

What we will be learning:

**Supervised Learning (regression, classification)**

**Unsupervised Learning (clustering, afterwards)**

Regression Model

$\hat{y}$  stands for estimated value of  $y$      $y$  = target value (unknown?)

$f$ , function is also called THE MODEL

what is the math formula we're going to use to compute  $f$ ?

$f_{w,b}(x) = wx + b$  OR DIRECTLY  $f(x)$

### **Constructing a Cost Function**

$w, b$ : parameter/coefficient/weight (it determines your function)

$w$  = slope

$(\hat{y} - y) = \text{error}$

We use  $J$  to refer to our cost function

We made a "squared error cost function". (Least squares estimation)

Cost functions differ from people to people, or depending on many things like our problem.

\*goal of linear regression is to: minimize  $J(w, b) \rightarrow w$  is an array,  $b$  is a constant

We saw how value of  $J$  differs depending on parameters,  $w$  and  $b$ . We need to try them to find values of  $w$  and  $b$ , that minimizes  $J$ . But we can't try them on our own, one by one.

SOLUTION: algorithms that will do this task for us

**\*\*\*GRADIENT DESCENT      (and many many more)**

You can try to use it to minimize any function.

\*for linear regression, with squared error cost function, you always get a bow shape

for cost ( $J$ ) in the end

Start with initial values. (0 if convenient) Keep changing "parameters to try to

reduce  $J$ .      THERE MAY BE MORE THAN 1 MINIMUMS

FOR FURTHER EXPLANATION, THINK OF THE EXAMPLE OF HIGHHILLS AND VISUALIZE IT

you may end up in different valleys, if there are multiple local minimums

\*gradient descent continuously tries to get a lower value. So if you find a minimum, you are done and have no place to go.

## GRADIENT DESCENT

$$\text{ALGORITHM } \rightarrow \quad \underset{dw}{w} = w - \alpha \cdot d/J(w,b) \quad \underset{db}{b} = b - \alpha \cdot d/J(w,b)$$

$\alpha$  = learning rate, positive number between 0 and 1 (it's small)

REPEAT THIS UNTIL CONVERGE (the point where  $w$  and  $b$  no longer changes)

\*parameters are updated "*simultaneously*"

Use of derivative is explained in "G.D. Intuition" part. If your parabola is going right and up, then derivative is positive, so a number is being subtracted from  $w$ .  $w$  decreases.

If parabola is going left and up, then derivative is negative; therefore a positive number is being added to  $w$ .  $w$  increases. This step happens repeatedly and you get closer to the optimum value of  $w$ .

NEXT: WHAT IS "LEARNING RATE"? What happens when it's too small or too big?

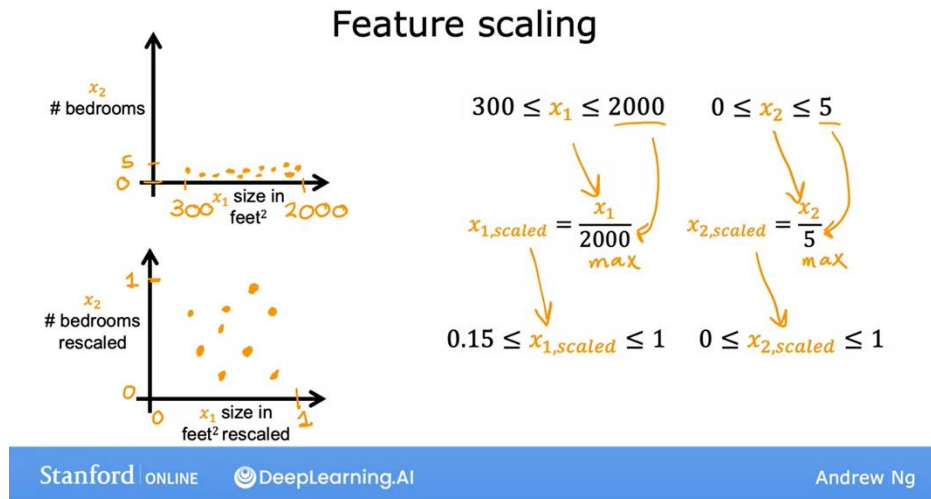
How am i supposed to choose a better learning rate?

## FEATURE SCALING

Parameters will most likely be in wide ranges and they will differ. One parameter will be in range(5) while other belongs to the range(2000).

So you need to scale them to work with both of them.

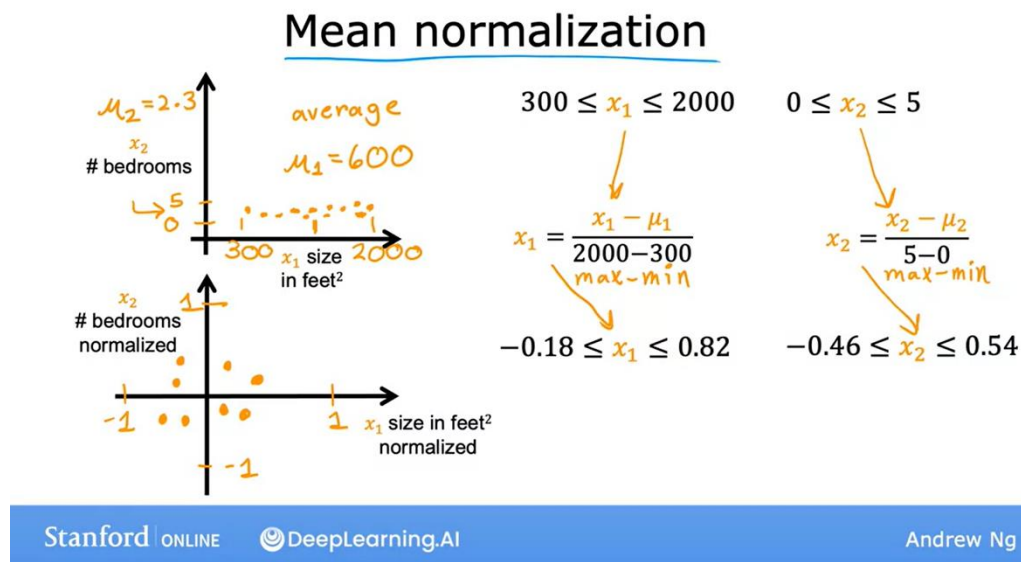
Solution: Divide the parameter with the max value. So maximum value of it is 1.



## ANOTHER OPTION: MEAN NORMALIZATION

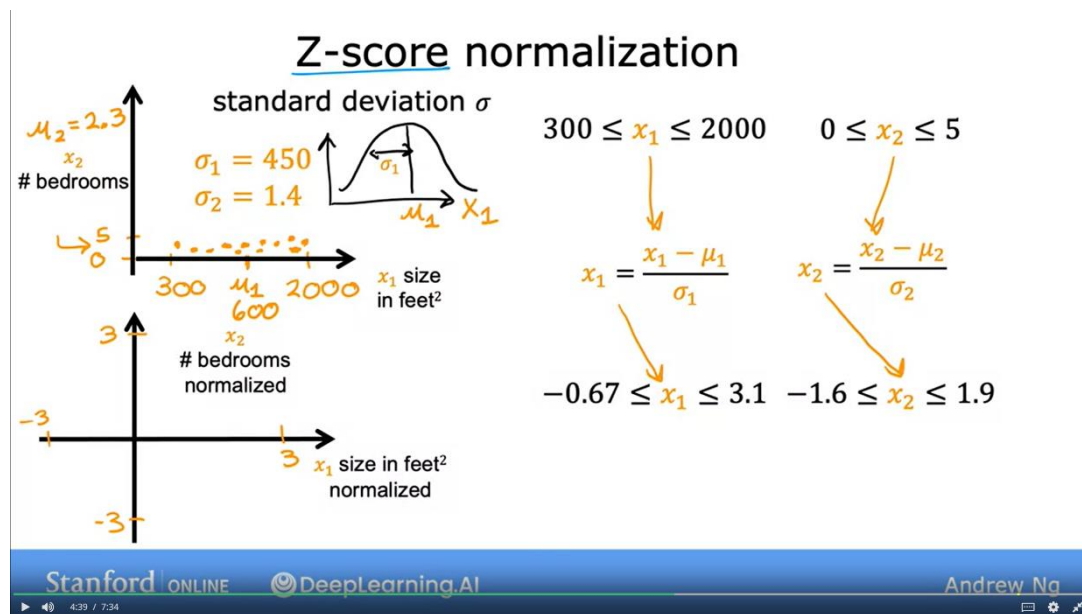
In this option, both parameters will be scaled to a state in which they are centered around (0,0). This option enables both negative and positive values.

- 1) Find mean value of the parameter.



## ANOTHER OPTION: Z-SCORE NORMALIZATION

For that, you will need standard deviation. This method also makes data center around (0,0).



RULE OF THUMB in feature scaling: Aim for  $-1 < x < 1$  for each feature. These values can loose a little bit. It's alright.

Too small values are a problem, just like the ones that are too big.

Also  $98,6 < x < 105$  is a problem as well.

**In the next one:** How to recognize if the grad. desc. is descending?

## Checking gradient descent for convergence

One of the key choices: Choice of learning rate

A two dimensional graph in which one axis is J and the other one is iterations: Learning curve

→ Desired behavior: J decreases. If it increases, even for a single iteration, sth is wrong.

After a number of iterations, J might cease to decrease, meaning it converges. Number of iterations needed to converge varies a lot. Telling this in advance is hard. So we use a graph.

ANOTHER WAY TO DECIDE IF THE TRAINING IS DONE:

### AUTOMATIC CONVERGENCE TEST

Declare an "epsilon" value. If the amount J decreases is less than "epsilon", declare convergence.

An.NG: I usually use graphs, due to lack of insight in choosing the "epsilon" value.

# Choosing the learning rate

Too small: Runs slow      Too big: Will diverge

If J does not decrease continuously: Bug in your code or too large learning rate

Try different values. Increase them slowly, like

... 0,001 → 0,003 → 0,01 → 0,03 → 0,1 → 0,3 → 1 ...

At some point you will find a value which is too large. Descend slowly from that value.

Optional lab:

Feature scaling, essentially **dividing each positive feature by its maximum value, or more generally, rescale each feature by both its minimum and maximum values using  $(x - \min) / (\max - \min)$ . Both ways normalizes features to the range of -1 and 1, where the former method works for positive features which is simple and serves well for the lecture's example, and the latter method works for any features.**

Mean normalization:  $x_i := (x_i - \mu_i) / (\max - \min)$

Z-score normalization which we will explore below.

**\*\*The scaled features get very accurate results much, much faster!** (Shown in lab with time values compared.)

**\*\*Any predictions using the parameters learned from a normalized training set must also be normalized.**

**\*\*The result is that updates to parameters during gradient descent can make equal progress for each parameter.**

**Feature Engineering:** Intuition about the problem is used to design new features, by transforming or combining original features.

e.g: we have width and depth. We can add an another feature like w.d (area)

\*It makes learning algorithm's job easier and faster.

## POLYNOMIAL REGRESSION

Linear re. is pretty useful but mostly it won't suffice

**\*\*\*Quadratic equa. Eventually comes down. This behaviour is mosly wrong. Can use cubic maybe?**

These are just a few examples.

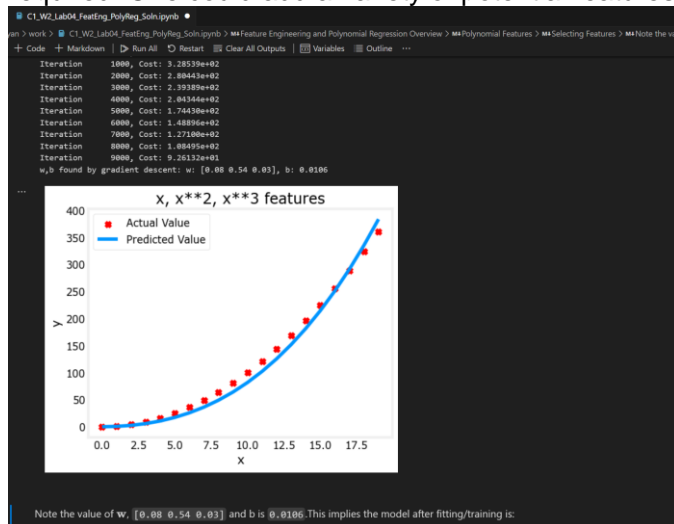
**If you are using powers in the equation like this, feature scaling becomes increasingly important.** That's because when you get a feature's square, it's range increases dramatically.

Getting square root of  $x$  might also work well in some occasions.

(in the 2. course we cover how to choose and scale appropriately)

## Selecting Features

Above, we knew that an  $x^2$  term was required. It may not always be obvious which features are required. One could add a variety of potential features to try and find the most useful



In the lab above, we used  $x$ ,  $x^2$ ,  $x^3$  instead of just  $x^2$ . But gradient descent itself has emphasized that  $x^2$  is the most suitable feature, by increasing it and decreasing others.

**Gradient descent is picking the 'correct' features for us by emphasizing its associated parameter**

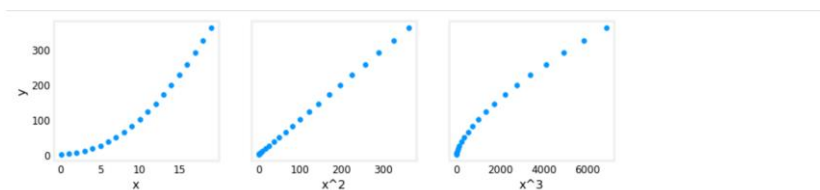
Let's review this idea:

- Initially, the features were re-scaled so they are comparable to each other
- less weight value implies less important/correct feature, and in extreme, when the weight becomes zero or very close to zero, the associated feature is not useful in fitting the model to the data.
- above, after fitting, the weight associated with the  $x^2$  feature is much larger than the weights for  $x$  or  $x^3$  as it is the most useful in fitting the data.

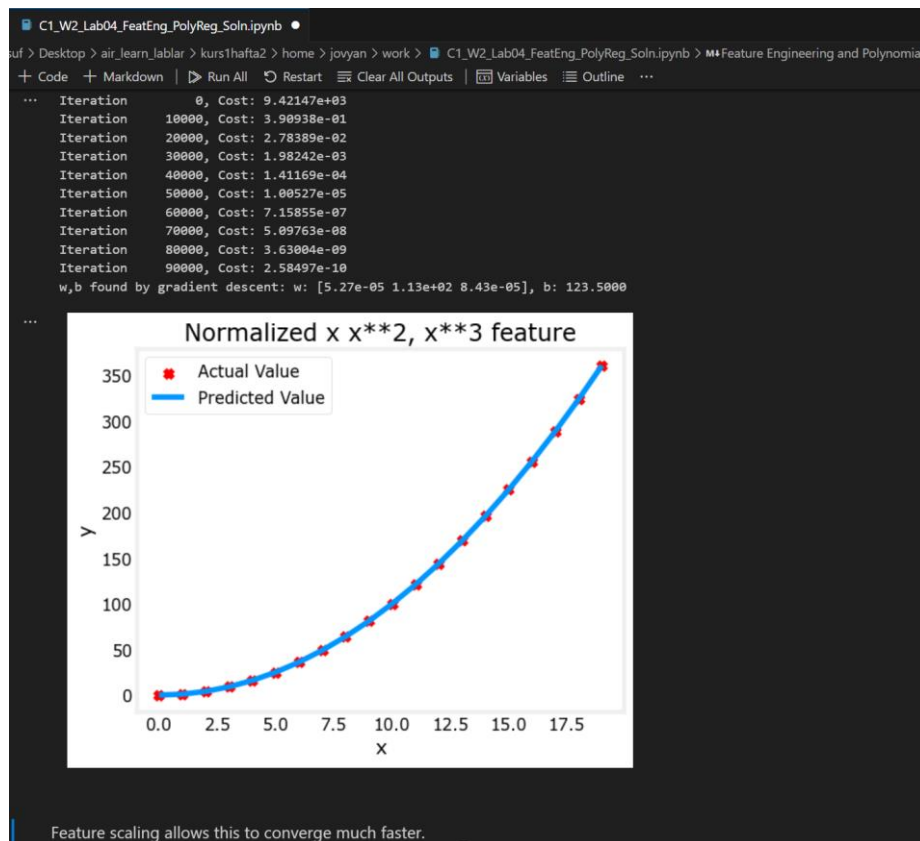
## An Alternate View

Above, polynomial features were chosen based on how well they matched the target data. Another way to think about this is to note that we are still using linear regression once we have created new features. Given that, the best features will be linear relative to the target. This is best understood with an example.:

**IN OTHER WORDS:** More suitable features will show a more “linear like” graph. That’s because it’s weight increases and effects the cost function more.



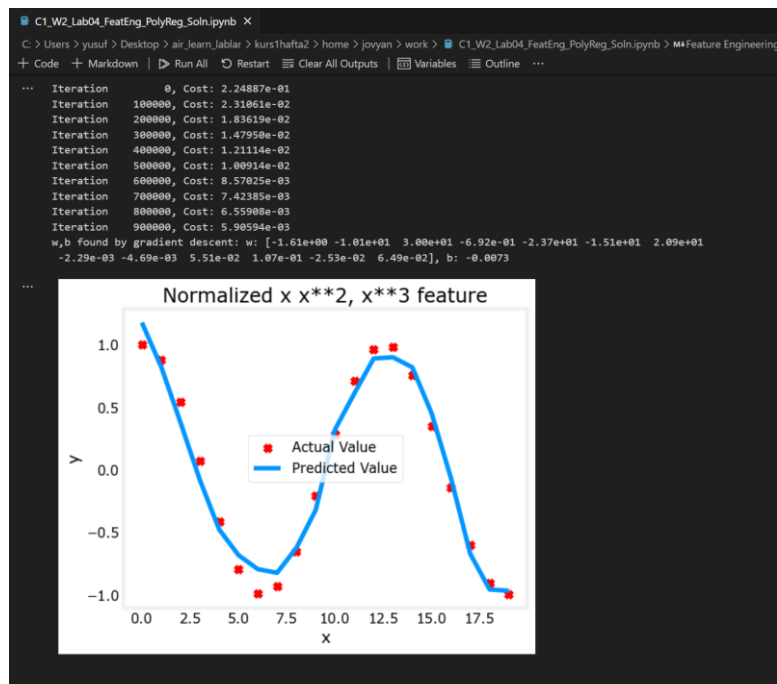
## VERY IMPORTANT EXAMPLE ABOUT FEATURE SCALING



## LOOK AT THAT RATE OF DESCEND

It converges soo much faster.

**Note:** With feature engineering, even very complex functions can be modeled.



# Lessons From Programming Lectures

- 1) **Get familiar with your data.** Inspect your  $x, y$  and their shapes.  
Shape means number of training examples in your dataset.
- 2) **Visualize your data**  
It is often useful to understand the data by visualizing it.  
For this dataset, you can use a scatter plot to visualize the data, since it has only two properties to plot (profit and population).  
\*Many other problems that you will encounter in real life have more than two properties (for example, population, average household income, monthly profits, monthly sales). When you have more than two properties, you can still use a scatter plot to see the relationship between each pair of properties.
- 3) **Your goal is to build a linear regression model to fit this data.**  
With this model, you can then input a new city's population, and have the model estimate your restaurant's potential monthly profits for that city.
- 4) Compute cost, gradient descent, learning parameters using batch grad.desc.  
**Recall batch refers to running all the examples in one iteration.**

## WEEK 3: CLASSIFICATION

$Y$  can only be certain values.

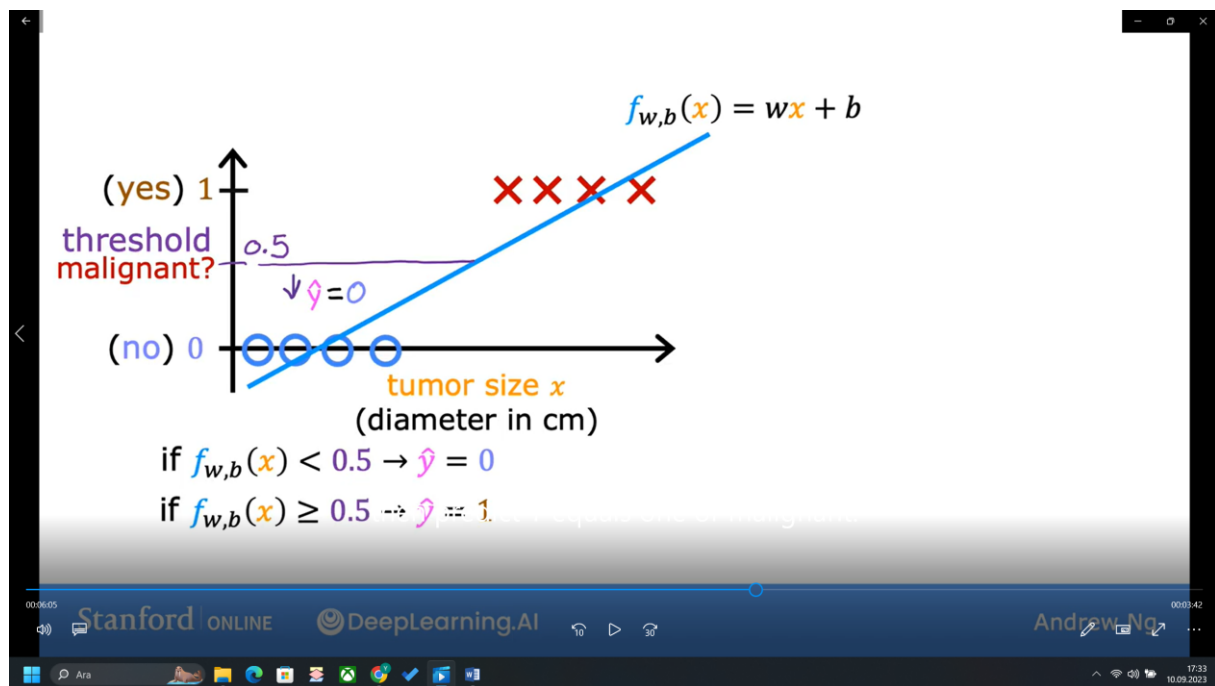
Yes or no → binary classification

class and category mean basically same thing. Classes are “yes” or “no” or etc.

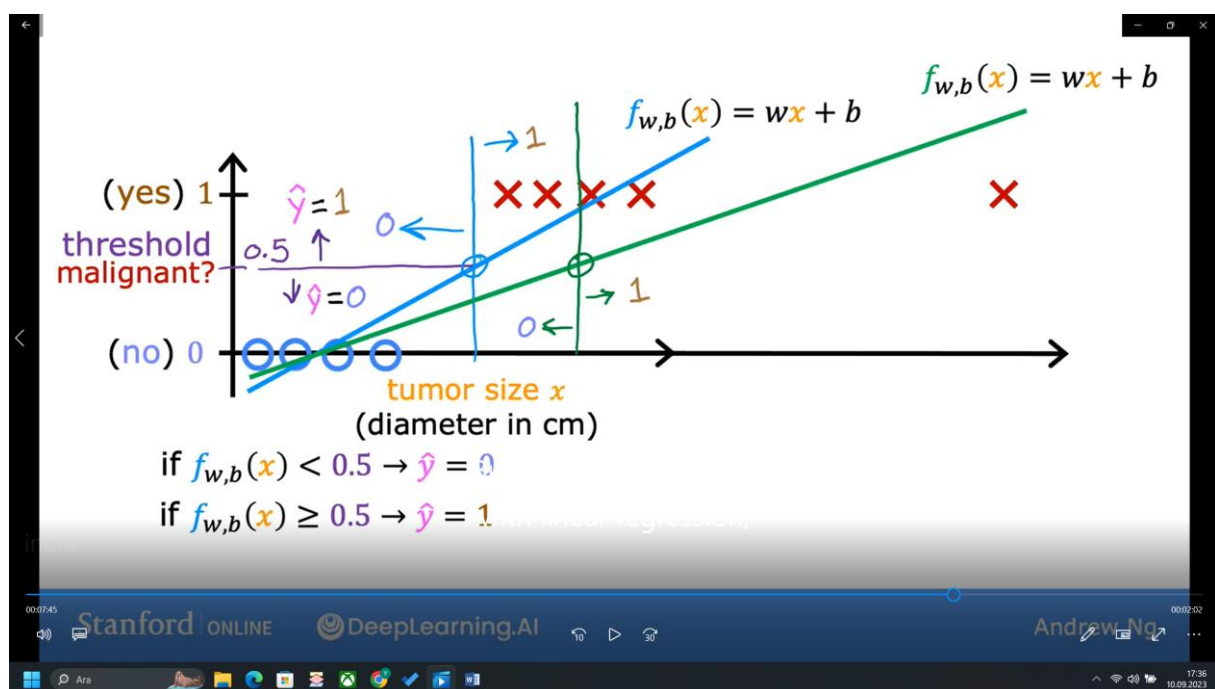
0 and 1 do not have semantic meanings. They refer to “absence” or “presence” of the thing we are looking for.

\*Linear regression does not work in these problems.





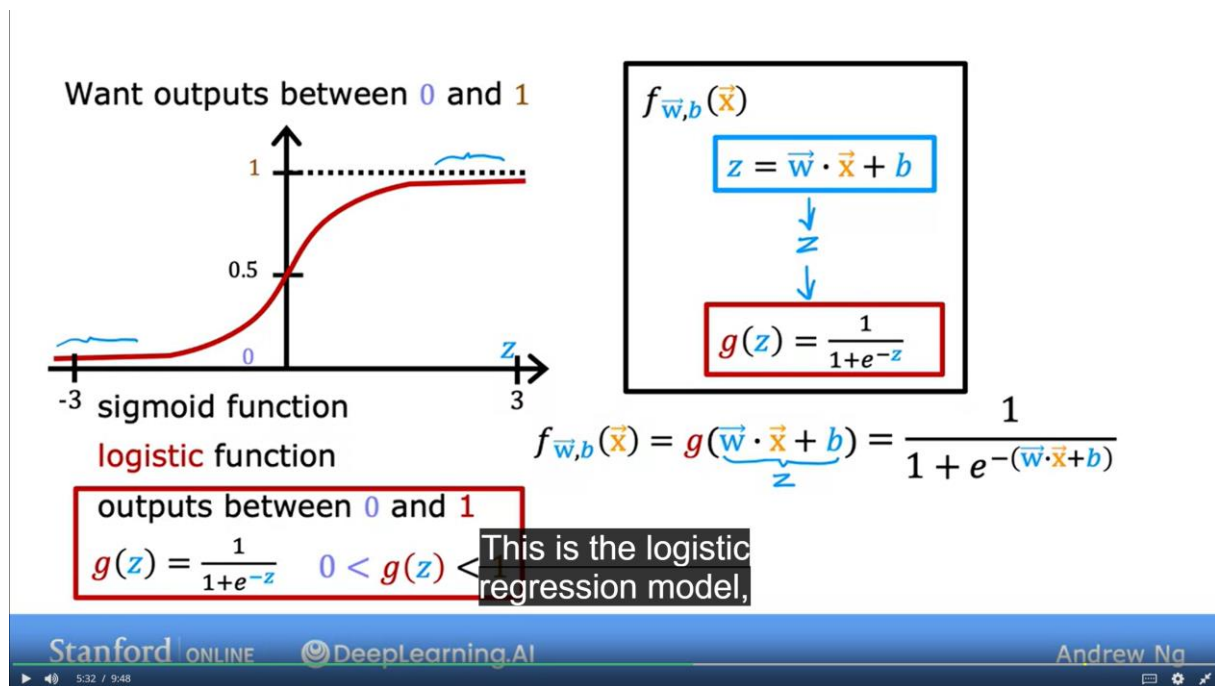
But still, in a situation you need to use it, you may apply rules, like picking thresholds and if number > 0.5, it's 1.



This might, however, work well in the first example and work poorly in the second one.

**\*\*Threshold = Decision boundary.** In linear regression it changes in a poor way.

## Sigmoid Function and Logistic Regression



Get  $z$  from linear function, pass it to the sigmoid function.

It gets  $x$ , outputs numbers between 0 - 1. Output actually means the probability for the class to be 1.

$F(x) = 0.7 \rightarrow$  %70 chance that  $y$  is 1.  $Y$  has got to be either 1 or 0.

## Interpretation of logistic regression output

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

"probability" that class is 1

$$f_{\vec{w},b}(\vec{x}) = P(y = 1 | \vec{x}; \vec{w}, b)$$

Probability that  $y$  is 1,  
given input  $\vec{x}$ , parameters  $\vec{w}, b$

Example:

$x$  is "tumor size"  
 $y$  is 0 (not malignant)  
or 1 (malignant)

$$P(y = 0) + P(y = 1) = 1$$

$f_{\vec{w},b}(\vec{x}) = 0.7$   
70% chance that  $y$  is 1 because you may see this in other places.

Quote from Andrew NG: "For a long time, a lot of Internet advertising was actually driven by basically a slight variation of logistic regression. This was very lucrative for some large companies, and this is basically the algorithm that decided what ad was shown to you and many others on some large websites."

FROM LAB:

In the case of logistic regression,  $z$  (the input to the sigmoid function), is the output of a linear regression model.

- In the case of a single example,  $z$  is scalar.
- in the case of multiple examples,  $z$  may be a vector consisting of  $m$  values, one for each example.

**Note:** NumPy has a function called `exp()`, which offers a convenient way to calculate the exponential of all elements in the input array.

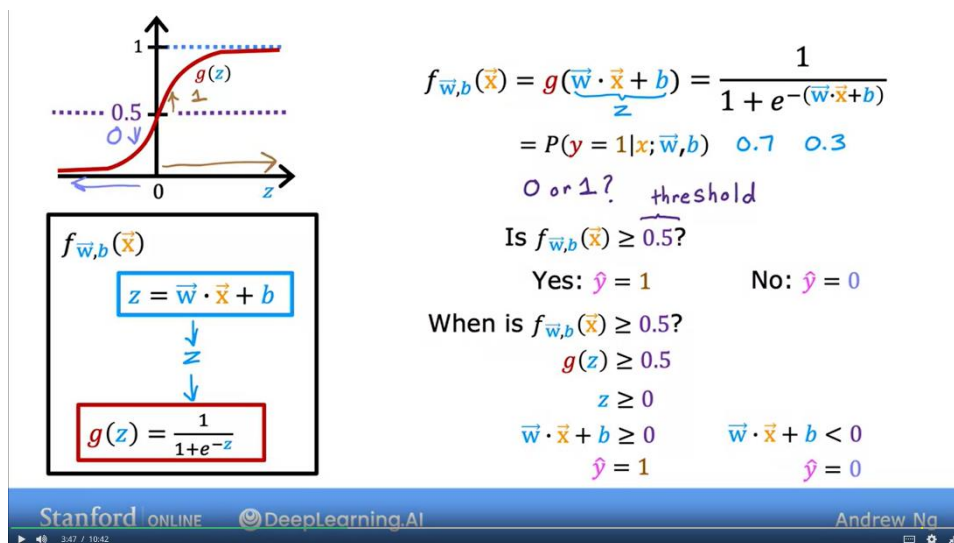
## Decision boundary

- 1) compute  $z$
- 2) apply sigmoid function to the  $z$

Threshold  $\rightarrow 0.5$  is a common choice. Whether the prediction will be 0 or 1, inspect the formulas, depend on the sign of  $z$  ( $w \cdot x + b$ ).

\*Tanh can also be used for the interval  $(-1, 1)$ . (CAN STILL REGULARIZE IT TO INTERVAL 0-1)

(it's what i did for my differential equations project)

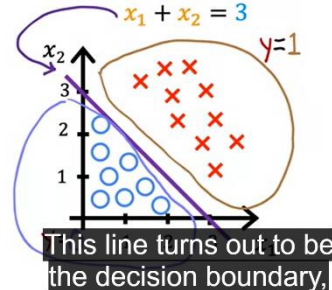


On decision boundary, you are neutral about whether y is 0 or 1.

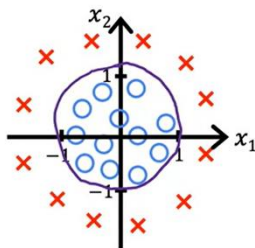
## Decision boundary

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

Decision boundary  $z = \vec{w} \cdot \vec{x} + b = 0$   
 $z = x_1 + x_2 - 3 = 0$   
 $x_1 + x_2 = 3$



## Non-linear decision boundaries



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

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

this turns out to be the circle.

Just like polynomial regression, non-linear functions are used in logistic regression

You can also see more complex decision boundaries.

\*If your parameters are like “x1, x2, x3...” your decision boundary will always be linear.

\*Let’s say you are creating a tumor detection algorithm. Your algorithm will be used to flag potential tumors for future inspection by a specialist. What value should you use for a threshold?



NOPE

High, say a threshold of 0.9?



YUP

Low, say a threshold of 0.2?

**Correct: You would not want to miss a potential tumor, so you will want a low threshold.** A specialist will review the output of the algorithm which reduces the possibility of

a 'false positive'. The key point of this question is to note that the threshold value does not need to be 0.5.

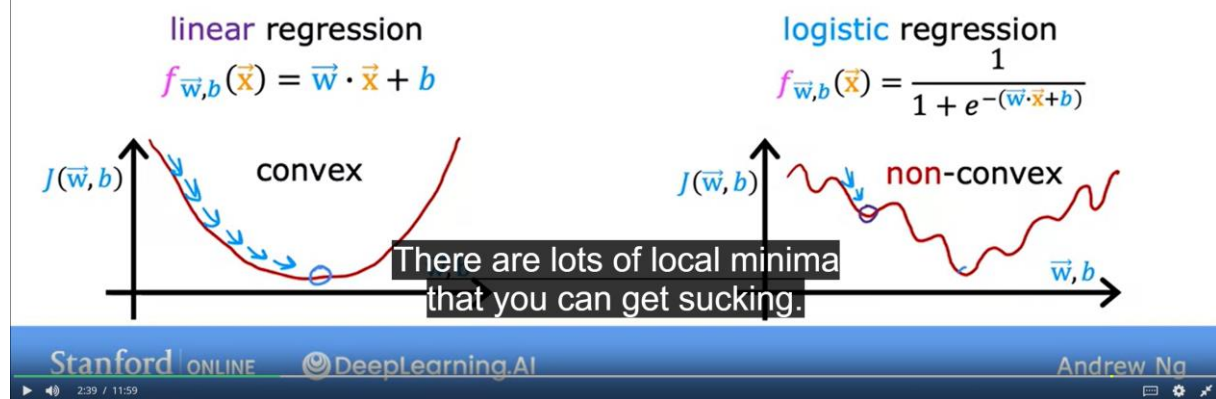
As we've seen in the lectures, by using higher order polynomial terms (eg:  $f(x)=g(x^{20}+x^{1-1})$ ), we can come up with more complex non-linear boundaries.

## Cost function for logistic regression

Mean squared based cost function does not work for log.reg.

### Squared error cost

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



Loss: Is the cost function working on a single training example.  $\rightarrow L$

**Cost function is the mean value of loss functions.** Loss function measures how well you are performing on one training example. Their sum gives you the cost value(entire training set).

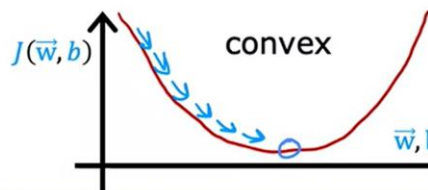
## Squared error cost

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

$$\text{loss} \quad L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})$$

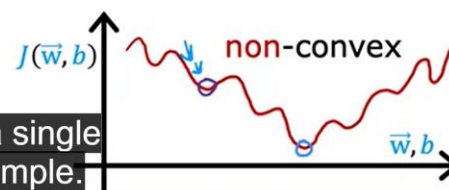
linear regression

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



logistic regression

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

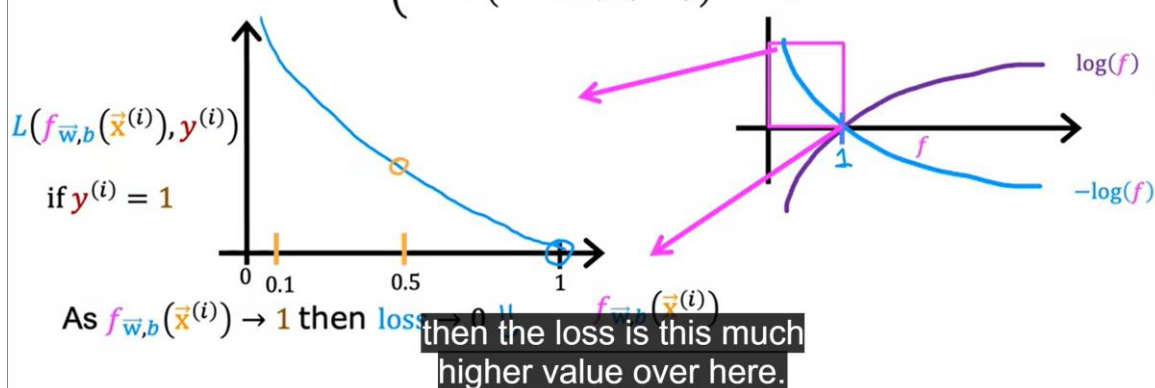


the loss on a single training example.

\*In loss function, we use log.

## Logistic loss function

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

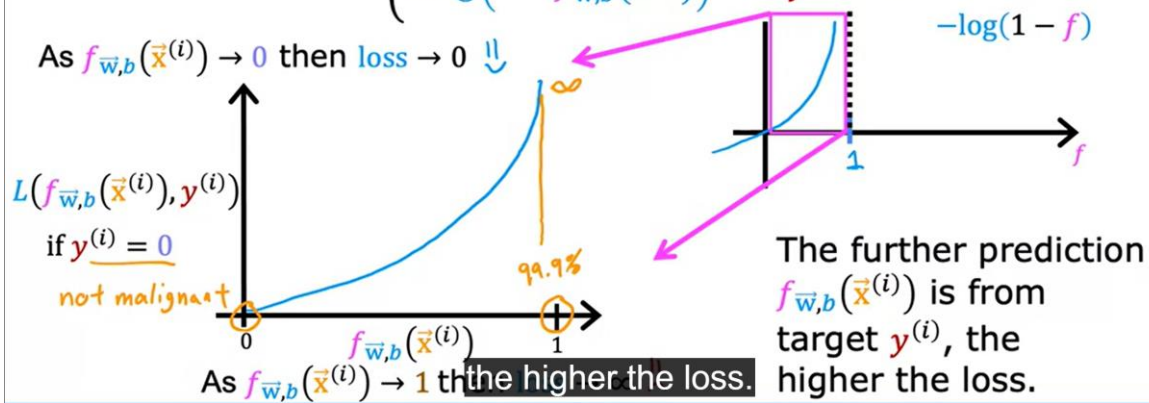


Logarithmic function provides the following: If predicted label is wrong, loss of that example is a lot higher than the desired amount.



## Logistic loss function

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



Stanford ONLINE

DeepLearning.AI

Andrew Ng

## Cost

cost

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{L(f_{\vec{w},b}(\vec{x}^{(i)}), y^{(i)})}_{\text{loss}}$$

$$= \begin{cases} -\log(f_{\vec{w},b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w},b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases} \quad \text{convex}$$

the overall cost function will be convex and thus you can

Stanford ONLINE

DeepLearning.AI

Andrew Ng

This choice of loss function makes that cost function is convex. Therefore gradient descent can find the global minimum. **(Proving that it's convex is beyond the scope of this course)**

## Cost

$$\begin{aligned}
 \text{cost} \\
 J(\vec{w}, b) &= \frac{1}{m} \sum_{i=1}^m L(\underbrace{f_{\vec{w}, b}(\vec{x}^{(i)})}_{\text{loss}}, y^{(i)}) \\
 &= \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \text{ convex} \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \text{ global minimum} \end{cases} \quad \text{can reach a} \\
 &\quad \text{find } w, b \text{ that minimize cost } J
 \end{aligned}$$

**w and b for logistic regression.**



**Loss** is a measure of the difference of a single example to its target value while the **Cost** is a measure of the losses over the training set

The loss function above can be rewritten to be easier to implement.

$$loss(f_{\vec{w}, b}(\mathbf{x}^{(i)}), y^{(i)}) = (-y^{(i)} \log(f_{\vec{w}, b}(\mathbf{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\mathbf{x}^{(i)})))$$

This is a rather formidable-looking equation. It is less daunting when you consider  $y^{(i)}$  can have only two values, 0 and 1. One can then consider the equation in two pieces:

when  $y^{(i)} = 0$ , the left-hand term is eliminated:

$$\begin{aligned}
 loss(f_{\vec{w}, b}(\mathbf{x}^{(i)}), 0) &= (-0) \log(f_{\vec{w}, b}(\mathbf{x}^{(i)})) - (1 - 0) \log(1 - f_{\vec{w}, b}(\mathbf{x}^{(i)})) \\
 &= -\log(1 - f_{\vec{w}, b}(\mathbf{x}^{(i)}))
 \end{aligned}$$

and when  $y^{(i)} = 1$ , the right-hand term is eliminated:

$$\begin{aligned}
 loss(f_{\vec{w}, b}(\mathbf{x}^{(i)}), 1) &= (-1) \log(f_{\vec{w}, b}(\mathbf{x}^{(i)})) - (1 - 1) \log(1 - f_{\vec{w}, b}(\mathbf{x}^{(i)})) \\
 &= -\log(f_{\vec{w}, b}(\mathbf{x}^{(i)}))
 \end{aligned}$$

## Simplified Cost Function for Logistic Regression

A better way to write loss and cost so the implementation is easier with gradients.

It's the formula on top.  $y$  can only be 0 or 1, so these values can make the other part of the formula disappear.

Multiply the " $y=1$ " part with  $y$  and the " $y=0$ " part with  $(1-y)$ . Each value makes the other part be 0.

Now we write the simplified Cost function. (below)

This formula is derived from statistics, using:

**maximum likelihood estimation.** (used to find optimum parameters)



## Simplified cost function

$$\begin{aligned}
 \text{loss} \quad L(f_{\vec{w},b}(\vec{x}^{(i)}), y^{(i)}) &= -y^{(i)} \log(f_{\vec{w},b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\vec{w},b}(\vec{x}^{(i)})) \\
 \text{cost} \quad J(\vec{w}, b) &= \frac{1}{m} \sum_{i=1}^m [L(f_{\vec{w},b}(\vec{x}^{(i)}), y^{(i)})] \quad \text{convex (single global minimum)} \\
 &= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\vec{w},b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w},b}(\vec{x}^{(i)}))] \\
 &\quad \text{maximum likelihood (don't worry about it!)}
 \end{aligned}$$

Now that we know these, let's learn how to implement gradient descent with this algorithm.

Our goal: Find parameters  $w, b$ .

Way to do so: Use gradient descent

Form of the gradient descent algorithm is the same as linear regression. WHAT CHANGES IS THE DEFINITION OF THE FUNCTION  $F(x)$ :

## Gradient descent for logistic regression

repeat { looks like linear regression!

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

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

} simultaneous updates

Same concepts:

- Monitor gradient descent (learning curve)
- Vectorized implementation
- Feature scaling

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

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

In the upcoming optional lab,

Time for the labs.

# The problem of overfitting

Next videos: There is a regularization technic to minimize this problem.

**Overfitting:** Model fits for the training data **more than needed**. It learns the noise etc. And this causes the model to perform poorly in test data.

You can even find parameters that make Cost equal to 0. Because there is no error. IT'S BAD THOUGH.

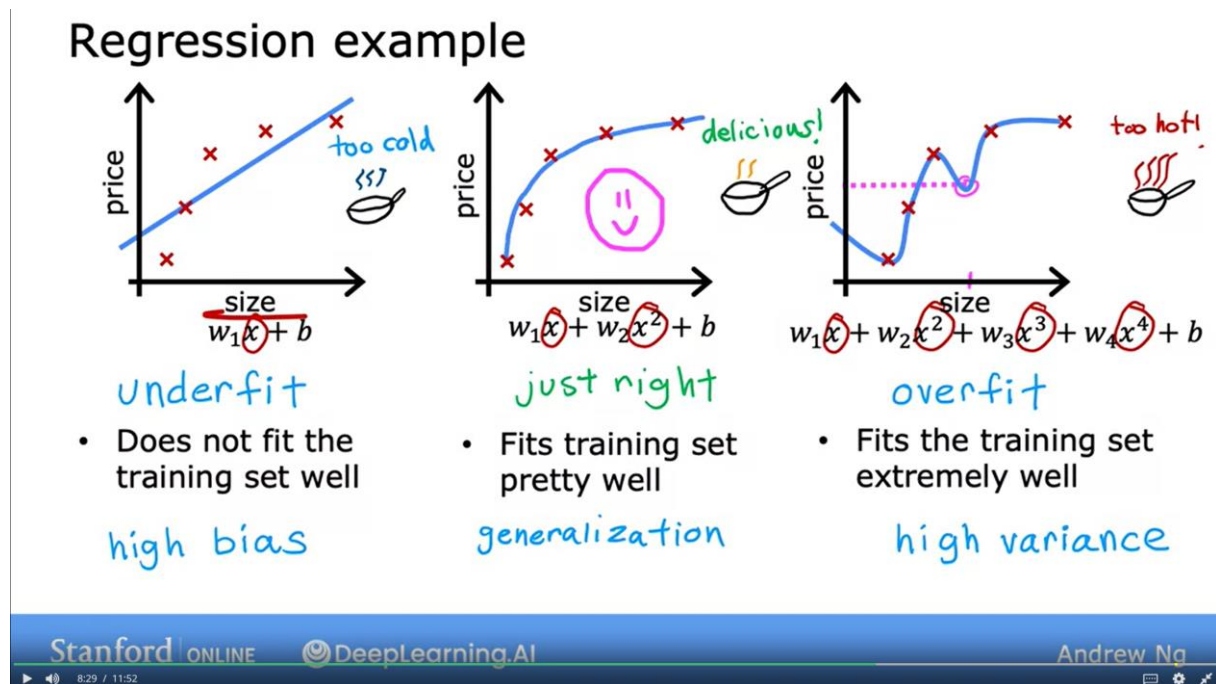
Model has high variance: If even one point of data was a bit different, function that has been fit might have been completely different.

**Underfitting:** Model is unable to fit. Might be too simple for example. Model has high bias.

**If you succeed:** Model is generalized.

**BIAS:** Model is unable to capture features of the data.

e.g: you assumed data is linear, you try to fit a line but data is more complicated.

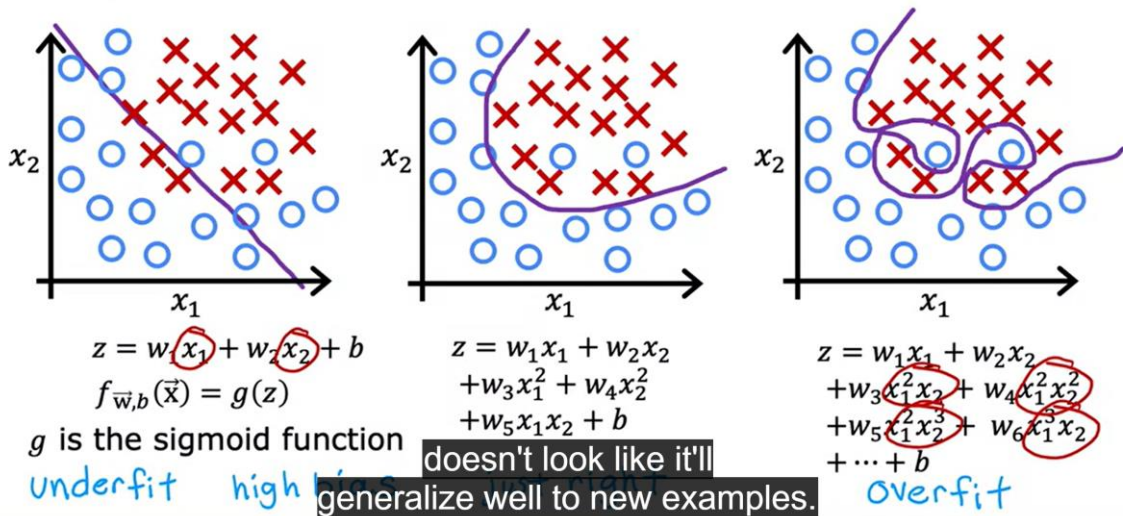


HAVING TOO MUCH FEATURES CAN CAUSE OVERFITTING.

Too few can cause underfitting.

\*Overfitting can occur in classification as well.

# Classification

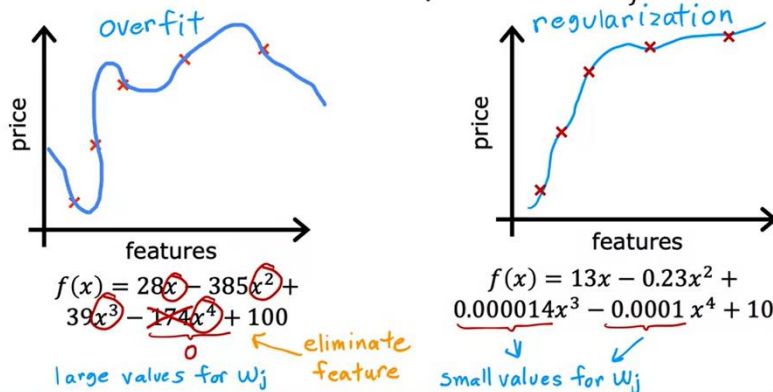


## ADDRESSING OVERFITTING

- 1) Collect more training data: This makes the function less wiggly. This way you can still use more features and still do fine.  
**You can not always get more data.**
- 2) See if you can use **fewer features**: Having a lot of features but insufficient data might result in overfitting because model can not generalize with that less data. Try using less features (subset of features). This is called **feature selection**. You use your intuition for this. Though by using this, you are throwing away some information and possibly harming your performance. **COURSE 2 INVOLVES FEATURE SELECTION ALGORITHMS.**
- 3) **REGULARIZATION**

## Regularization

Reduce the size of parameters  $w_j$



- Regularization lets you keep your features, but it prevents your features from having too much effect, which is what sometimes causes overfitting.
- We generally regularize parameters  $w_1, w_2, \dots, w_n$ . Regularizing  $b$  does not change much.  
Andrew: "I usually don't"

**Summarizing:**

- 1) Collect more data
- 2) Select features
- 3) Reduce size of parameters → regularization

**3. option is sth we almost always use for learning algorithms, especially for neural networks.**

Note: Regularization basically makes everything better. To sum it up, we implemented it with linear regression and logistic regression.