

ML: The science of getting computers to learn without being explicitly programmed

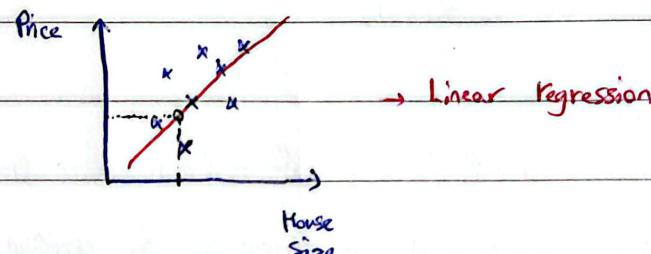
Supervised Learning: Algorithm that learn $x \rightarrow Y$ (input to output mapping)

* Learns from being given right answers

| E.g.: | Input (X) | Output (Y) | Application |
|-------|--------------------|------------------------|--------------------|
| | email | spam(0/1) | spam filtering |
| | audio | text transcript | speech recognition |
| | ad, user info | click? (0/1) | online advertising |
| | Images, radar info | position of other cars | self-driving car |

Regression: a particular type of supervised learning which predicts a number from infinitely many possible outputs

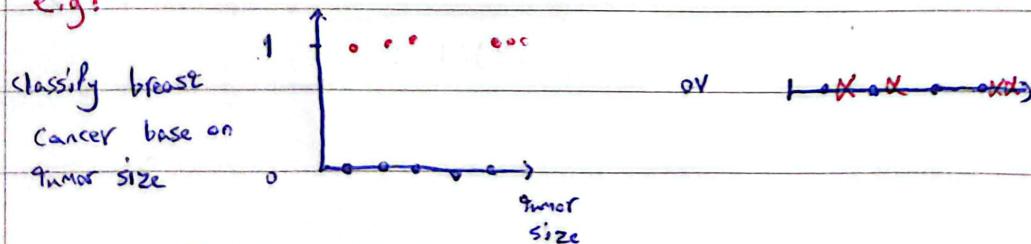
E.g.:



Classification: a particular type of supervised learning which predict categories.

* Classification predicts small finite limited set of possible output categories

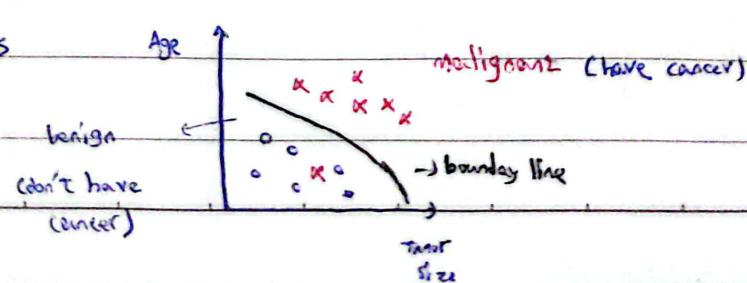
E.g.:



* Class = Category

* Categories can be non-numeric (cat or dog)

E.g.: Two inputs

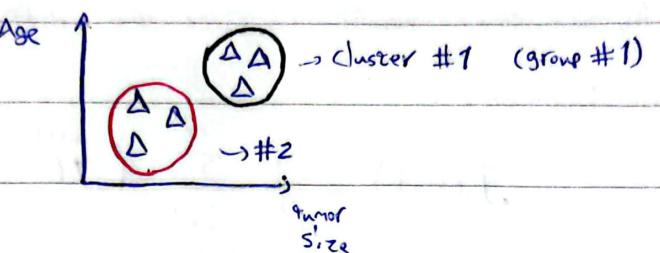


* learning algorithm has to decide how to fix a boundary line through data

Unsupervised Learning: We've given data that isn't associated with any output label y and our job is to find some structure or pattern or just something interesting.

e.g. Clustering: Places the

unlabeled data into different clusters



e.g. 2: Anomaly Detection: Find unusual data points (events) \rightarrow Fraud detection in financial system

e.g. 3: Dimensionality reduction: Compress data using fewer numbers (while losing as little info as possible)

Terminology and Process of supervised learning

Training set: Data used to train the model $\xrightarrow{\text{eg}}$

\hat{y} : estimate or the prediction for y

(\hat{y})

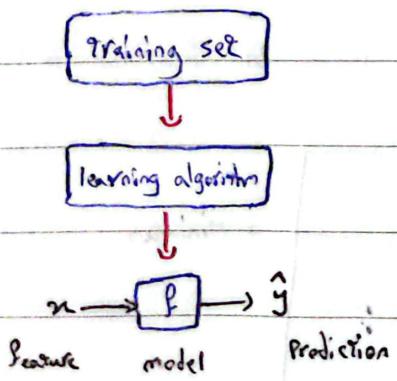
* \hat{y} may or may not be the actual true value

| x | y | Notation |
|--------------|------------------|-------------------------|
| size in feet | Price in 1000\$' | |
| (1) 2104 | 400 | x_1 'input' variable |
| (2) 1416 | 232 | 'feature' |
| : | : | |
| (47) 3230 | 870 | y_1 'output' variable |
| | | or 'target' |
| | | |
| | $m=47$ | |

m = number of training examples

$(x^{(i)}, y^{(i)})$ \Rightarrow i th training example

(x, y) \Rightarrow single training example



e.g.: size \rightarrow f \rightarrow Price

How to represent f ?

$$f_{w,b}(x) = w_1 x + b \rightarrow \text{Linear regression with one variable}$$

or
 $f(x)$

* This is also called univariate (one variable) linear regression.

Cost function: a function that would compute the difference between y and \hat{y}

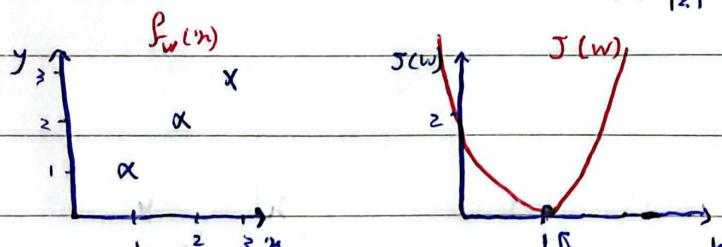
Squared error cost function:

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

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

* Imagine $b=0 \rightarrow f_w(x) = w_1 x$ and $J(w) = \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2$

e.g.:

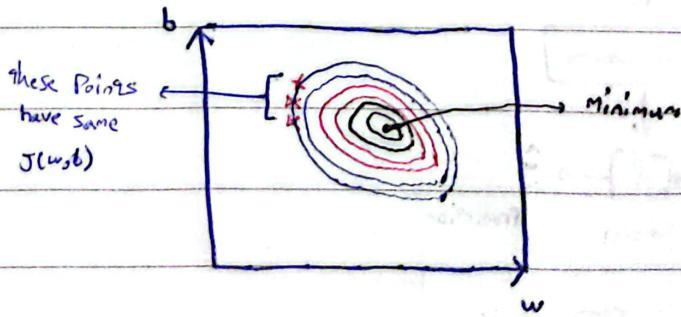


* Our goal is to minimize cost func or $J(w) \rightarrow$ so this point is the min point

$\rightarrow w=1 \rightarrow f_w(x) = x$ is the line that would have least cost and would fit the data pretty well.

In general: If $b \neq 0 \rightarrow f_{w,b}(x) = w_1 x + b$ and $J(w, b)$ is a 3D function, with axis w, b and $J(w, b)$.

Contour plot of $J(w, b)$:



Gradient descent: an algorithm that we can use to minimize cost func ($J(w_1, \dots, w_n, b)$)

Repeat until convergence (until you reach ^{the} point at a local minimum where w and b no longer

change much \Rightarrow with each

iteration)

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

α : learning rate

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

* Simultaneously update w and b

$$w = \text{tmp_}w$$

$$b = \text{tmp_}b$$

}

* if α is too small \rightarrow Gradient descent may be slow

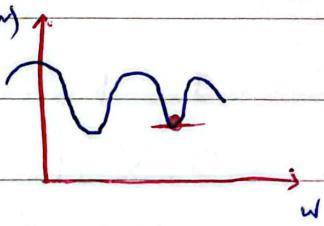
a) overshoot, never reach minimum

* if α is too large \rightarrow Gradient descent may:

b) Fail to converge, diverge

* If we are already at a local minimum, since the tangent line has slope $\approx 0 \rightarrow$ then $w = w - \alpha \cdot 0$

$\rightarrow w = w$, so gradient descent won't do anything!

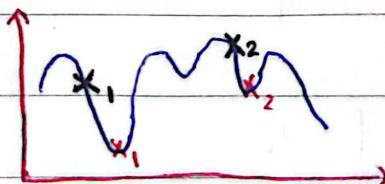


* as we get near a local minimum, derivative becomes smaller (moving towards 0), and the update steps become smaller \rightarrow so we can reach min without decreasing α

* Depending on where you initialize the parameters w and b , you can end up at different local minima

x_i : starting point i

x_i' : ending point i after
we run algorithm



* squared error cost func is a convex func. Convex funcs are bowl-shape funcs and have single local min (which is also global min as well).

e.g.: Linear regression model: $f_{w,b}(x_i) = w_i + b$ cost func: $J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$

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

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

* This version of gradient descent is called "Batch gradient descent"

* Batch: Each step of gradient descent uses all of the training examples (there are versions that with only a subset of them)

Multiple Features (Variables)

x_j j^{th} feature

n = Number of features

$\vec{x}^{(i)}$ feature of i^{th} training example

$x_j^{(i)}$ value of feature j in i^{th} training example

$$\vec{x}^{(1)} = [1, 2, 3, \dots, 5]$$

Model:

$$f_{w,b}(\vec{x}) = w_1 x_1 + \dots + w_n x_n + b$$

or

$$\vec{w} = [w_1, w_2, \dots, w_n] \quad b \text{ is a number} \quad \vec{x} = [x_1, x_2, \dots, x_n]$$

$$\rightarrow f_{w,b}(\vec{x}) = \vec{w} \cdot \vec{x} + b \quad \rightarrow \text{multiple linear regression}$$

dot Product

Vectorization :

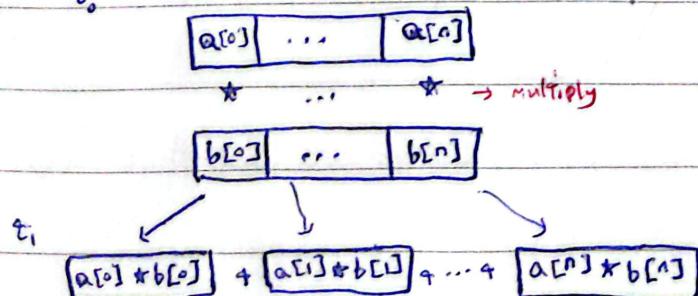
* Use np.dot (a, b) over for-loop.

benefits

1 - Shorter code

2 - Faster (cause it uses parallel hardware) !

! :



* efficient \rightarrow scale to large datasets

Gradient descent with vector notation (multiple features)

Repeat {

$$\text{for } i=1 \quad w_i = w_i - \alpha \frac{1}{m} \sum_{j=1}^m (f_{\vec{w}, b}(\vec{x}^{(j)}) - y^{(j)}) x_i^{(j)}$$

⋮

$$\text{for } i=n \quad w_n = w_n - \alpha \frac{1}{m} \sum_{j=1}^m (f_{\vec{w}, b}(\vec{x}^{(j)}) - y^{(j)}) x_n^{(j)}$$

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

}

An alternative to gradient descent: Normal equation

Normal equation:

- a) Only for linear regression
- b) Solve for w, b without iterations

* Gradient descent is still the recommended method

Disadvantages:

- 1. Doesn't generalize to other learning algos
- 2. Slow when number of features is large ($> 10^4$)

Feature and Parameter Values

Imagine: $y = w_1 z_1 + w_2 z_2 + b$

| | |
|------------------------------------|-------------------------------|
| z_1 Range: 300 - 2000 (large) | z_2 Range: 0 - 5 (small) |
|------------------------------------|-------------------------------|

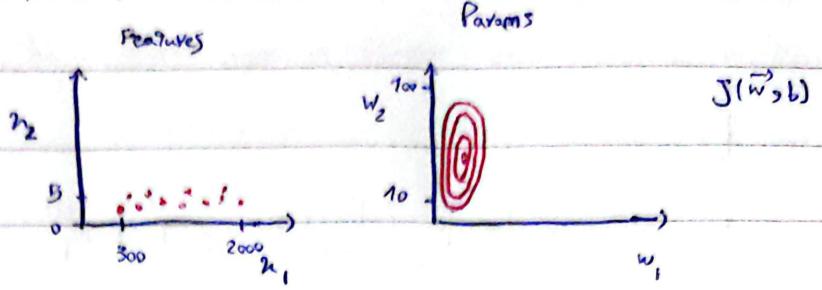
Training example: $z_1 = 2000 \quad z_2 = 5 \quad \text{Price (y)} = 500$

Possible Params $\rightarrow w_1 = 0.1$ (small) $w_2 = 50$ (large) $b = 50$

* When a possible range of values of a feature is large, a good model will learn to choose small parameter value

* When a possible range of values of one feature is small, a good model will learn to choose large parameter value

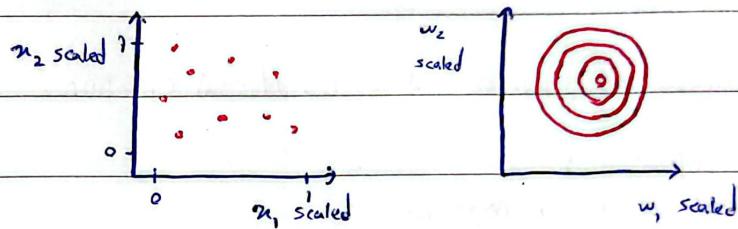
effect of this on gradient descent:



* If we run gradient descent with this, gd will end up bouncing back and forth for a long time before it can finally find its way to the min



→ What to do? Scale Features!



Feature scaling:

$$1 - \text{Divide by max: } 300 \leq x_1 \leq 2000 \rightarrow x_{1,\text{scaled}} = \frac{x_1}{2000} \rightarrow 0.15 \leq x_{1,\text{scaled}} \leq 1$$

$$2 - \text{Mean normalization: } x_{1,S} = \frac{x_1 - M_1}{\max - \min} \rightarrow -0.18 \leq x_{1,S} \leq 0.82$$

$$3 - Z\text{-score normalization: } x_{1,S} = \frac{x_1 - M_1}{\sigma_1} \rightarrow -0.67 \leq x_{1,S} \leq 3.1$$

* aim for a okay range $\rightarrow -1 \leq x \leq 1 \checkmark -3 \leq x \leq 3 \checkmark -100 \leq x \leq 100 \rightarrow$ too large

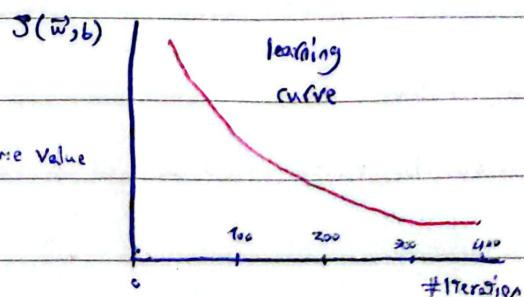
$-0.001 \leq x \leq 0.001 \rightarrow$ too small $98.6 \leq x \leq 104 \rightarrow$ too large

Make sure gd is working correctly

- $J(\vec{w}, b)$ should decrease after every iteration

- After 50 iterations, learning curve should converge toward some value

- # iterations needed varies



Another way: automatic convergence test:

Let ϵ be a small number. If $J(\vec{w}, b)$ decreases by $\leq \epsilon$ in one iteration \rightarrow declare convergence
(like 10^{-3})

* learning curve method is better due to difficulty of finding a good ϵ

Choosing the learning rate

We can make decisions about our learning rate (whether it's too large or too small) by looking at the learning curve or # iterations.

* a mistake in learning curve might also be caused by a bug in our program. We can identify that by using a small enough α and if $J(\vec{w}, b)$ increases in one iteration \rightarrow there is a bug!

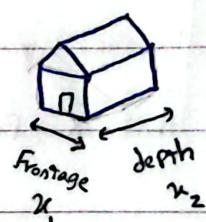
- Values of α to try: 0.001 0.003 0.01 0.03 0.1 0.3 1
 $\underbrace{}_{\alpha=3} \quad \underbrace{}_{\alpha=3}$

* Andrew tend to use the largest possible α that is reasonable and gd runs well for it

Feature engineering: Using intuition to design new features, by transforming or combining original features.

$$\text{e.g.: } f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x_2 + b$$

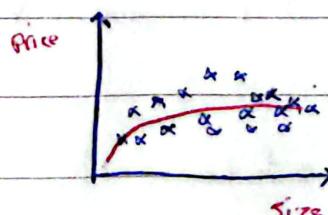
$$\text{area} = x_1 \times x_2 \rightarrow x_3 = x_1 x_2$$



$$\rightarrow f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

Polynomial regression: Sometimes a straight line doesn't fit data - see very well so we need to use a non-linear one.

* feature scaling is important here



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

Classification:

- **Binary classification:** There are only two possible outputs (two possible classes)

or
categories

| | | |
|------|-------|------|
| e.g: | No | Yes |
| | False | True |
| | 0 | 1 |

Negative class Positive class

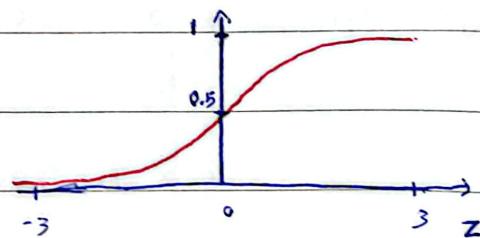
★ What happens if we use linear regression for classification? Sometimes you might get lucky and it may work but often it will not work well

Logistic regression:

- **Sigmoid function (logistic function):**

- outputs between 0 and 1

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



$$z = \vec{w} \cdot \vec{x} + b \rightarrow f_{\vec{w}, b}(\vec{x}) = g(\vec{w} \cdot \vec{x} + b) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}} \rightarrow \text{logistic regression} \rightarrow \begin{matrix} \text{used} \\ \text{for} \\ \text{classification} \end{matrix}$$

Interpretation of logistic regression output: $f_{\vec{w}, b}(\vec{x}) = P(Y=1 | \vec{x}; \vec{w}, b) \rightarrow$ Probability that y is 1, given input \vec{x} , Parameters \vec{w}, b

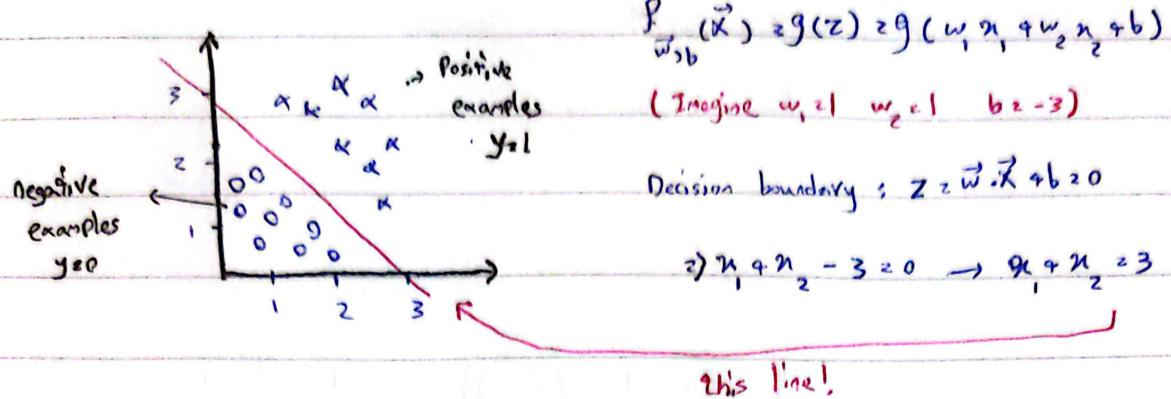
Decision boundary: We can set a threshold above which you predict $\hat{y}=1$, and below which you might say $\hat{y}=0$. A common choice $\rightarrow 0.5$

$$f_{\vec{w}, b}(\vec{x}) \geq 0.5 \text{ or } g(z) \geq 0.5 \text{ or }$$



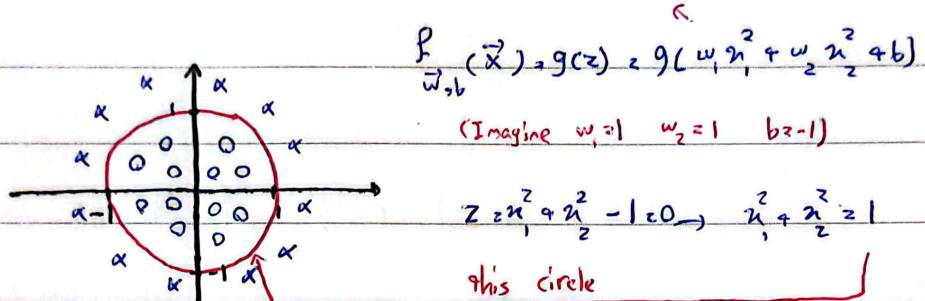
$$z \geq 0 \text{ or } \vec{w} \cdot \vec{x} + b \geq 0 \rightarrow \hat{Y}=1 \text{ or } \rightarrow \hat{Y}=0$$

e.g.: a classification problem with 2 features



Non-linear decision boundaries:

use Polynomials

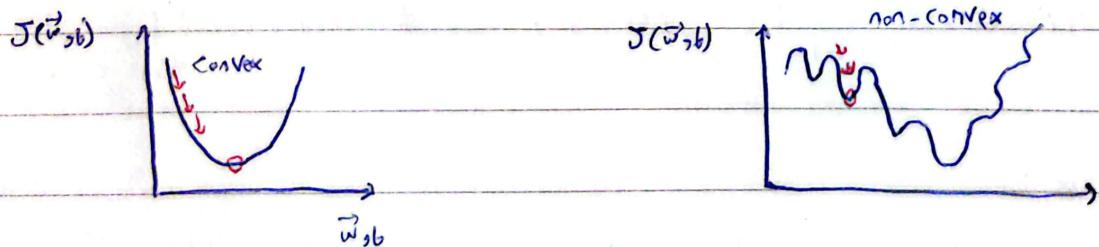


Cost function for logistic regression:

Squared error cost function is not good for logistic regression cause there are a lot of local minimas in the plot \rightarrow gd will do poorly. We need a func that would make the plot convex again.

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (P_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 \rightarrow \text{squared error}$$

linear regression logistic regression



* Loss func for squared error: $L(P_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \frac{1}{2} (P_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2$

Subject:

Date:

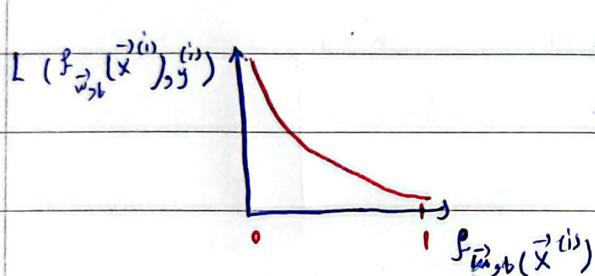
* Loss func \rightarrow measure how well you're doing on 1 training example

* Cost func = \sum loss func \rightarrow for all training examples

Logistic loss func :

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

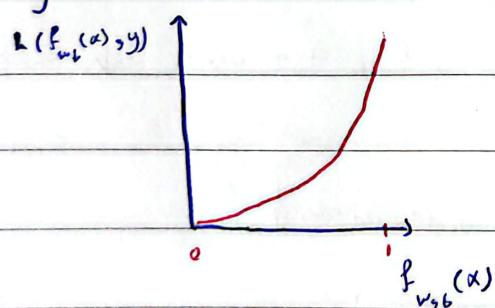
if $y^{(i)} = 1$



* As $f_{\vec{w}, b}(\vec{x}^{(i)}) \rightarrow 1$ then loss $\rightarrow 0$

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

if $y^{(i)} = 0$



* As $f_{\vec{w}, b}(\vec{x}) \rightarrow 0$ then loss $\rightarrow \infty$

* As $f_{\vec{w}, b}(\vec{x}) \rightarrow 1$ then loss $\rightarrow 0$

cost func : $J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})$

* Is convex: proof beyond scope of this course

Simplified loss func :

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

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

* we use this loss func cause it's related to Maximum likelihood

Gradient descent for logistic regression

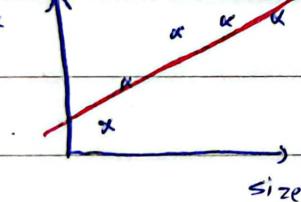
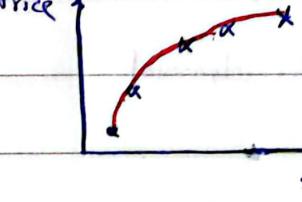
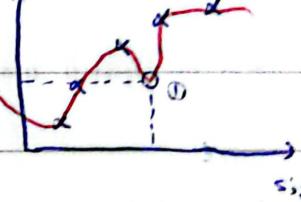
$$\text{cost: } J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1-y^{(i)}) \log(1-f_{\vec{w}, b}(\vec{x}^{(i)}))]$$

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

$$*\frac{\partial}{\partial b} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \quad *g'(z) = g(z)(1-g(z))$$

* Same concepts as for linear regression → learning curve - vectorize - Feature Scaling etc

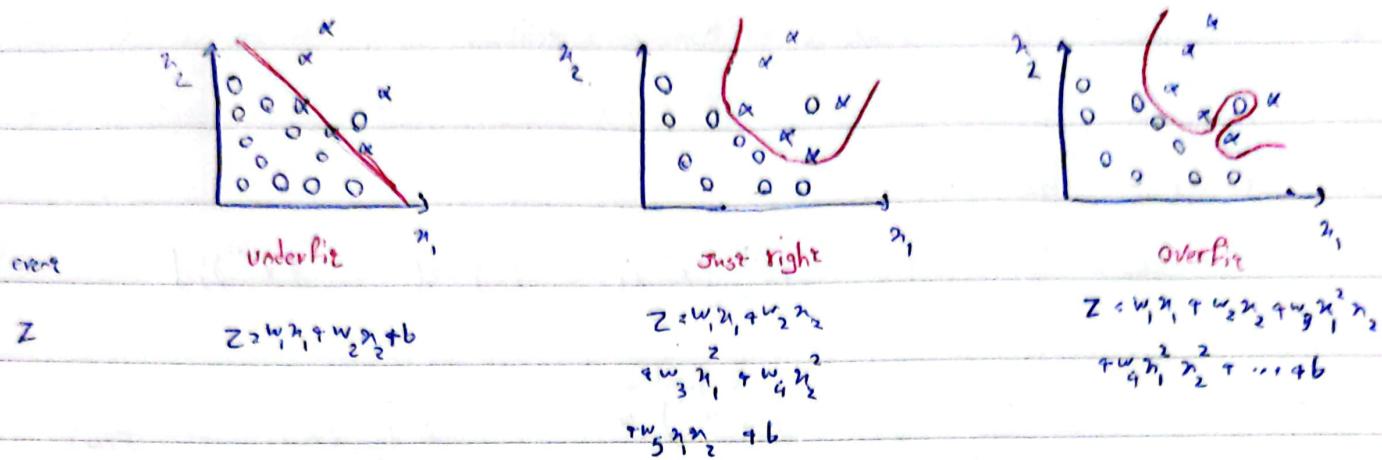
Underfitting & Overfitting - for linear regression:

| |  |  |  |
|----------------------------------|---|--|---|
| Event | Underfit | just right | overfit |
| e.g. model | $w_1 x + b$ | $w_1 x + w_2 x^2 + b$ | $w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$ |
| another name (Technical term) | (high bias) | (generalization) | (high variance) |
| explanation | Does not fit the training set well | Fits training set pretty well | Fits the training set extremely well |

* The middle model will do well on training examples that are not in the training set

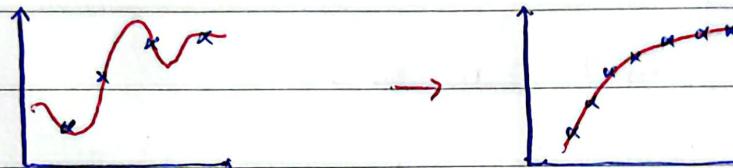
* The overfit model will have cost func = 0 but it's not good for prediction → e.g.: ①. So we can't expect this model will do well on examples that are not in the training set

For Classification:



Addressing overfitting is both for regression & classification

1 - Collect more training examples



* It would still have high order polynomial

2 - Feature Selection: try selecting and using only a subset of the features

* disadvantage: useful features could be lost

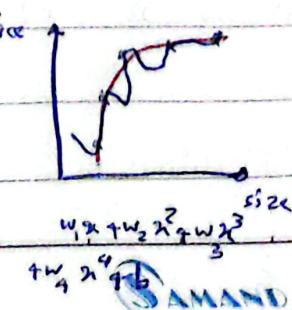
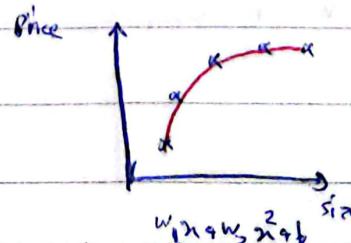
3 - Regularization: reduce the impact of some of the features without eliminating them (like in feature selection) \rightarrow reduce the associated w

$$\text{e.g.: } f(n) = 28n - 385n^2 + 39n^3 - 179n^4 + 100 \rightarrow 13n - 0.23n^2 + 0.000941n^3 - 0.0001n^4 + 10$$

* Reduce b doesn't make that much difference \rightarrow Andrew usually don't reduce it

Cost func with regularization

\rightarrow We need to make w_3 and w_4 really small



Date:

Subject:

★ We can do this with a modified cost func: $\frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 + 10^3 w_3 + 10^3 w_4$

$$\min J = \min (10^3 w_3 + 10^3 w_4) = \min (w_3, w_4)$$

★ we can use any big number instead of 10^3

In general: we don't know which w to penalize so we penalize all of them:

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 + \underbrace{\lambda \sum_{j=1}^n w_j^2}_{\text{regularization term}} + \underbrace{\frac{\lambda}{2m} b^2}_{①}$$

★ λ = regularization parameter \rightarrow we need to choose it like we choose learning rate

★ n = number of features $\quad \Phi$ penalize b ① can be excluded \rightarrow cause it has little impact
minimize

★ mean squared error term \rightarrow fit data well $\quad \Phi$ minimize regularization term \rightarrow keep w_j small
 \rightarrow λ balances both if we choose it wisely $\quad \Phi$ to reduce overfitting

$\lambda \rightarrow 0$ then no regularization \rightarrow overfit

★ if

$\lambda \rightarrow \infty$ then $w_j \rightarrow 0$ so $f_{w,b}(x) \approx b$ \rightarrow Underfit

Gradient descent for regularized linear regression

repeat {

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

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

we didn't regularize
 b

} simultaneous update

★ Some math intuition on how this regularized version shrinks w_j

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

Subject:

Date:

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

$w_j (1 - \alpha \frac{\lambda}{m})$ usual update
 $\text{like before regularization}$

$$\rightarrow 1 - \alpha \frac{\lambda}{m} . \text{ Imagine } \alpha = 0.01, \lambda = 1, m = 50 \rightarrow 1 - \alpha \frac{\lambda}{m} = 0.9998$$

$\rightarrow 0.9998 w_j \rightarrow$ in every iteration we are shrinking w_j just a little bit

For logistic regression:

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

★ gradient descent for this regularized version is the same for linear regression (regularized) and the only difference is that $f_{w,b}(x)$ is ^{the} sigmoid func