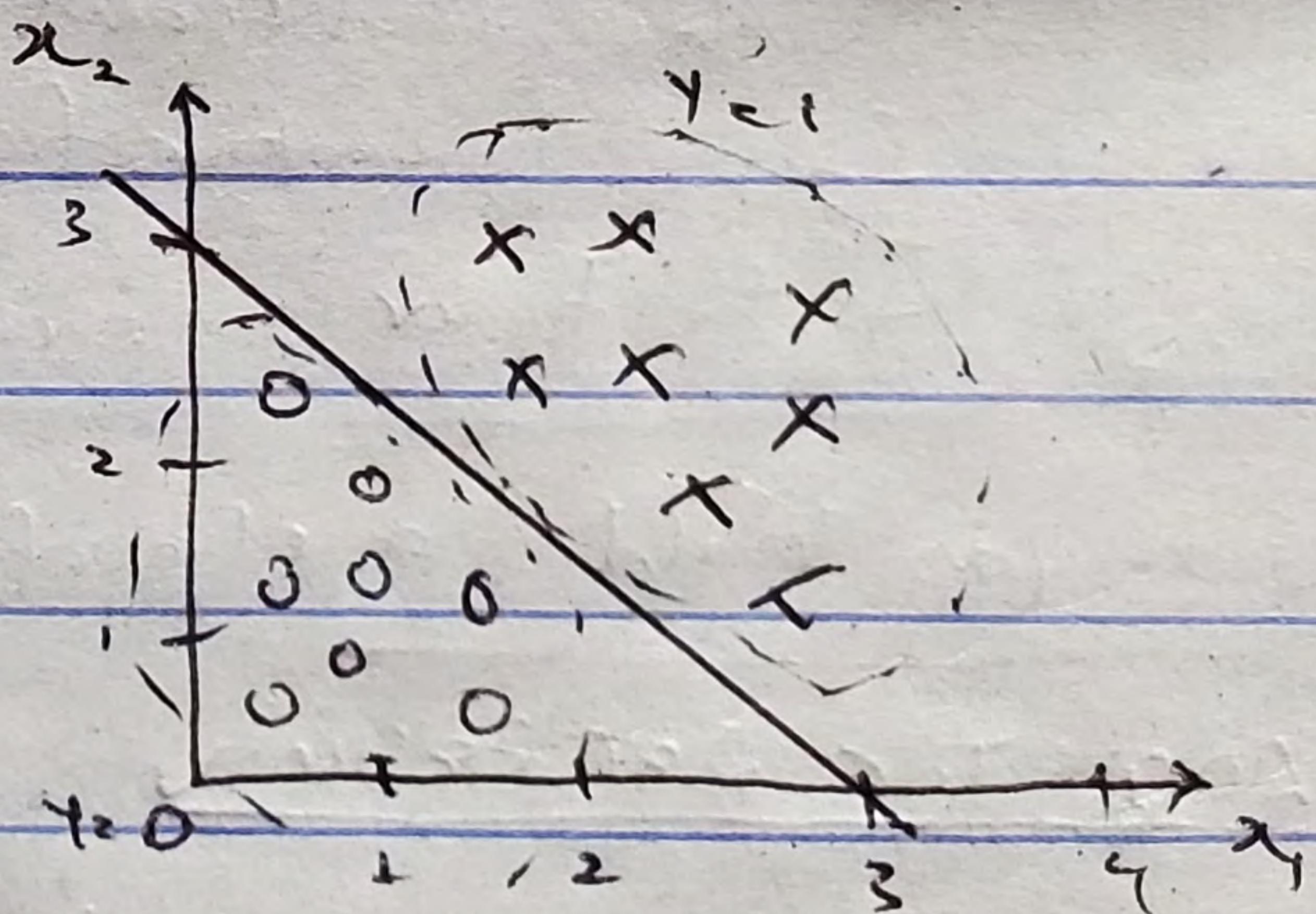
$$\begin{aligned} z &\geq 0 \\ \vec{w} \cdot \vec{x} + b &\geq 0 \approx \vec{w} \cdot \vec{x} + b < 0 \\ \hat{y} &= 1 \quad \hat{y} = 0 \end{aligned}$$

$$f_{\vec{w}, b}(\vec{x}) \geq g(z) = g(w_1 x_1 + w_2 x_2 + b)$$

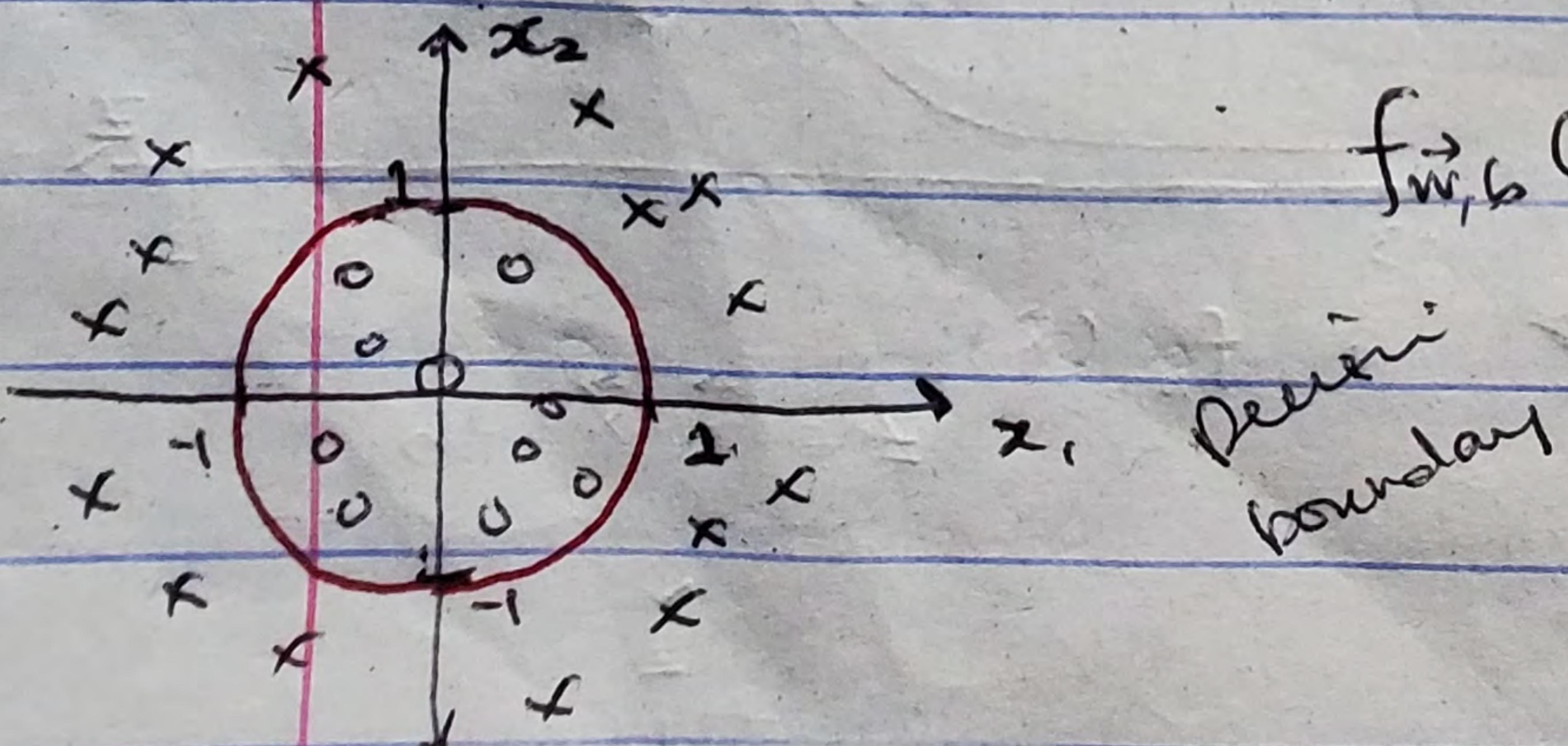
$$z = \vec{w} \cdot \vec{x} + b \geq 0$$

$$\text{Decision boundary: } z = x_1 + x_2 - 3 = 0$$

$$x_1 + x_2 = 3$$



NON-LINEAR DECISION BOUNDARIES



$$f_{\vec{w}, b}(\vec{x}) = g(z) = g(w_1 x_1^2 + w_2 x_2^2 + b)$$

$$z = x_1^2 + x_2^2 - 1 = 0$$

$$x_1^2 + x_2^2 = 1$$

WK3 Cost Function for Logistic Regression

Training set:

turner size cm's x_i	Patient's age x_n	malignant? y	$i = 1, \dots, m$
10	52	1	$j = 1, \dots, n$ - features
2	73	0	target $y = 0$ or 1
5	55	0	
12	49	1	
...	

$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$

Research Question: Given the training set, how can you choose

parameters $\vec{w} = [w_1, w_2, \dots, w_n]$ and b ?

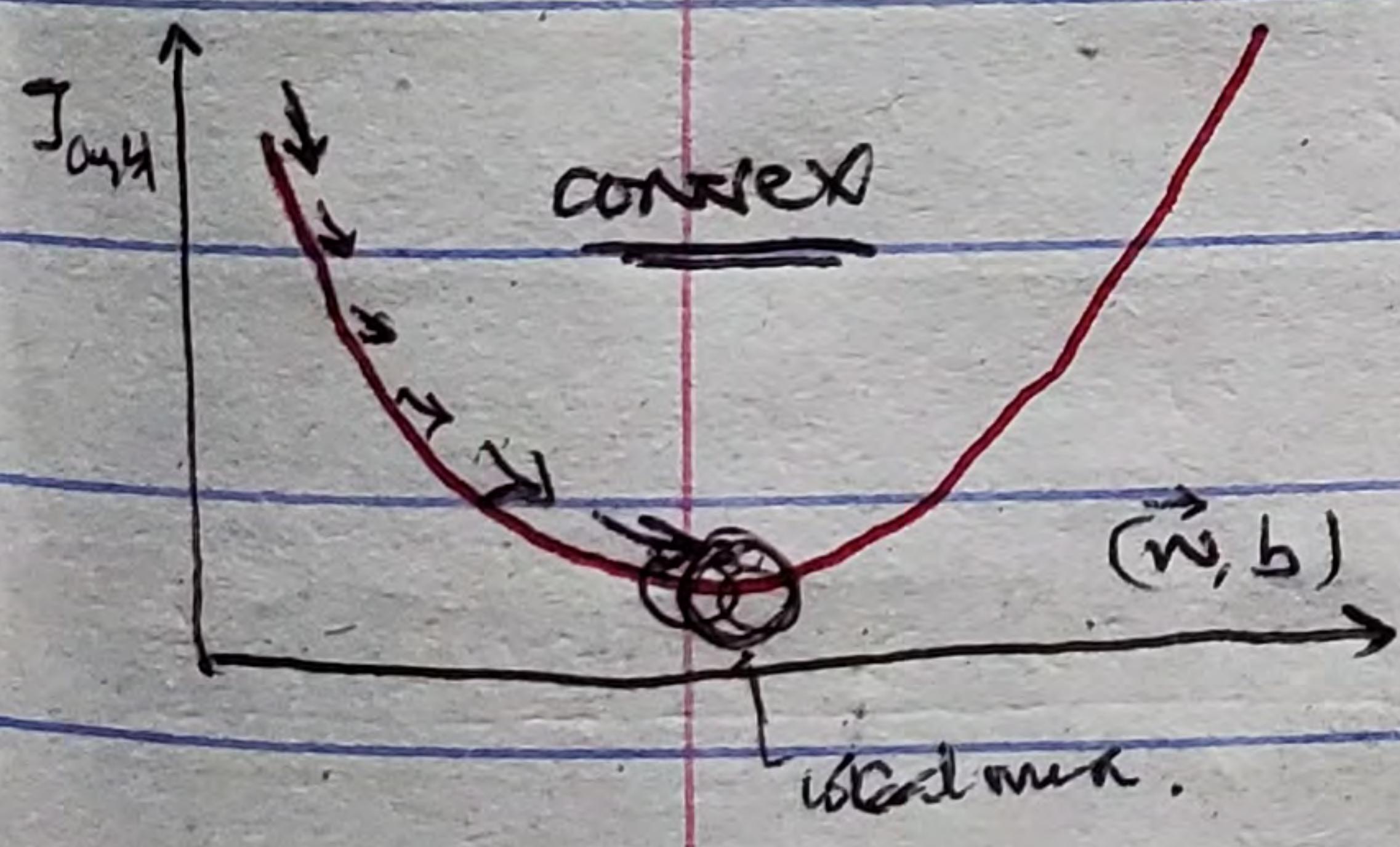
recall that for linear regression, we used the SQUARED ERROR COST FUNCTION

$$J_{(\vec{w}, b)} = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} [f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}]^2$$

If we were to use gradient descent to find the true minima & plot the cost function for each value of \vec{w}, b (considering b is constant); we'd get

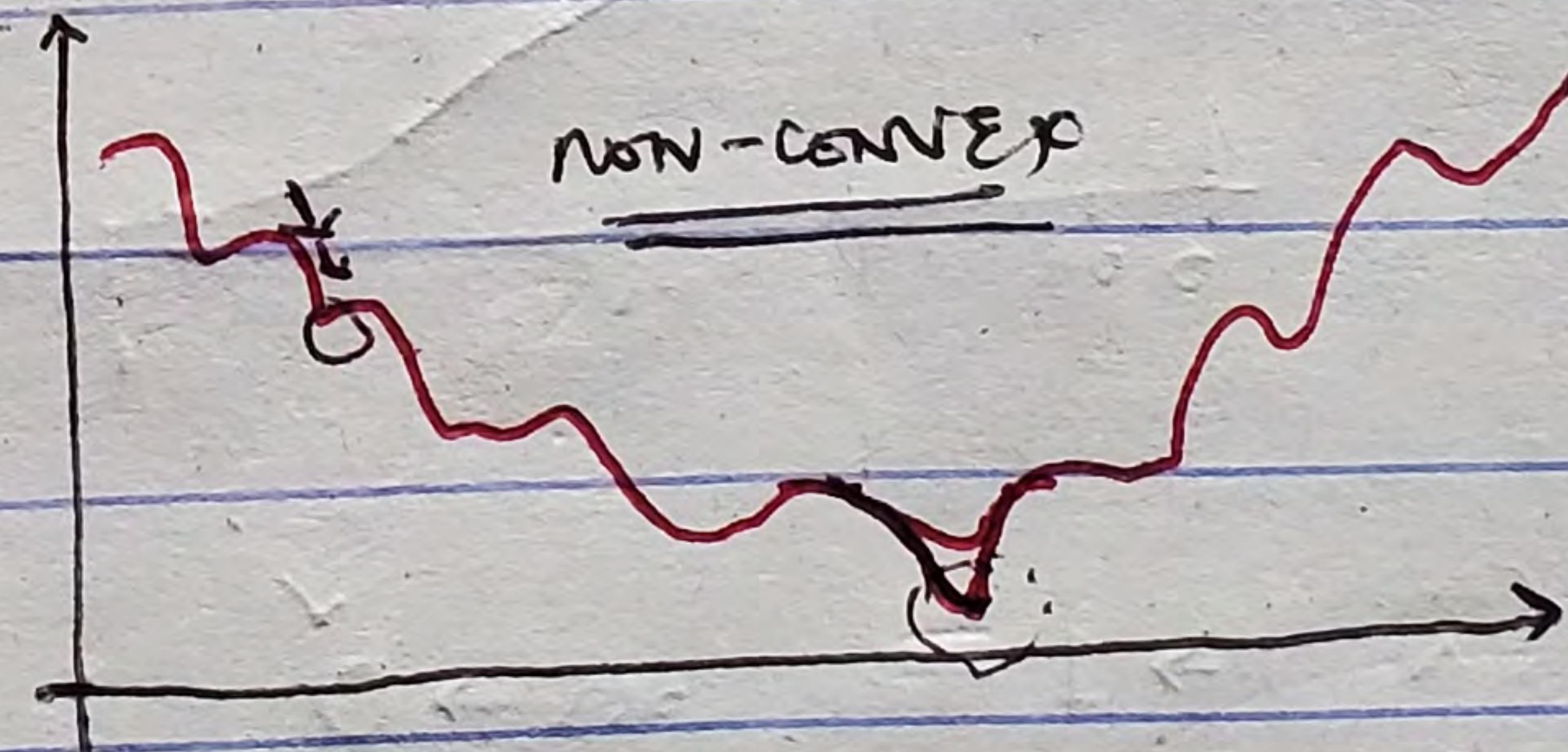
linear regression

$$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$



Logistic regression

$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$



Using Squared error cost function results in a non-convex shape & it gets lots of false local minima & prevents gradient descent from reaching the true minimum.

∴ Squared error cost function isn't good for logistic regression

To get our new cost function: Let's first look @ the Squared Cost function

↪ change its definition a bit.

Looking @ the expression; $J_{(\vec{w}, b)} = \frac{1}{m} \sum_{i=1}^m \left[\frac{1}{2} [f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}]^2 \right]$

Made the summation

Loss; $L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})$

that is the loss
on a single training
example

denoted by

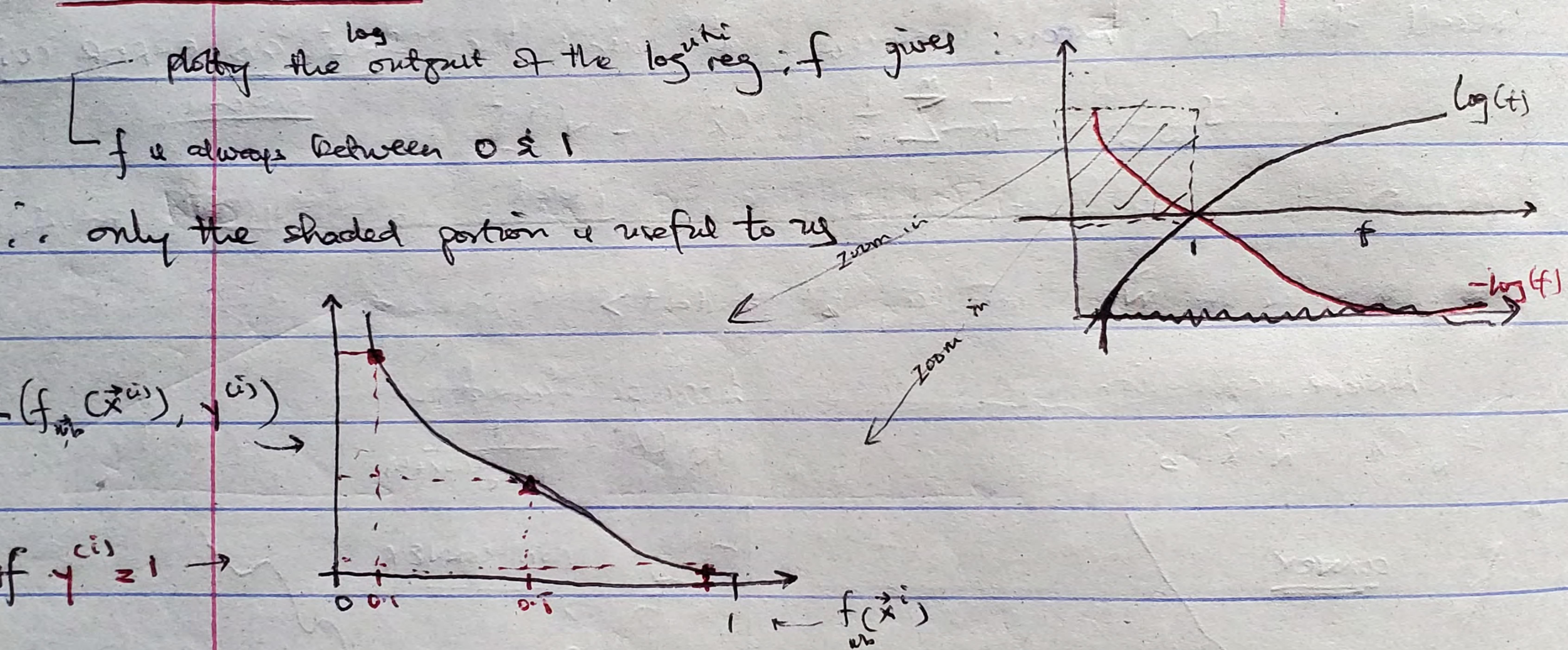
Definition of loss, L ; $L(f_{w,b}(\vec{x}^{(i)}), y^{(i)})$ is a function of the prediction of the learning algorithm $f_{w,b}(\vec{x})$ as well as of the true label $y^{(i)}$

- In squared error cost function; the loss given the predictor $f_{w,b}(\vec{x})$ and the true predictor y is equal to $\frac{1}{2}$ of the squared difference.

- Choosing a different form for the loss function will help us keep the overall cost function = $\frac{1}{m} \times$ sum of loss functions, to be a convex function.

- Recall the loss function applies forces \vec{y}' and tells us how well we're doing on that example -

$$\text{Logistic Loss function: } L(f_{w,b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{w,b}(\vec{x}^{(i)})) & \text{for } y=1 \\ -\log(1-f_{w,b}(\vec{x}^{(i)})) & \text{for } y=0 \end{cases}$$



- A) As $f_{w,b}(\vec{x}^{(i)}) \rightarrow 1$ then loss $\rightarrow 0$ \rightarrow If the algorithm predicts a probability close to 1 & the true label, as $y^{(i)} = 1$, then the loss ≈ 0 .
- If the output, $f_{w,b} = 0.5$ when $y^{(i)} = 1$ then the loss is a bit higher.

(If the output, $f_{w,b} = 0.1$, then $y^{(i)} = 1$, then loss is really high)

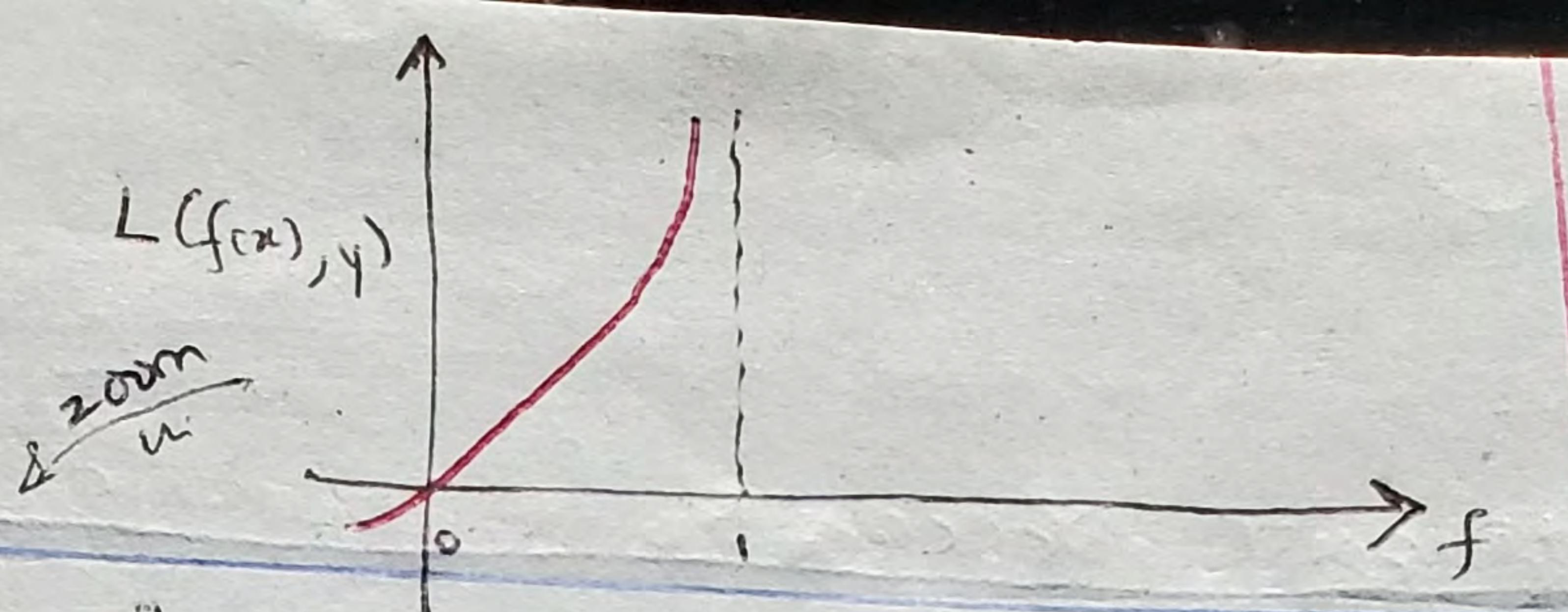
- B) As $f_{w,b}(\vec{x}^{(i)}) \rightarrow 0$ then loss $\rightarrow \infty$

For $y^{(i)} = 1$ - the loss is lowest when $f_{w,b}(\vec{x}^{(i)})$ predicts close to true label $y^{(i)} = 1$
Same for $y^{(i)} = 0$

* Now let's look @ the second part of the loss func

$$-\log[1 - f_{w,b}(\vec{x}^{(i)})] \text{ if } y^{(i)} = 0$$

when the function is plotted



When $f(x) \neq 0$, then the loss, $L \geq 0$. The larger the $|f(x)|$ gets, the bigger the loss. As the prediction, $f(x)$ approaches 1, the loss approaches ∞

As $f_{\vec{w}, b}(\vec{x}^{(i)}) \rightarrow 1$ then $\text{loss} \rightarrow \infty$ (bad boy in)

\therefore for $y^{(i)} = 0$, the further the prediction $f_{\vec{w}, b}(\vec{x}^{(i)})$ is from the ~~target~~ $y^{(i)}$, the higher the loss

⇒ From our loss functions, we can see that we're ~~encouraging~~ ^{incentivizing} our algorithm to predict the right value by giving it a lower loss & punishing it for predicting a false value by giving it a higher loss

⇒ The Logistic loss function gives us a convex shape when we plot the cost function $J(\vec{w}, b)$. Note that the cost, $J(\vec{w}, b)$ is the average of the loss for all the training examples.

Simplified Cost Function for Logistic Regression

Recall our initial loss func, $L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$

→ it can be re-represented as;

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log[f_{\vec{w}, b}(\vec{x}^{(i)})] - (1-y^{(i)}) \log[1 - f_{\vec{w}, b}(\vec{x}^{(i)})]$$

This re-representation is the same as above because $y^{(i)}$ can only take either 1 or 0

$$\therefore \text{ex, } J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m [L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})]$$

$$= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1-y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))]$$

* Note we chose this cost function using Maximum Likelihood Estimation which is an idea from statistics on how to efficiently find parameters for models.

WK3 GRADIENT DESCENT for Logistic Regression,

Cost functn : $J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) \right]$

repeat

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$$

$$\frac{\partial}{\partial w_j} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m [f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}] x_j^{(i)}$$

$$\frac{\partial}{\partial b} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m [f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}]$$

Although the algorithm written looks the same for Linear Regress & Logistic Regress

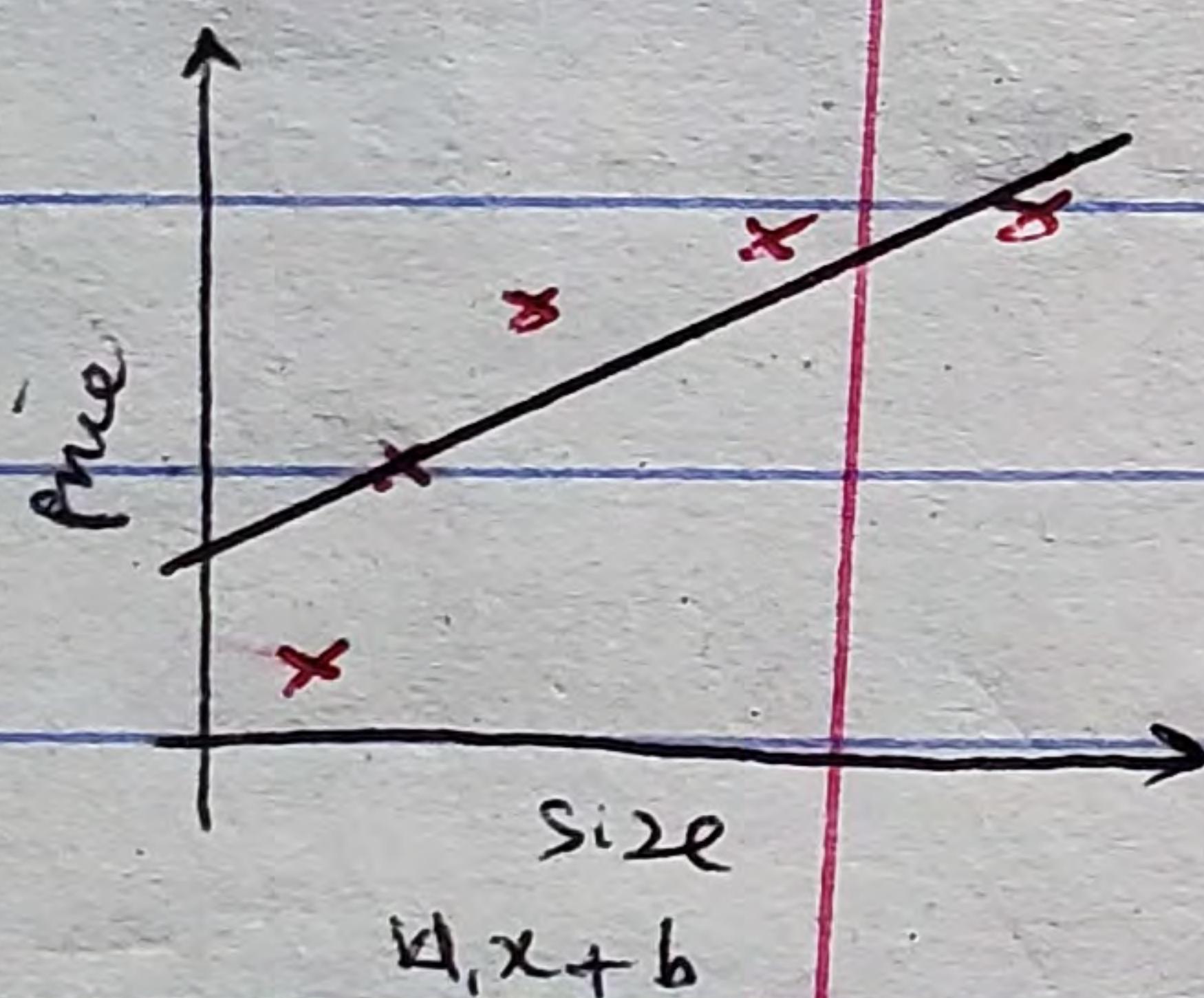
recall that

Linear regressn, $f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$

Logistic regressn, $f_{\vec{w}, b}(\vec{x}) = g(\vec{w} \cdot \vec{x}) \text{ or } g(z) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$

- Monitor gradient descent using a learning curve.

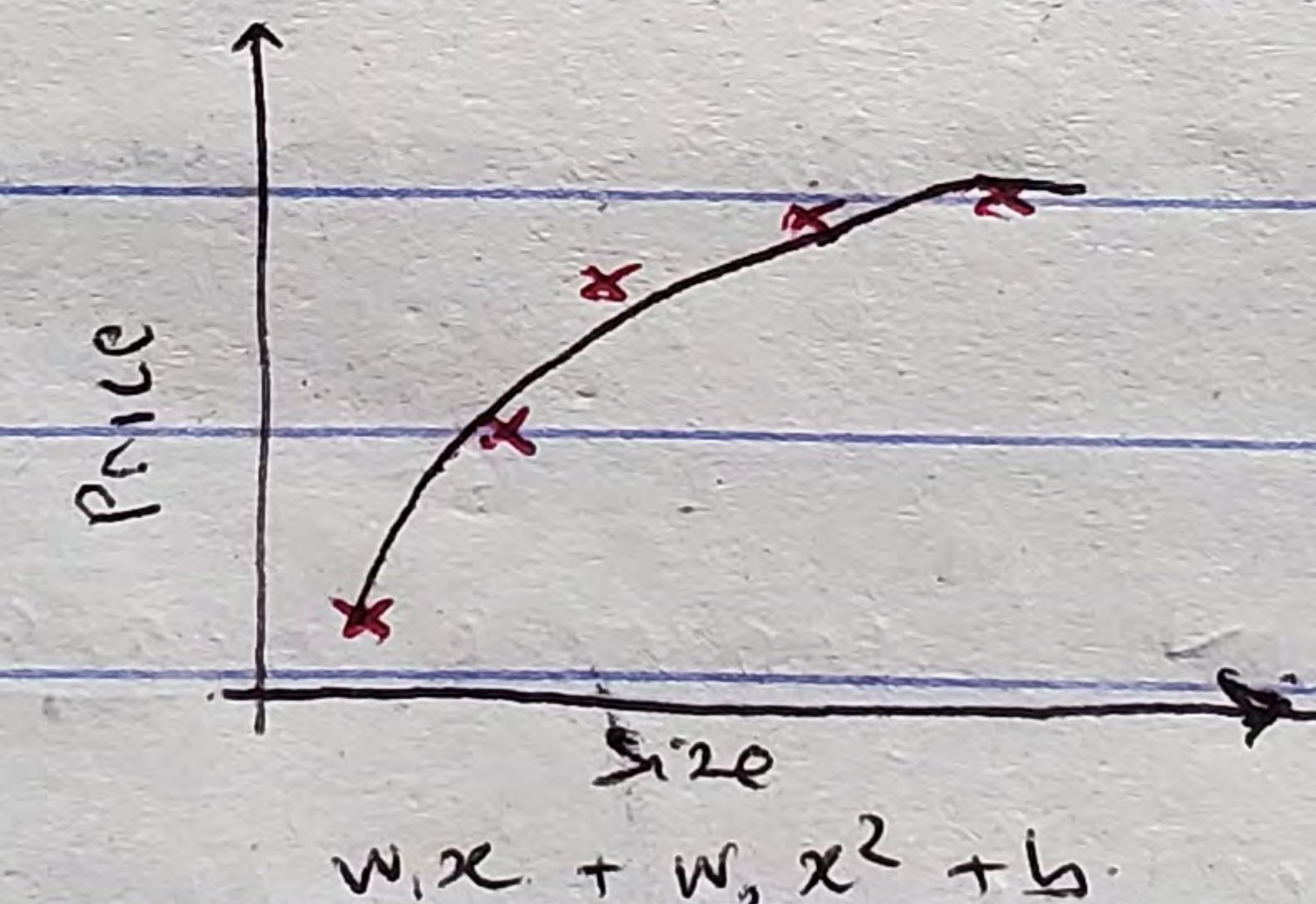
WK3 THE PROBLEM OF OVERFITTING



underfit
or high bias

Doesn't fit
the transeet
well

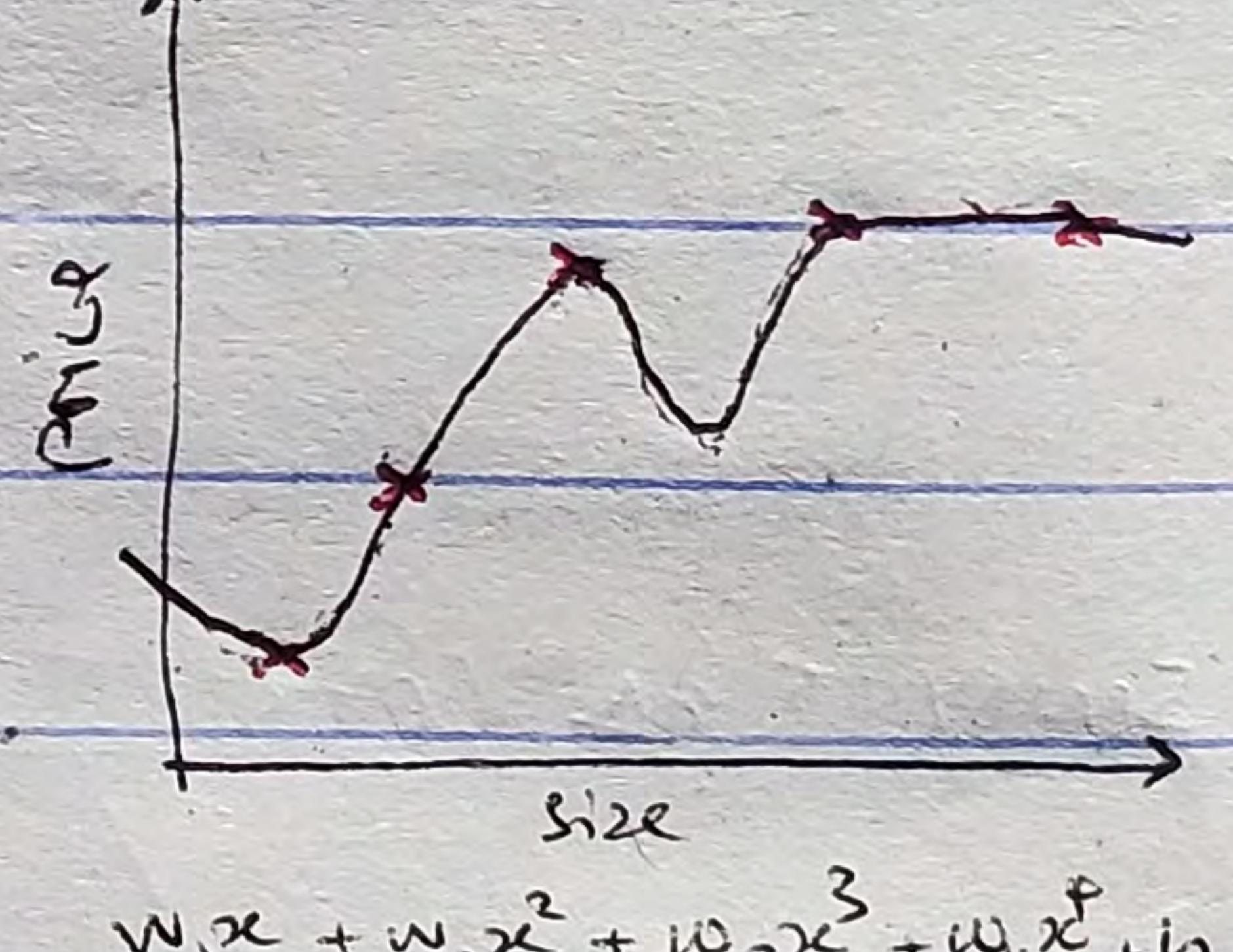
the model has a preconceived notion of
what the relationship should be despite
the data showing otherwise



- Fits training set pretty well

Model is **generalized**, i.e. it makes
good predictions even on brand
new examples it hasn't seen before

This is **GENERALIZATION**



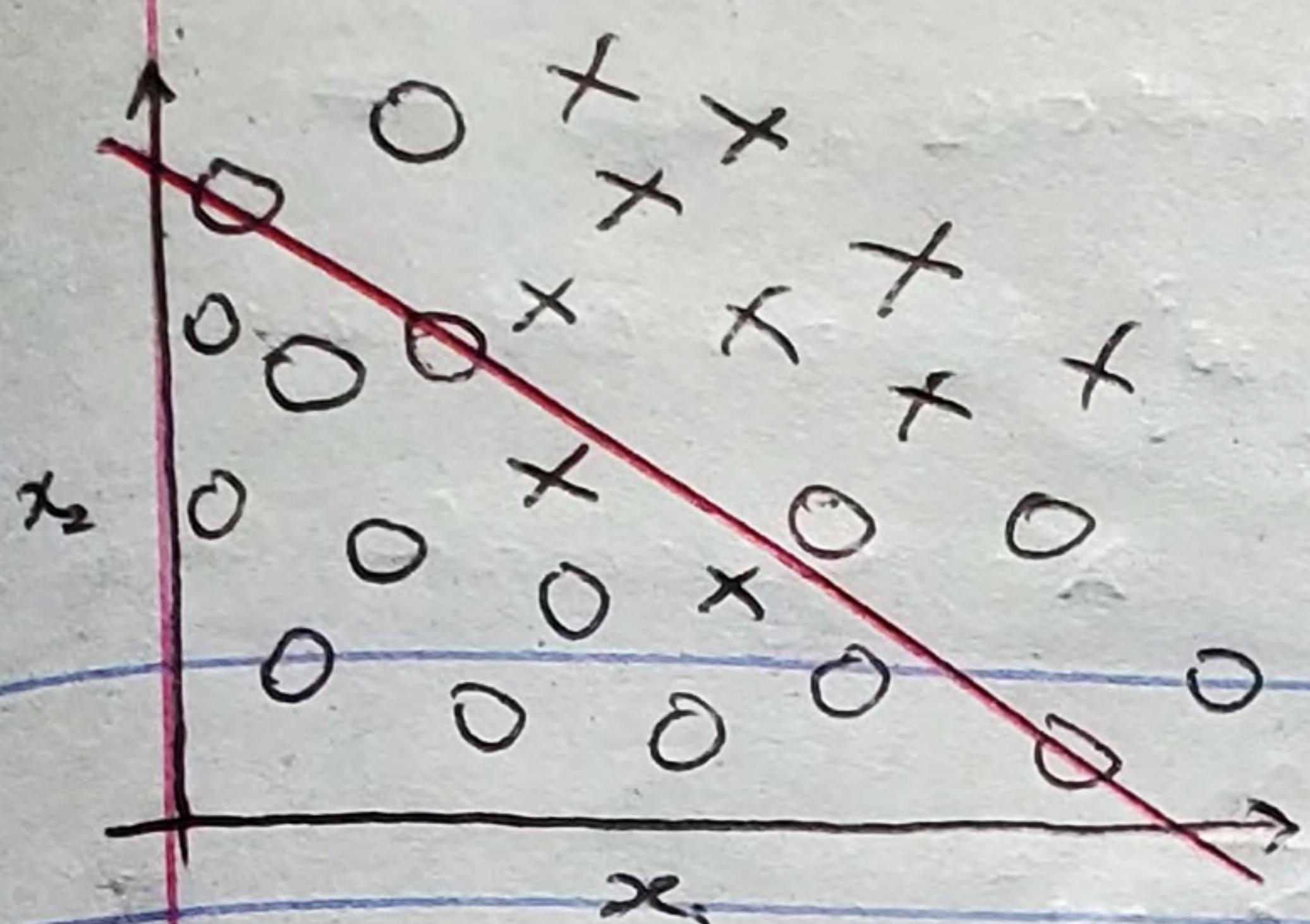
Tits the training set extremely
well

It fits the model too well that

it isn't generalized

OVERFITTING or
HIGH VARIANCE

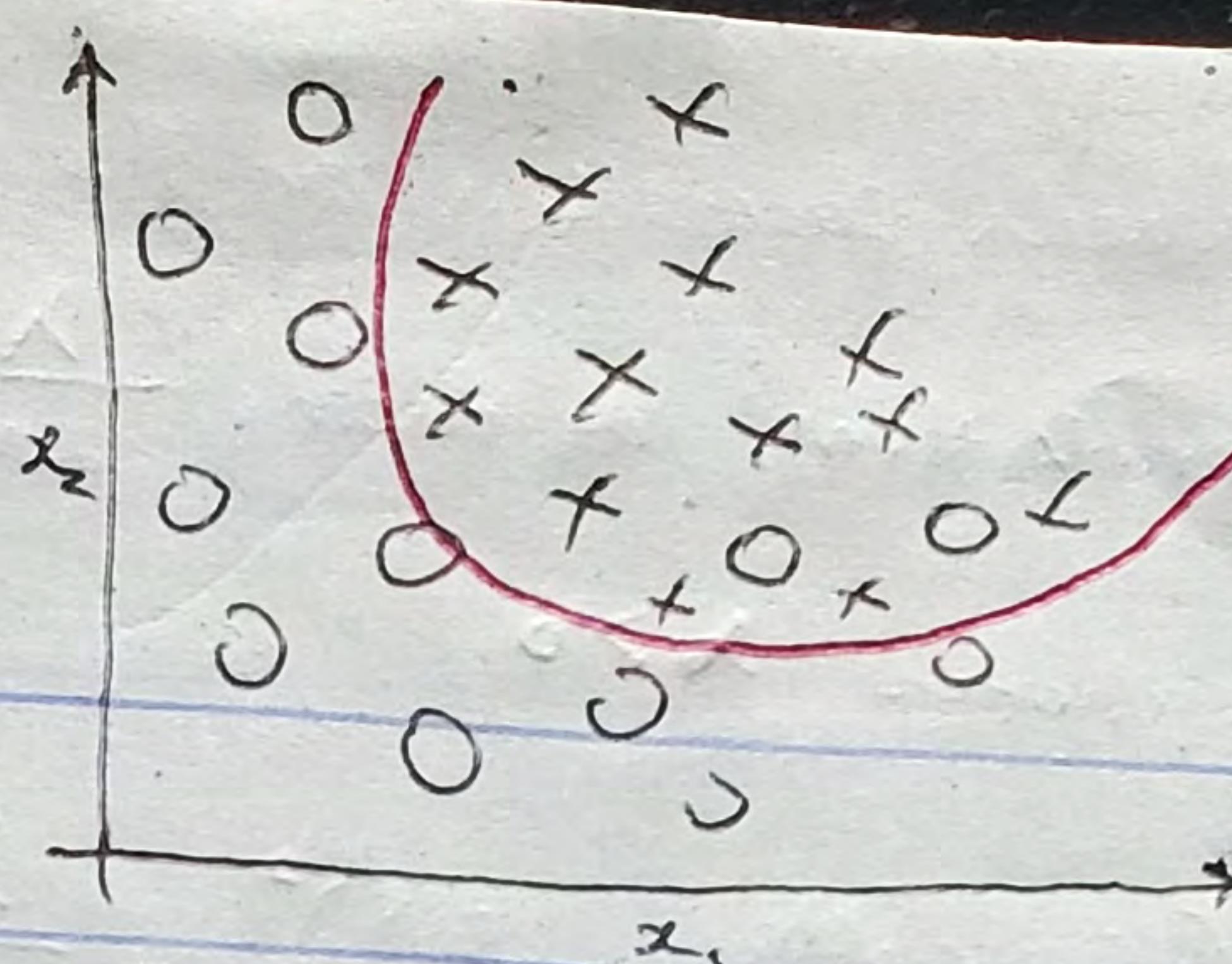
- In ML, people use underfit & highbias interchangeably & overfit & high variance interchangeably.
- The in overfitting, the algorithm is trying very hard to fit every example and it turns out that if the training set was a little different, the function that the algorithm fits would be totally different instead that results in highly variable predictions
 - the same applies in classification



$$z = w_1 x_1 + w_2 x_2 + b$$

$$f_{w,b}(x) = g(z)$$

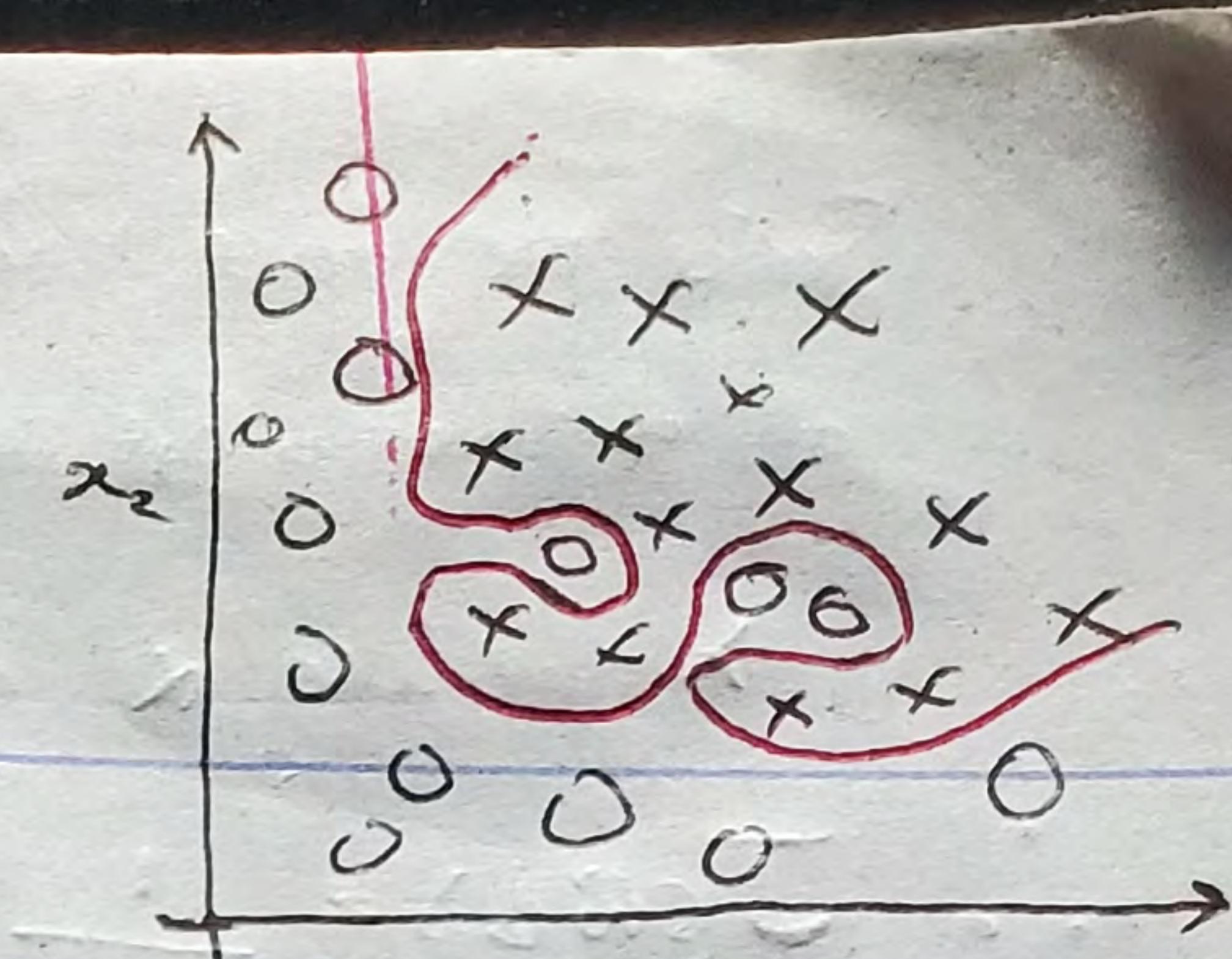
UNDERFIT



$$z = w_1 x_1 + w_2 x_2 + w_3 x_1^2$$

$$+ w_4 x_2^2 + w_5 x_1 x_2 + b$$

Just right ✓



$$z = w_1 x_1 + w_2 x_2 + w_3 x_1^2$$

$$+ w_4 x_2^2 + w_5 x_1^3$$

$$+ w_6 x_2^3 + \dots + b$$

OVERFIT ✗

A: ADDRESSING OVERFITTING

#1. Gathering more training data helps address overfitting [more training data for overfitting]

#2. Reducing the No. of features used also addresses overfitting

↳ Lots of features + Insufficient data → overfitting
 Feature selection
 Select features + Insufficient data → Just right

#3. Regularization: A way to more gently reduce the impact of some of the features when doing something as harsh as eliminating it outright. It encourages the learning algorithm to shrink the values of the parameters without necessarily demanding that they're set to exact 0.

→ It lets you keep all your features but prevents them from having an overly large effect which what mostly cause overfitting

WR3 REASIDE 2: REGULARIZATION

• If you have a lot of features, n features. You wouldn't know which features are important and which should be penalized so regularization penalizes all the w_i parameters. This results in a ~~smooth~~ & simpler function less prone to overfitting

Example: Let's say we have these features:

size x_1	bedrooms x_2	floors x_3	age x_4	avg income x_5	...	dist to coffee shop x_{100}	price y
---------------	-------------------	-----------------	--------------	---------------------	-----	----------------------------------	--------------

we have n features, $n \geq 100$!. Reg parameters = $w_1, w_2, w_3, \dots, w_{100}, b$

Since we don't know which parameters are the important ones, let's ~~penalize~~
all of them a bit and shrink them by adding $\frac{\lambda}{2m} \sum_{j=1}^n w_j^2$ to our cost function
 λ = regularization param

λ , just like α , can be chosen $\lambda > 0$

$$\therefore J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(x^{(i)}) - y^{(i)})^2 + \left[\frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \right] \quad \begin{array}{l} \text{regularization term} \\ \text{term (ii)} \end{array}$$

\downarrow mean squared error

λ is divided by $2m$ to scale it down so as to help us choose a good value

for λ . Therefore even if m increases, the same value of λ is also more likely to work because of the extra scaling ($\frac{1}{2m}$).

• Realizing the b term parameter doesn't really change anything.

$$J_{\text{MSE}} = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(x^{(i)}) - y^{(i)})^2$$

• In the modified cost function, we want to minimize the original cost \approx Mean square error cost
+ additionally the regularization term

— This has 2 goals \rightarrow (i) Minimize the mean squared error encourage the algorithm to fit the data well by minimizing the ~~cost~~ loss

(ii) Minimizing the regularization term encourages the algorithm to ~~keep~~ the parameters w_j small which tends to reduce overfitting.

The value of λ determines how you balance the two goals.

~~underfit~~ Increasing λ tends to ~~cause~~ your learning algorithm to decrease the size of w_j

~~overfit~~ Decreasing λ tends to cause your learning algorithm to increase the size of w_j

WK3 REGULARIZED LINEAR REGRESSION

We aim find parameters \vec{w} ; $\pm b$ that minimize the cost function $J(\vec{w}, b)$

$$\min_{\vec{w}, b} J(\vec{w}, b) = \left[\frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \right]$$

Baronally, we used gradient descent for the mean squared error portion only now we have the regularization term -

repeat {

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} w_j$$

$$b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$$

$$\rightarrow = \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \quad \text{this stays the same because we don't regularize } b$$

Let's look @ w_j (note that I've substituted the term for $\frac{\partial}{\partial w_j} J(\vec{w}, b)$ in & I'm grouping like terms)

$$w_j = w_j - \alpha \frac{\lambda}{m} w_j - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$= w_j (1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)}$$

recall that it was our initial unregularized value for $\alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$

\therefore the only change in the gradient descent for regularized linear regression is that instead of

$$w_j = \alpha \frac{\partial}{\partial w_j} J_{\text{unreg}}$$

we now have $w_j (1 - \alpha \frac{\lambda}{m}) - \alpha \frac{\partial}{\partial w_j} J_{\text{unreg}}$ [I hope you understand this I do]

Note: I'm differentiating between unregularized & regularized linear regress.

\therefore @ every update, instead of having w_j , you have e.g.; $w_j (1 - 0.01 \frac{1}{50})$
 this effectively shrinks w_j on every iteration after each $\rightarrow 0.9998 w_j$

WK3 REGULARIZED LOGISTIC REGRESSION

~~Main~~ Our modified cost function that uses regularization. π ;

$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

\rightarrow Minimizing the cost func., $\min_{\vec{w}, b} J(\vec{w}, b)$ penalizes parameters \vec{w} & prevents them from being ^{too} large

\rightarrow For gradient descent, our gradient $\frac{\partial}{\partial w_j} J(\vec{w}, b)$

$$\frac{\partial}{\partial w_j} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m [(f_{\vec{w}, b}(\vec{x}_j^{(i)}) - y^{(i)}) x_j^{(i)}] + \frac{\lambda}{m} w_j$$

It's the exact same equation as the linear regression except for the fact that $f(x) = g(z)$