

- ① Run out of work in AI
- ② New capabilities for computers

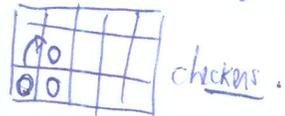
Some examples

- ① Database mining, large datasets, web automation, web click data, medical records, social networks
- ② Applications can't program by hand
 - autonomous car, helicopter
 - handwritten recognition
 - NLP
 - vision
- ③ Self-customizing
 - Amazon
 - Netflix product recommendations
- ④ Understanding human learning.

What is ML

① field of study that gives computers the ability to learn without being explicitly programmed. (Arthur Samuel)

② Tom Mitchell: well-posed learning program \Rightarrow a computer program is said to learn from exp (E) w.r.t task (T) and some perf. measure (P) if its perf. on (T) as measured by (P), improves with exp. (E).



checkers.

Types of ML

- SL vs SSL vs unsup. learning
- Rein learning, recommen. systems

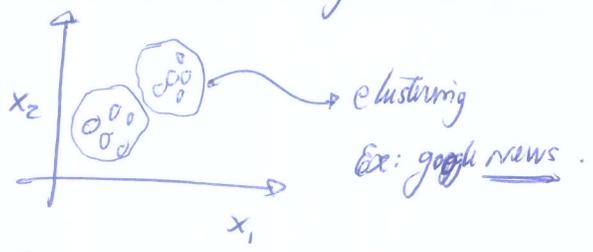
Practical advice: how to use and where to apply ML.

Amorim

Supervised Learning

- Regression \Rightarrow predict cont. output.
- classif \Rightarrow predict discrete value (class) output.

unsupervised Learning find some structure on the data



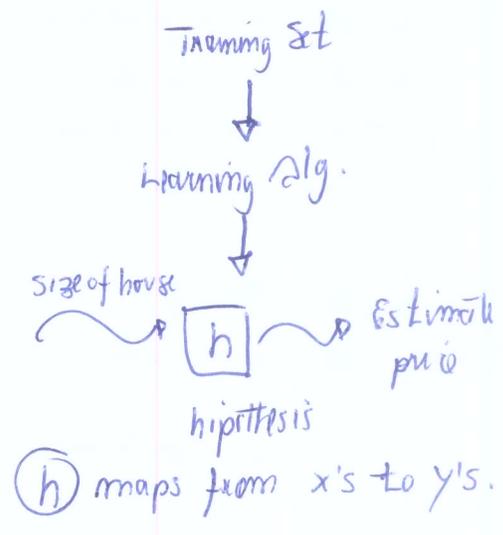
- Example: Genes x Individuals
- datacenters
 - cohesive groups of people on social networks
 - Group customers
 - Astronomical data analysis (how galaxies are formed)

Cocktail party problem \Rightarrow overlapping voices.

Language - Octave
 - matlab
 - R } first prototyping then \rightsquigarrow other languages.

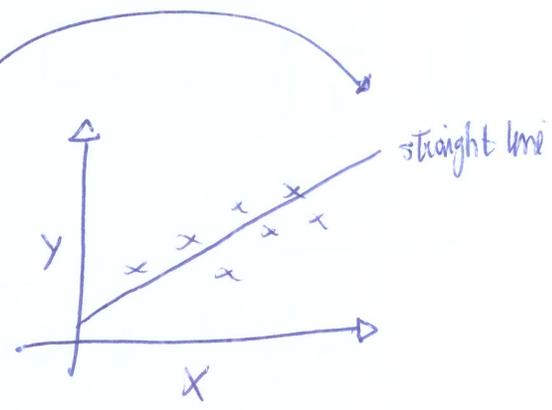
Notation

m : nbr of training examples
 x 's input variable/features
 y 's out. variable/target variable
 (x, y) one training example.
 $(x^{(i)}, y^{(i)})$ i th training example.



How to represent h ?

Initial choice $h_{\theta}(x) = \theta_0 + \theta_1 x$
 shorthand $h(x)$

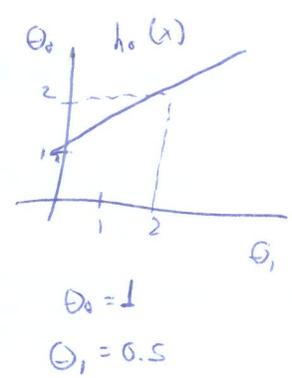
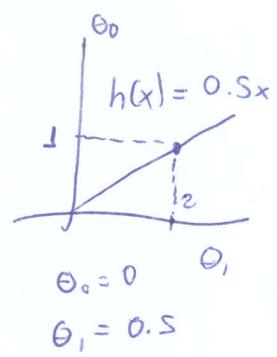
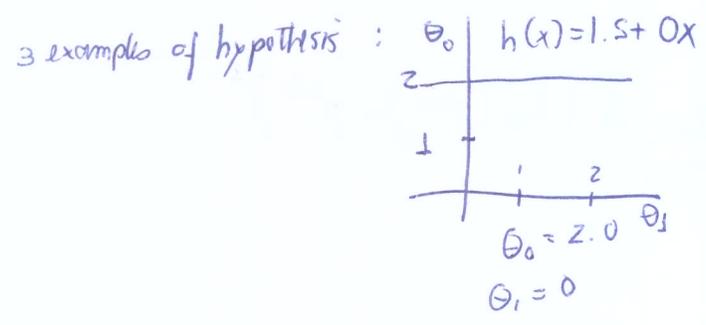


why linear functions \Rightarrow simple building block.

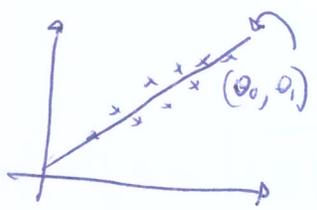
Here we have an example of a linear regression with one variable or univariate linear regression.
1 variable.

Defining a cost function

hypothesis $h_{\theta}(x) = \theta_0 + \theta_1 x$ θ s parameters



In linear regression



$h_{\theta}(x)$ must be close to the training data.

More formally

minimize $\theta_0, \theta_1 (h_{\theta}(x) - y)^2$

$$\min_{\theta_0, \theta_1} \left(\sum_{i=1}^M (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right)$$

$\frac{1}{2M}$ # Examples only for make it in the training easier.

Derive

$$\min_{\theta_0, \theta_1} \left(\frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right) = J(\theta_0, \theta_1).$$

$$h_{\theta}(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$$

with this, we define a cost function $J(\theta_0, \theta_1)$.

so we can simply say $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$ — cost function called squared error cost function. most used for Regression.
SSD

Cost function intuition

① hypothesis $h_{\theta}(x) = \theta_0 + \theta_1 x$.

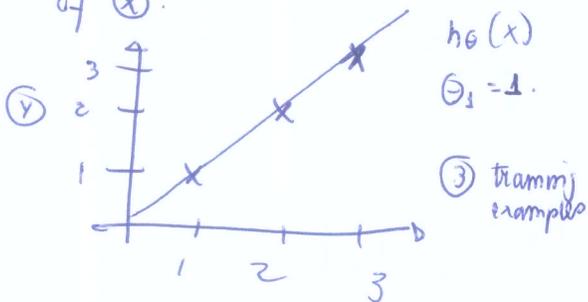
② Params θ_0, θ_1

③ Cost function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

④ Goal: minimize $J(\theta_0, \theta_1)$
 θ_0, θ_1

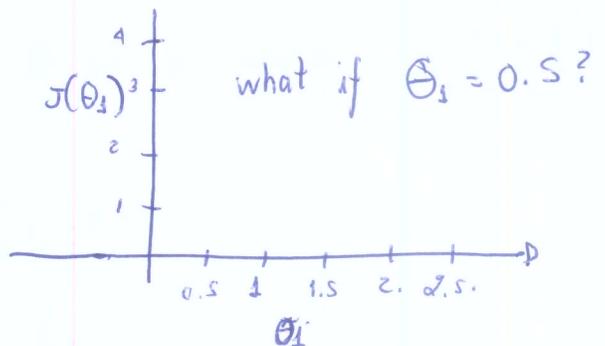
If we simplify $h_{\theta}(x) = \theta_1 x$. $\theta_0 = 0$.

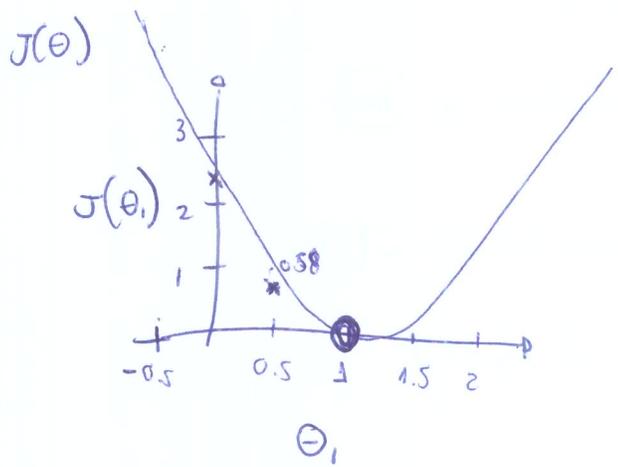
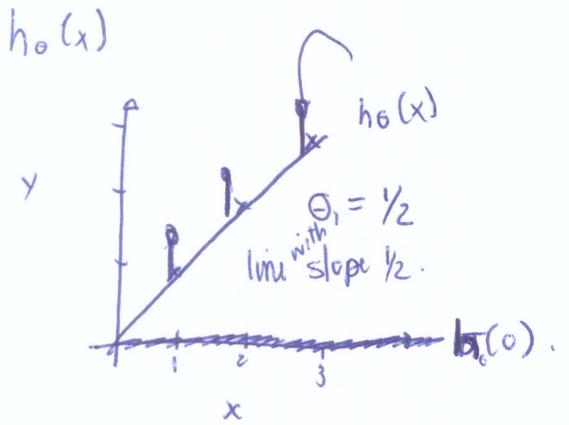
$h_{\theta}(x)$ for fixed θ_1 , this is a function of x .



what is $J(\theta_1)$? $\frac{1}{2m} (\theta_1 (x^{(1)} - y^{(1)}))^2 = \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0$

$J(\theta_1)$ is a function of the param θ_1 .





$$J(0.5) = \frac{1}{2M} \left[(0.5-1)^2 + (1-2)^2 + (1.5-3)^2 \right]$$

$$= \frac{1}{2 \times 3} (3.5) \approx 0.58$$

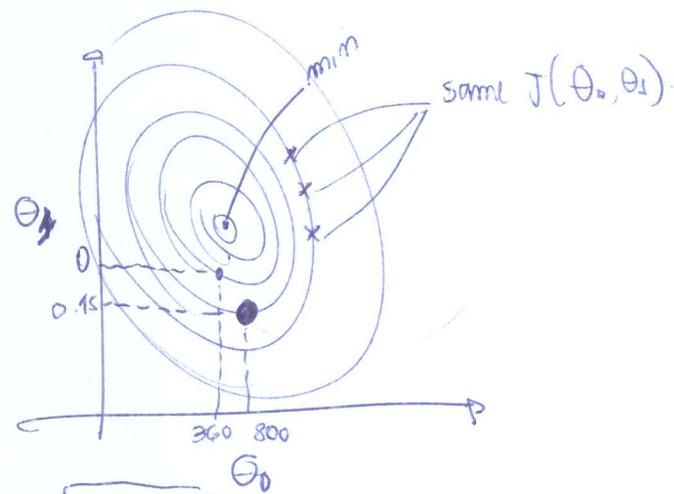
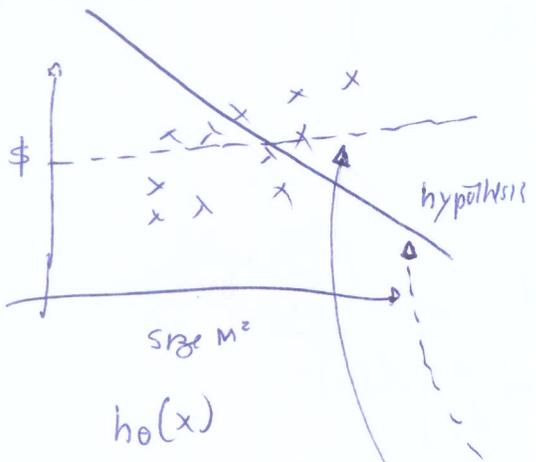
$J(0) = ?$

$$J(0) = \frac{1}{2m} (1^2 + 2^2 + 3^2)$$

$$\frac{1}{2m} \cdot 14 = 2.3$$

for each value of θ_1 , we have a different hypothesis for $h_0(x)$ (different line).

Cost function (contour plots)



θ_0	θ_1
800	1.5

$\theta_0 = 0$
$\theta_1 = 350$

Gradient Descent

Used everywhere in **ML**

1) Have some function: $J(\theta_0, \theta_1)$

2) $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

3) outline

(a) start with some θ_0, θ_1

(b) keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$ until we end up at a minimum.

what we need to define is α (the learning rate).

Implementation

partial derivatives

Correct form

$$\begin{aligned} \theta_0 &\leftarrow \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \theta_1 &\leftarrow \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \end{aligned}$$

$\left. \begin{matrix} \theta_0 \leftarrow \theta_0 \\ \theta_1 \leftarrow \theta_1 \end{matrix} \right\} \text{simultaneously update both } \theta_0 \text{ and } \theta_1$

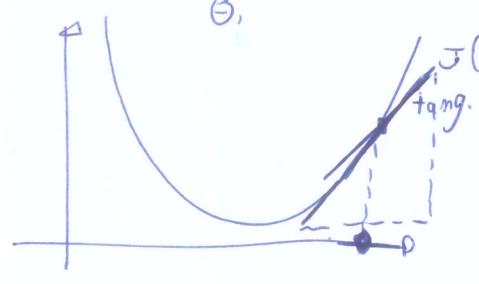
Incorrect form

$$\begin{aligned} \theta_0 &\leftarrow \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \theta_1 &\leftarrow \theta_1 \end{aligned}$$

⋮

Intuition about the Gradient Descent

Suppose we have $\min_{\theta_1} J(\theta_1)$ $\theta_1 \in \mathbb{R}$.



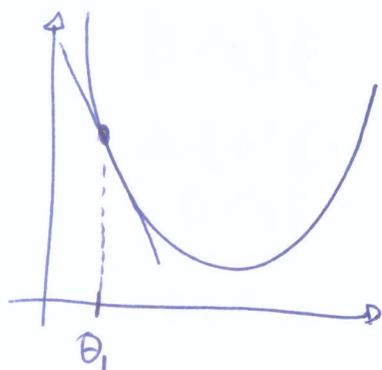
$\theta_1 = \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$

derivative $\left(\frac{\partial}{\partial \theta_1} \right)$ partial derivative.

slope of the line tangent to function.

≥ 0

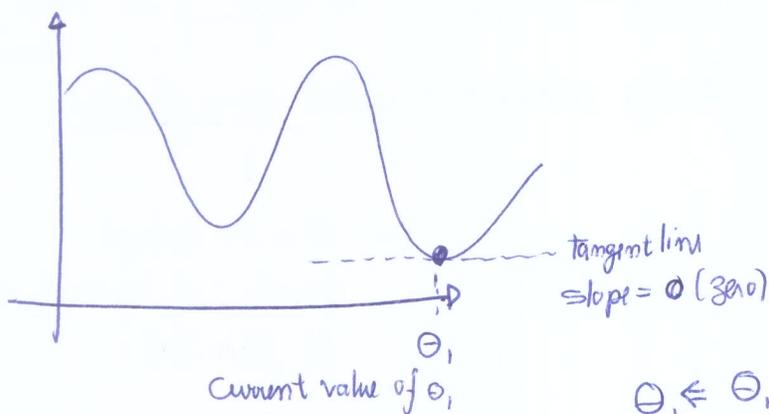
$$\theta_1 \leftarrow \theta_1 - \alpha \cdot (\text{positive number}) \downarrow \theta_1$$



$$\frac{dJ(\theta_1)}{d\theta_1} \leq 0.$$

$$\theta_1 \leftarrow \theta_1 - \alpha (\text{neg. number}) \uparrow \theta.$$

What if θ_1 is already at a minimum ?



$$\theta_1 \leftarrow \theta_1 - \alpha \cdot \frac{dJ(\theta_1)}{d\theta_1}$$

zero

$$\theta_1 \leftarrow \theta_1 - \alpha \cdot 0$$

$$\theta_1 \leftarrow \theta_1.$$

θ_1 stays unchanged that's why GD can converge to a local min even with α fixed.

GD automatically takes smaller steps as we approach a local min due to the derivative term. So, no need for decreasing α over time.

GD for linear regression

how to deriv $\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \cdot \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2.$

math

$$= \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \sum_{i=1}^m \left(\theta_0 + \theta_1 \cdot x^{(i)} - y^{(i)} \right)^2.$$

$$\frac{\partial}{\partial x} (x^2 - z)^2$$

$$2(x^2 - z) \cdot 2x$$

$$4x(x^2 - z).$$

$$\theta_0 \Rightarrow j=0: \frac{1}{m} \cdot \sum_{i=1}^m \left(h_{\theta}(x^{(i)}) - y^{(i)} \right).$$

$$\theta_1 \Rightarrow j=1: \frac{1}{m} \cdot \sum_{i=1}^m \left(h_{\theta}(x^{(i)}) - y^{(i)} \right) \cdot x^{(i)}.$$

~ The cost function for ~~GD~~ linear regression is always a bowl shape

~ The technical term for this is that it is called convex function



- is a bowl shape
- doesn't have local opt except for the global one.

The algorithm we just defined is what we call Batch Gradient Descent.

Batch: each step of ~~GD~~ uses all training examples.

~ there are other versions that use only a ~~few~~ subset of the training.

~ for the case of linear regression, there is a closed form solution (normal equations) but Gradient descent will scale better for larger datasets than normal equations.

~ The derivative is just the slope of the cost function J .

Generalization of Gradient Descent

Recalling, the GD algorithm is given by

$$\left. \begin{aligned} & \text{repeat until convergence} \\ & \theta_0 \leftarrow \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \\ & \theta_1 \leftarrow \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)} \end{aligned} \right\} \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_i}$$

→ It can be susceptible to local optima.

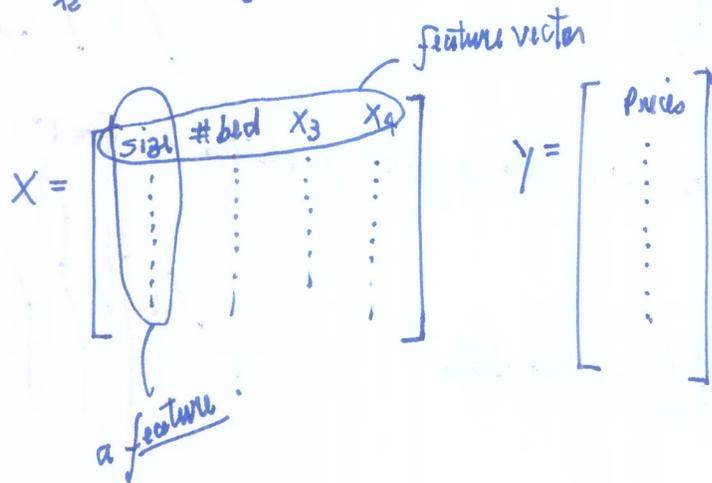
→ But the cost function for the linear regression is always convex (bow-shape)
→ only the global optimum.

There are two extensions of GD that we must consider

① In $\min J(\theta_0, \theta_1)$, solve for θ_0, θ_1 exactly, without iterative GD algorithm. One advantage is that α is not necessary anymore but in some cases (large high dim) the GD may be more appropriate.

② Longer ~~and~~ number of features. For instance
size, # bedrooms, # floors, # age of home \rightsquigarrow y (price).
 x_1 x_2 x_3 x_4

Notation



Increasing the number of features

we have considered so far linear functions of several observed ~~variables~~ features

x_1, x_2, \dots

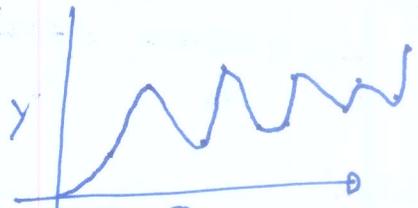
Suppose now that we have only one feature (x_1) but we would like our predictor to be a non linear function of (x) $\hat{y}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 \dots$
we can simply define new features $x_2 = x^2, x_3 = x^3$ just like we did with $x_0 = 1 = x^0$.

The predictor then becomes $\hat{y}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots$
It is : it still fits the linear regression model but in a new feature space with additional features that are deterministic functions of our observations.

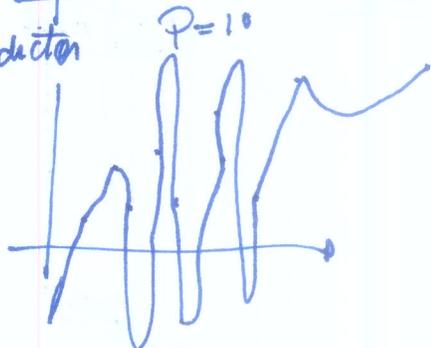
Applying the mse estimation $\hat{y}(x) = \sum_{i=1}^n \theta_i x_i$

Overfitting

In our polynomial fits for 11 data points, a more complex model $\hat{y}(x)$ a $p(x)$ with degree 10 and 11 coef. fits the data perfectly

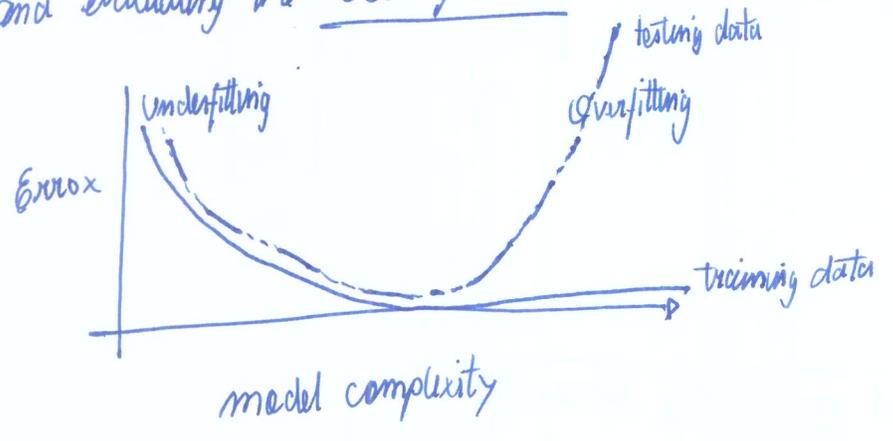


but it does not look like a good predictor



we fall on overfit to the data

We can see the generalization or test error simply by gathering more data and evaluating the cost function



Continuing Regression with multiple variables

size house x_1 , # bedrooms x_2 , floors x_3 , years x_4 , Price y

$x^{(i)}$ training example.

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

Previously we had a model/hypothesis: $\hat{y}(x) = \theta_0 + \theta_1 x$

Let's update or upgrade our model

$$\hat{y}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_4 x_4$$

Example $\hat{y}(x) = 80 + 0.1 x_1 + 0.01 x_2 + 3 x_3 - 2 x_4$
 each year decreases price.

for convenience let's define $x_0 = 1$.
 zeroth feature always with the value 1.

Then $x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$ $\theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$
 0-index 0-index vector

$h_{\theta}(x) = \hat{y}(x)$

$\hat{y}(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_m x_m$
 $= \theta^T x$

$$\begin{bmatrix} \theta_0 & \dots & \theta_m \end{bmatrix} \begin{bmatrix} x_0 \\ \vdots \\ x_m \end{bmatrix}$$

 $\theta^T x$

$\theta^T x$

This is called multivariate linear regression

multiple variables/features

How do we do the GD with multiple variables/features?

Hypothesis $\hat{y}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_m x_m$
Params $(\theta_0, \dots, \theta_m)$. let's think it as a vector $\vec{\theta} \in \mathbb{R}^{m+1}$

Cost function $J(\theta_0, \dots, \theta_m) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}(x^{(i)}) - y^{(i)})^2$
 $J(\vec{\theta})$

Gradient Descent Repeat
 $\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\vec{\theta})$
update at the same time for all $j \in \{1, \dots, m\}$

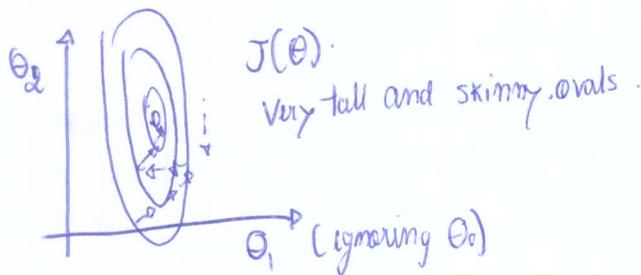
How to perform the partial derivatives?

GD for $m \geq 1$. Repeat
 $\theta_j \leftarrow \theta_j - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (\hat{y}(x^{(i)}) - y^{(i)}) x_j^{(i)}$ only thing that changes
 $x_0^{(i)} = 1$ by definition

Gradient Descent in practice 1: feature scaling

- ↳ what do we do if we have some features dominating the others?
- ↳ we need to scale the features in order GD can converge more quickly.
- ↳ we need to make sure the features are in the same range of values

Example
 $x_1 = \text{size (0-2000 m}^2\text{)}$
 $x_2 = \text{mbr of bedrooms (1-5)}$

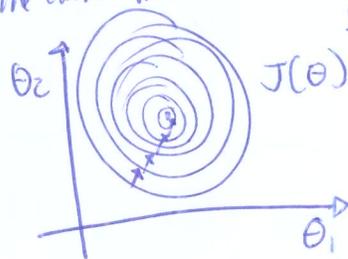


if we run GD in here, it may take a while before converging

A nice workaround here would be to normalize the data then the contours would look like more circles

$$x_1' = \frac{\text{size (m}^2\text{)}}{2000} \quad x_2' = \frac{\# \text{ bedrooms}}{5}$$

$$\left. \begin{array}{l} 0 \leq x_1 \leq 1 \\ 0 \leq x_2 \leq 1 \end{array} \right\}$$



Feature scaling ↳ get every feature into approximately a $-1 \leq x_i \leq 1$ range

↳ x_0 is already in the range ($x_0=1$)

↳ $0 \leq x_1 \leq 3$ (is OK)

$-2 \leq x_2 \leq 4$ (is OK)
but

$-20 \leq x_3 \leq 40$ X

$-0.0004 \leq x_4 \leq 0.0001$ X

↳ how to do it? training data

In addition to scale normalization, sometimes we perform a step further and do mean normalization (14)

- ~ Replace x_i with $x_i - \mu_i$ for zero meaning the features
- ~ Do not change $x_0 = 1$.

Example $x_1 = \frac{\text{size} - 1000}{2000}$

$x_2 = \frac{\# \text{ bedrooms} - 2}{S}$

if we do that, we will end up with values close enough to

$-0.5 \leq x_1 \leq 0.5, -0.5 \leq x_2 \leq 0.5$

General rule

$x_i \leftarrow \frac{x_i - \mu_i}{S_i}$

where μ_i is the mean of x_i in training set and S_i is the range of values in training set ($\max_i - \min_i$).

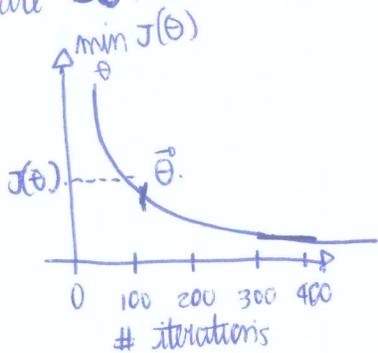
sometimes it will interesting to use the std of feature

(i) $x_i \leftarrow \frac{x_i - \mu_i}{\sigma_i}$

z-norm
z-score

Data standardization

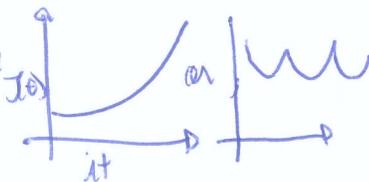
How do we know GD is working correctly and how to choose the learning rate α .



what you expect.

$J(\theta)$ should decrease after iteration.
 # of iterations vary according to the application.
 An automatic convergence would be possible, for instance if $J(\theta)$ in one iteration is only slightly lower (10^{-3} for instance) than $J(\theta)$ in the previous iteration. But this is difficult.

if you have something like



then choose a smaller α .



Reason of increase: big α overshoots the minimum.

Two general observations:

- ① for small enough α , $J(\theta)$ should decrease on every iteration...
- ② for small enough α or too small α , GD can be slow to converge.

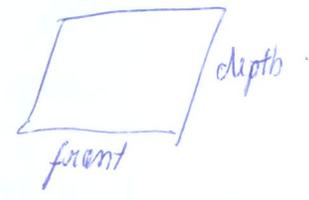
But... how to choose α ?

try... 0.001, 0.01, 0.1, 1, ... always plotting $J(\theta) \times \# \text{ iterations}$
 or
 0.001, 3×0.001 , 0.01, 3×0.01 , ..., ...
 or
 make it linked to the nbr of iteration $\frac{\alpha}{\text{iteration}}$.

Features and Polynomial Regression

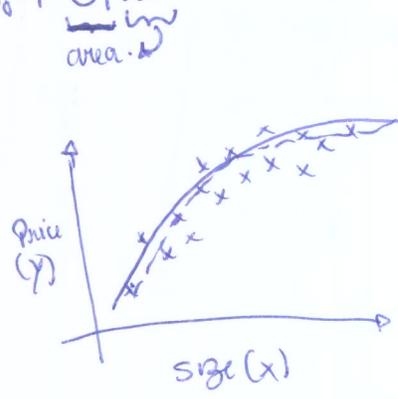
Housing prices $\hat{y}(x) = \theta_0 + \theta_1 \times \text{frontage} + \theta_2 \times \text{depth}$

we can define a new feature $x = \text{front} \times \text{depth}$ (area)



then $\hat{y}(x) = \theta_0 + \theta_1 x$
area

Polynomial Regression



$$\hat{y}_2(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

or

$$\hat{y}_3(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

for using our multivariate linear regression model, we just define new features

$$\hat{y}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

$$\hat{y}(x) = \theta_0 + \theta_1 x + \theta_2 x + \theta_3 x \text{ where } x_1 = \text{size}, x_2 = \text{size}^2 \text{ and } x_3 = \text{size}^3$$

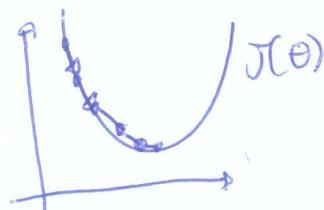
Note: if we do this, feature ~~is~~ scaling is **paramount**

- size: $1 - 10^3$
 - size²: $1 - 10^6$
 - size³: $1 - 10^9$
- mind*

How to choose the features? There are some algorithms for that.

The Normal Equations

Gradient Descent



The Normal Equation is a method for solving for θ analytically (one step).

minimize $\theta \in \mathbb{R}^{m+1}$ $J(\theta_0, \dots, \theta_m) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}(x^{(i)}) - y^{(i)})^2$

as we want to minimize $J(\theta)$, we solve for $\frac{\partial J(\theta)}{\partial \theta_j} = \dots$ set to 0 (derive and equal to zero) for every j .

solving for $\theta_0, \theta_1, \dots, \theta_m$.

For example: for a given problem, put the feature vector in a matrix X and the outcomes (target) in a vector y .

$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1400 & 2 & 3 & 40 \\ 1 & a & b & c & d \\ 1 & e & f & g & h \end{bmatrix}$
 $\mathbb{R}^{m \times (m+1)}$

$y = \begin{bmatrix} 900 \\ 230 \\ 500 \\ 200 \end{bmatrix}$
 \mathbb{R}^m

$\hat{\theta} = (X^T X)^{-1} X^T y$ gives the min.

General Case

m examples $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$, m features

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ \vdots \\ x_m^{(i)} \end{bmatrix} \in \mathbb{R}^{m+1}$$

we will design what we call a design matrix $X =$

$$X = \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \in \mathbb{R}^{m \times (m+1)}$$

what does it mean $\hat{\theta} = (X^T X)^{-1} X^T y$.

① $(X^T X)^{-1}$ is inverse of matrix $X^T X$.

when using normal equations, feature scaling is not necessary!

So, what should we use, GD or Normal Equations?

<u>GD</u>	<u>Normal Equation</u>
① Need to choose α .	① No α .
② Needs many iterations	② No iterations
③ works well with a high number of features m	③ Needs to compute $(X^T X)^{-1}$
	④ slow if m is very large.

** $X^T X \in \mathbb{R}^{m \times m}$, therefore normally it costs $O(m^3)$ for inverting on most naive implementations

What is large?

- m in hundreds, go with normal equation
- m in thousands ($< 5k$), OK with normal eq.
- $m > 10k$ definitely go with GD or some alternatives.

finally the normal equation will not work for more complex problems such as classification, that's why GD is important and should be always thought as one option.

Normal Equations and non-invertibility

18

↳ when computing $\vec{\Theta} = (X^T X)^{-1} X^T Y$, what if $(X^T X)^{-1}$ is non-invertible/singular/degenerate?

↳ R and Octave and Matlab have workarounds (robust) inverse functions called pseudo-inverse that does the right thing.

↳ normally there are two main causes for degeneration:

① Redundant features (linearly dependent)

↳ solution: dim. reduction
- feature deletion/selection

Example $x_1 = \text{size in m}^2$ $1 \text{ km} = 1000 \text{ meters}$
 $x_2 = \text{size in km}^2$
so $x_1 = (1000)^2 \cdot x_2$

② Too many features ($\# \text{ examples } (m) \ll \# \text{ features } (n)$)

↳ solution: - dim reduction
- dim/feature selection/deletion
- Regularization

Example $\left. \begin{array}{l} m = 10 \text{ examples} \\ n = 100 \text{ dimensions} \end{array} \right\} \text{params } \Theta \in \mathbb{R}^{10 \times 101}$

Organization of the class on August 14th, 2013

slide 9, 11, 12, 13, 14, 15, 16, 17, 18, 19.

Logistic Regression and Classification

- ① Our outcome now is discrete-valued.
- ② ~~LR~~ LogR is one of the most used learning algorithms.

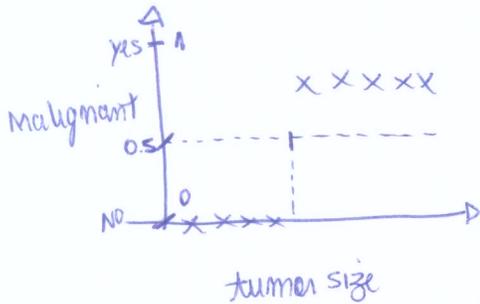
Some classification problems:

- spam x non-spam (e-mail)
- fraudulent x non-fraudulent (transactions)
- tumor malignant x benign
- young blind x not young blind (DR retinal)
- diab. retinopathy analysis

$y \in \{0, 1\}$

- positive class
- neg. class
- absence of something

binary classif. problem.



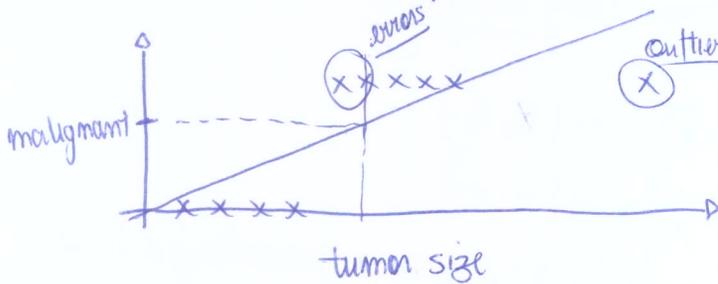
what if we approach the problem as a linear regression one?

$\hat{y}(x) = \Theta^T x$

Then we could threshold classifier's output at $\hat{y}(x)$ @ 0.5:

- if $\hat{y}(x) \geq 0.5$ predict ①
- otherwise predict ②

But let's change the problem:



often linear regression ~~is~~ isn't a good idea for classification.

Another problem with ~~LR~~ LogR linear regression is that even with training examples with $y \in \{0, 1\}$ it can output $\hat{y}(x) > 1$ or $\hat{y}(x) < 0$.

LogR solves this problem such that $0 \leq \hat{y}(x) \leq 1$ always.

we will use it as a classif. algorithm. The name is historical.

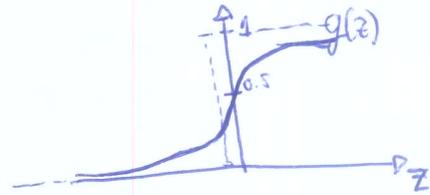
LogR

Hypothesis Representation

LogR model want $0 \leq \hat{y}(x) \leq 1$.

when we had LR, $\hat{y}(x) = \Theta^T x$, for LogR, we change it to $\hat{y}(x) = g(\Theta^T x)$,
Linear Regression

where $g(z) = \frac{1}{1 + e^{-z}}$
Sigmoid or logistic function
Synonyms.



By doing that, my new predictor will be

$$\hat{y}(x) = \frac{1}{1 + e^{-(\Theta^T x)}}$$

Now we need to fit parameters for Θ .

Interpretation of hypothesis output

$\hat{y}(x)$ = estimated prob that $y=1$ on input x .

Example $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$

Suppose my outcome is $y(x) = 0.7$

Tell patient that 70% chance of tumor being malignant.

mathematically, we have $\hat{y}(x) = P(y=1 | x; \Theta)$
 $= P(y=1 | x; \Theta)$

prob. of $y=1$ given x , parameterized by Θ .

Properties $\left\{ \begin{array}{l} P(y=0 | x; \Theta) + P(y=1 | x; \Theta) = 1 \\ P(y=0 | x; \Theta) = 1 - P(y=1 | x; \Theta) \end{array} \right.$

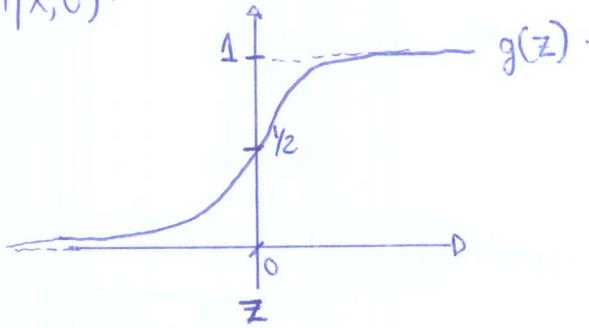
* y must take values in $\{0, 1\}$ binary value only.

Decision Boundary

Log R

$$\hat{y}(x) = g(\vec{\theta}^T x)$$

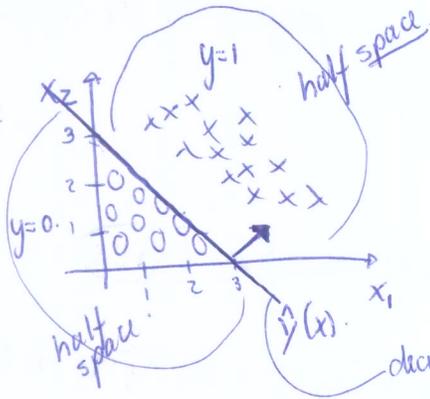
$$g(z) = \frac{1}{1 + e^{-z}}$$



Suppose predict $y=1$ if $\hat{y}(x) \geq 0.5$
 $y=0$, otherwise.

if we look at the sigmoid plot, we see $g(z) \geq 0.5$ when $z \geq 0$.
 then given our predictor $\hat{y}(x) = g(\vec{\theta}^T x)$, it will be ≥ 0.5 when $\vec{\theta}^T x \geq 0$.

Suppose we have training set



$$\hat{y}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

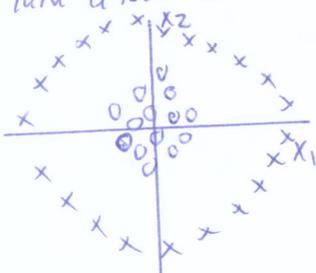
We still don't know how to choose θ s but

Suppose $\begin{cases} \theta_0 = -3 \\ \theta_1 = 1 \\ \theta_2 = 1 \end{cases} \quad \vec{\theta} = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$

Predict $y=1$ if $\underbrace{-3 + x_1 + x_2}_{\vec{\theta}^T x} \geq 0$ rewriting $\boxed{x_1 + x_2 \geq 3}$
 $\rightarrow y(x) = 0.5$ exactly.

The decision boundary is a property of the hypothesis including the params $\theta_0 \dots \theta_2$ and not of the dataset.

Let's take a look at a more complex example



how can we fit Log R params to solve this problem?

We can transform the feature space to higher order just like we did with polynomials.

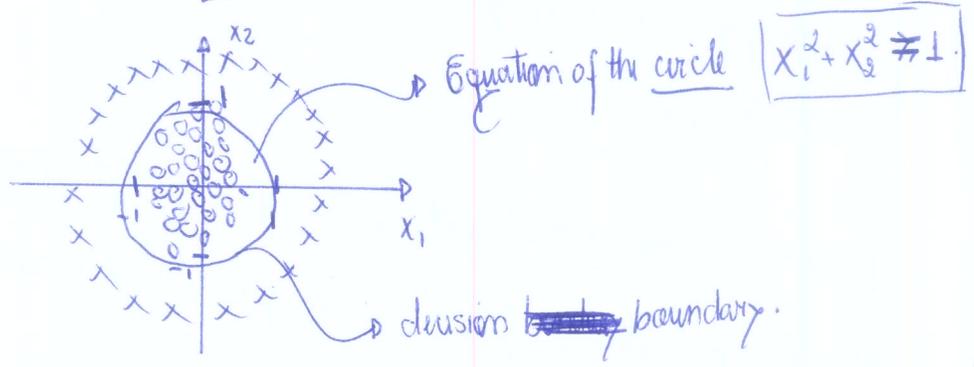
$$\hat{y}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 \underbrace{x_1^2}_{\text{added features}} + \theta_4 \underbrace{x_2^2}_{\text{added features}})$$

polynomial

Suppose we have found $\vec{\theta} = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$.

My predictor will predict $y=1$ if $-1 + x_1^2 + x_2^2 \geq 0$

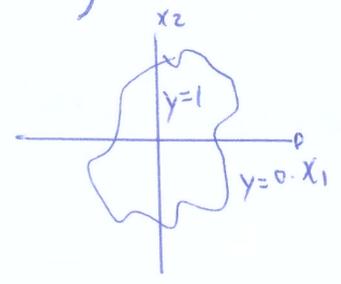
$x_1^2 + x_2^2 \geq 1$



Once again, the decision boundary is a property of the hypothesis/model and not of the dataset

Can we have even more complex decision boundaries? yes.

$\hat{y}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^2 x_2^2 + \theta_6 x_1^3 x_2 \dots)$



_____ x _____

But, how do we automatically fit the params for a model/predictor? (23)

training set $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

(m) examples $x \in \begin{bmatrix} x_0 \\ \vdots \\ x_n \end{bmatrix}_{\mathbb{R}^{n+1}}$, $x_0 = 1$, $y \in \{0, 1\}$.
 classif problem.

$$\hat{y}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

hypothesis

Question: how to choose parameters θ ?

Recall that in the linear Regression hypothesis, we had $J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (\hat{y}(x^{(i)}) - y^{(i)})^2$

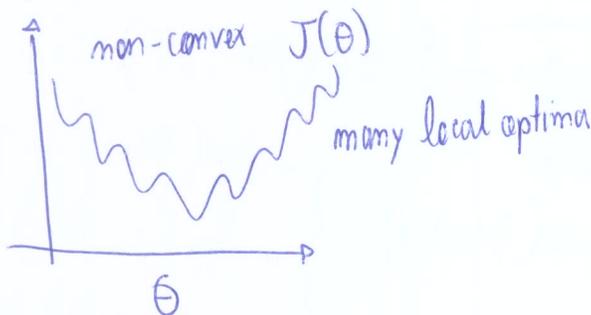
for LogR, let's call $\frac{1}{2} (\hat{y}(x^{(i)}) - y^{(i)})^2$ as cost $(\hat{y}(x^{(i)}), y)$.

$$\text{So Cost}(\hat{y}(x^{(i)}), y^{(i)}) = \frac{1}{2} \cdot (\hat{y}(x^{(i)}) - y^{(i)})^2$$

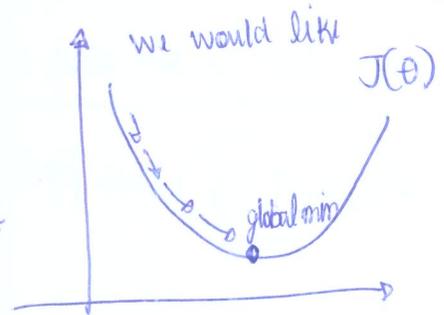
Simplifying,

$$\text{Cost}(\hat{y}(x), y) = \frac{1}{2} \cdot (\hat{y}(x) - y)^2$$

this cost function works fine for LR but we are focused on LogR and if we minimize it for LogR, it will be clear it isn't a convex function.



for LogR has non linearity $\frac{1}{1+e^{-\theta^T x}}$ complicated non linear function

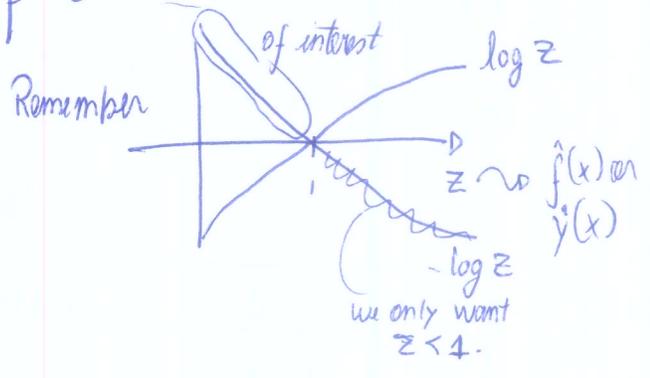
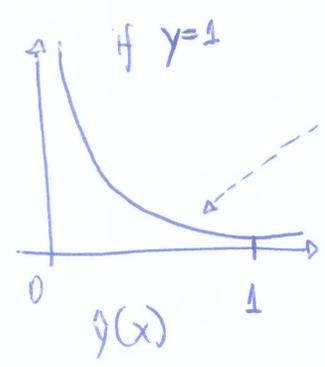


So we need to come up with a way of designing a different cost function that is convex.

Logistic Regression Cost function

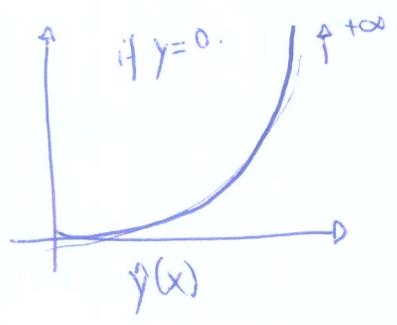
$$\text{Cost}(\hat{f}(x), y) = \text{Cost}(\hat{y}(x), y) =$$

$$\left. \begin{aligned} & -\log(\hat{y}(x)) \text{ if } y=1 \\ & -\log(1-\hat{y}(x)) \text{ otherwise } (y=0) \end{aligned} \right\}$$

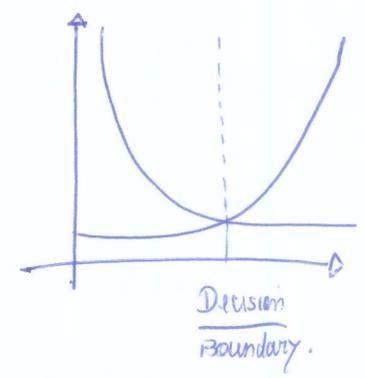


and prediction is $\hat{y}(x)=1$
 \sim Cost = 0 if $y=1, \hat{y}(x)=1$
 But as $\hat{y}(x) \rightarrow 0$, Cost $\rightarrow \infty$.

we penalize the learning alg. by a very large cost as it commits mistakes (go further away from the correct value/output).



putting together



simplified Cost function and Gradient Descent

~~Log~~ (L Reg) cost function

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(\hat{y}(x^{(i)}), y^{(i)})$$

where

$$\text{Cost}(\hat{y}(x), y) = \left\{ \begin{aligned} & -\log(\hat{y}(x)) \text{ if } y=1 \\ & -\log(1-\hat{y}(x)) \text{ otherwise} \end{aligned} \right.$$

$y \in \{0, 1\}$ binary problem

Let's compare the two cases in one.

$$\text{Cost}(\hat{y}(x), y) = -y \log(\hat{y}(x)) - (1-y) \log(1-\hat{y}(x)).$$

So now, we can design our cost function as



$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(\hat{y}(x^{(i)}), y^{(i)})$$

$$= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \cdot \log(\hat{y}(x^{(i)})) + (1-y^{(i)}) \cdot \log(1-\hat{y}(x^{(i)})) \right]$$

We choose it because it can be derived from statistics using principle of MLE max. likelihood estim. and it is convex.

~ To fit params $\vec{\theta}$: $\min_{\vec{\theta}} J(\vec{\theta})$.

~ To make a prediction for a new x , Output $\hat{y}(x) = \frac{1}{1 + e^{-\vec{\theta}^T x}}$ which means $P(y=1 | x; \vec{\theta})$ learned

Using GD for minimizing $J(\vec{\theta})$.

want $\min_{\vec{\theta}} J(\vec{\theta})$ Repeat $\theta_j \leftarrow \theta_j - \alpha \cdot \frac{\partial J(\theta)}{\partial \theta_j}$

$$** \frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (\hat{y}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

what? it look ~~like~~ exactly the same as before for LR. But the ~~change~~ change is now $\hat{y}(x)$ before $\hat{y}(x) = \theta^T x$. now $\hat{y}(x) = \frac{1}{1 + e^{-\theta^T x}}$ instead

In a vectorized form:

$$\begin{bmatrix} \theta_0 \\ \vdots \\ \theta_m \end{bmatrix}_{m \times 1} \leftarrow \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_m \end{bmatrix}_{m \times 1} - \alpha \left(\begin{bmatrix} z(X) \\ \vdots \\ y_m \end{bmatrix}_{m \times 1} - \begin{bmatrix} y_0 \\ \vdots \\ y_m \end{bmatrix}_{m \times 1} \right) \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{m \times 1}$$

$m \times m$ $m \times 1$ $m \times 1$ $m \times m$

No some adjustments are necessary.

$$\begin{bmatrix} f(A) \\ \vdots \\ \vdots \end{bmatrix}_{m \times m} \begin{bmatrix} B \\ \vdots \\ \vdots \end{bmatrix}_{m \times 1} = \begin{bmatrix} C \\ \vdots \\ \vdots \end{bmatrix}_{m \times 1}$$

$$\begin{bmatrix} f(A) \\ \vdots \\ \vdots \end{bmatrix}_{m \times 1} - \begin{bmatrix} Y \\ \vdots \\ \vdots \end{bmatrix}_{m \times 1} = \begin{bmatrix} C \\ \vdots \\ \vdots \end{bmatrix}_{m \times 1}$$

Logistic Regression - advanced optimization

So far, we have seen the opt. alg Gradient Descent and with it, we have a cost function $J(\theta)$ and want $\min_{\theta} (J(\theta))$.

Given θ , we compute

- ① $J(\theta)$
- ② $\frac{\partial J(\theta)}{\partial \theta_j} \quad \forall j=0, \dots, n_{\text{features}}$

and GD does:

Repeat:

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

but this is one alternative. If we show how to compute $J(\theta)$ and $\frac{\partial J(\theta)}{\partial \theta_j}$, we can use more sophisticated solutions such as

- Conjugate Gradient
- BFES
- L-BFGS.

Normally, such alternatives

- ① Do not need α (learning rate) manually selected.
- ② often faster than GD.

however, they are more complex.

Example on how to use them:

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \quad J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$

of course $\theta = \begin{bmatrix} 5 \\ 5 \end{bmatrix}$ minimizes it and that's what we want to find.

$$\frac{\partial J(\theta)}{\partial \theta_1} = 2(\theta_1 - 5)$$

$$\frac{\partial J(\theta)}{\partial \theta_2} = 2(\theta_2 - 5)$$

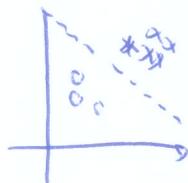
Normally, all we need to do is to write code for computing $J(\theta)$ and its derivatives.

Logistic Regression on multi-class problems - One vs All.

Exempli: digits, faces $\left\{ \begin{array}{l} \text{profile left} \\ \text{profile right} \\ \text{frontal} \end{array} \right.$ ag-group estimation, foldering/tagging photos, emails,

foldering $\left\{ \begin{array}{l} \text{family} \\ \text{business} \\ \text{gym} \end{array} \right.$ $y \in \{1, \dots, k\}$ classes.

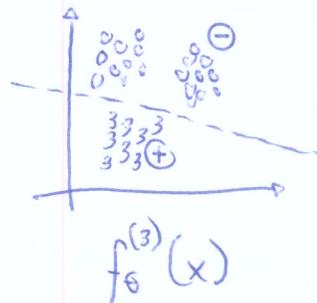
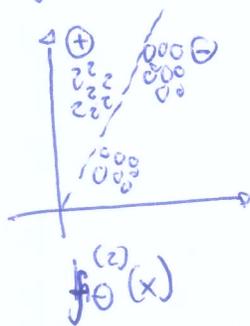
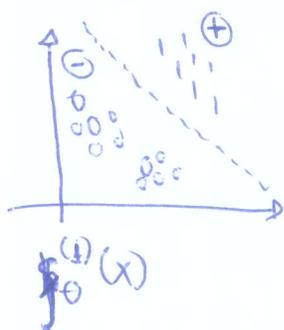
Binary Classif.



Multi-class



Let's do it one vs all (OVA): first we transform our problem into 3 problems.



what we did was ~~train~~ $f_{\theta}^{(i)}(x) = P(y=i|x;\theta)$ $i \in \{1, \dots, 3\}$.

Basically, we train a (LogR) classifier $f_{\theta}^{(i)}(x)$ for each class i to predict $P(y=i)$. For a new input x , for predicting, we select the class that

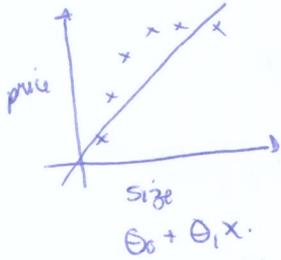
maximizes $\max_i \hat{y}_i$



Regularization and the problem of overfitting

Regularization will allow us to diminish overfitting problems.

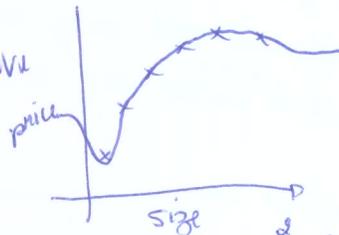
Lets return to the problem of housing



underfit, high bias

both mean the model is not fitting the data very well.

alternatively, we could have

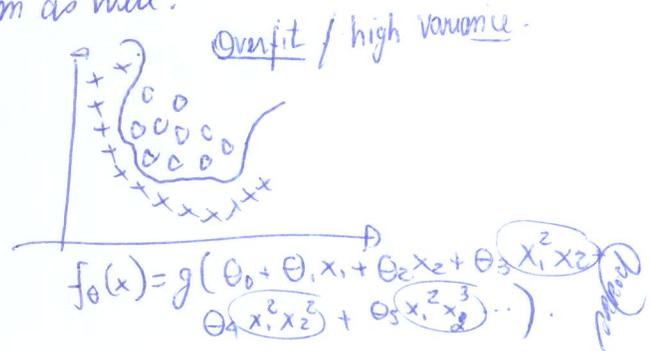
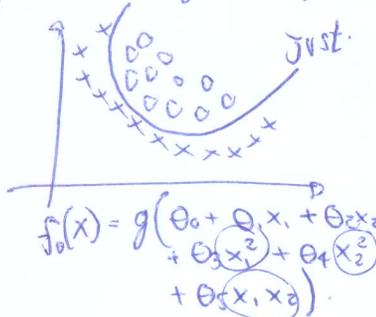
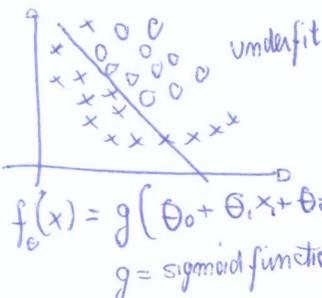


Overfit, "high variance"

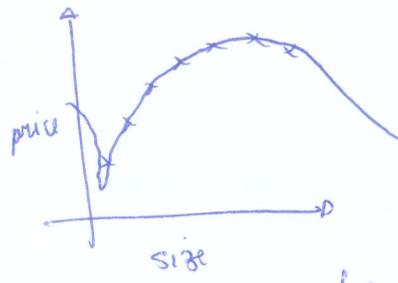
the model fitted the data perfectly and it is just to it and the hypothesis is too variable since there is not enough data to prove it

overfitting: if we have too many features, the learned hypothesis may fit ~~the~~ the training set very well ($J(\theta) = \frac{1}{2m} \sum_{i=1}^m f_{\theta}(x^{(i)}) - y^{(i)} \approx 0$) but fail to generalize to new examples (predict prices on new data).

The same thing can happen with logistic regression as well.



Addressing Overfitting



Plotting is one way but it wouldn't work for high-dimensional data

Two main options to deal with it

1 Reduce number of features

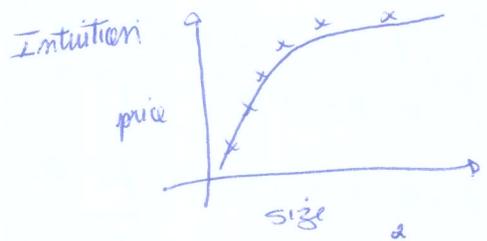
- manually which features to keep.
- model selection algorithm.
- automatic alg. that select which features to keep/throw away.

* however, sometimes this means you are throwing away important info about the problem.

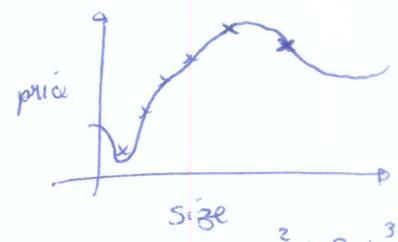
2 Regularization

Keep all the features but reduce magnitude/values/importance of parameters θ_j . works well ~~with~~ ~~many~~ when we have lots of features, ~~and~~ each ^{of which} contributing a bit for predicting y .

Regularization - Cost function



$$\theta_0 + \theta_1 x + \theta_2 x^2$$



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Suppose we penalize and make θ_3 and θ_4 really small

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (f_{\theta}(x^{(i)}) - y^{(i)})^2 + \boxed{1000 \theta_3^2 + 1000 \theta_4^2}$$

high penalization

this would give $\theta_3 \approx 0$
 $\theta_4 \approx 0$

Regularization Cost function

Besides the intuition, here's the general ideal:

- no small values for parameters $\theta_0, \theta_1, \dots, \theta_m$
 - Simpler hypothesis
 - less prone to overfitting.

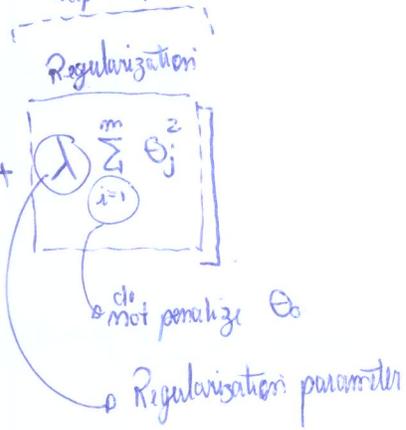
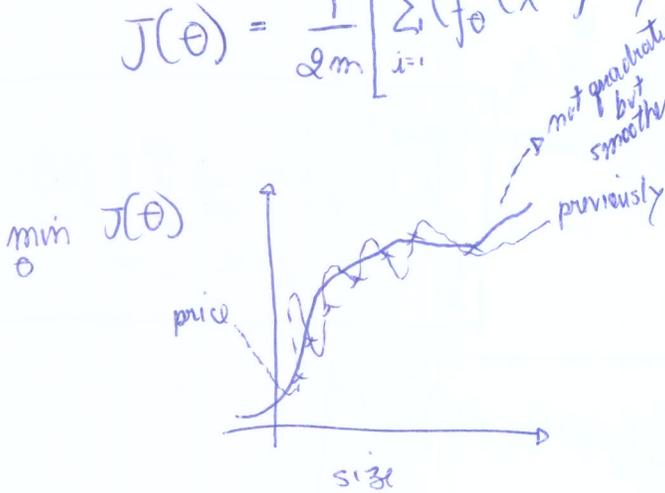
Let's see the housing example

- features: x_1, x_2, \dots, x_{100}
- params: $\theta_0, \theta_1, \dots, \theta_{100} \in \mathbb{R}$

it is difficult to pick one or some to penalize beforehand. If it was easy, it would be the same as selecting the less important features! keep the params small.

So, here's what we do:

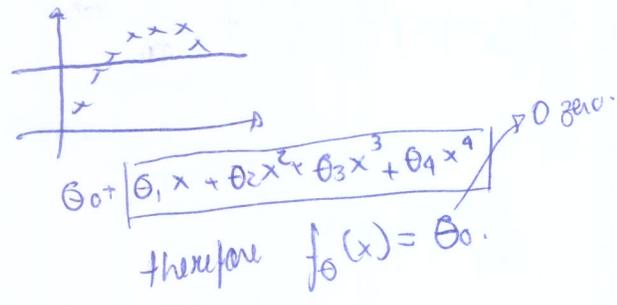
$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (f_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^m \theta_j^2 \right]$$



λ controls the tradeoff between good fitting to training data

vs
Keep the parameters small.
and with both we want to keep the hypothesis simple and avoid overfitting.

Question: what if λ is too high (eg. $\lambda = 10^{10}$) in a regularized linear regression?
 ↳ all θ_s except θ_0 will tend to zero
 and the model will be biased



underfitting
 too strong bias/pre-conception
 the a straight line in θ_0 will fit well in spite of data on the contrary.

Be carefull when choosing λ .

Regularized Linear Regression

So far, for linear regression, we have seen GD and Normal Equations for fitting the parameters.

Here's the optimization function for regularized linear regression

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (f_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^m \theta_j^2 \right]$$

\Downarrow
 $\min_{\theta} J(\theta)$

Previously, the GD for minimizing this was

Repeat \downarrow

$$\theta_j \leftarrow \theta_j - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (f(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$j = 0, \dots, m.$

}

if we separate θ_0 :

Repeat \downarrow

$$\theta_0 \leftarrow \theta_0 - \frac{\alpha}{m} \sum_{i=1}^m (f(x^{(i)}) - y^{(i)}) x_0^{(i)}$$
$$\theta_j \leftarrow \theta_j - \frac{\alpha}{m} \sum_{i=1}^m (f(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

}

$\frac{\partial J(\theta)}{\partial \theta_0}$

Repeat \downarrow

$$\theta_0 \leftarrow \theta_0 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (f(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)}$$
$$\theta_j \leftarrow \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (f(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} + \frac{\lambda}{m} \cdot \theta_j \right]$$

} regularization

$\frac{\partial J(\theta)}{\partial \theta_j}$ regularized

We can regroup things and then

$$\Theta_j \leftarrow \Theta_j - \alpha \cdot \left[\frac{1}{m} \sum_{i=1}^m (f(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} + \frac{\lambda}{m} \cdot \Theta_j \right]$$

$$\Theta_j \leftarrow \underbrace{\Theta_j \left(1 - \frac{\alpha \lambda}{m} \right)}_{< 1} - \frac{\alpha}{m} \sum_{i=1}^m (f(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

small small
target

for instance $1 - \frac{\alpha \lambda}{m} < 1$
 $1 - 0.01 = 0.99 \cdot \Theta_j$ shrinks Θ_j a little bit.

what about normal equations?

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \begin{matrix} \text{one training example} \\ \\ \\ \end{matrix}$$

$(m+1) \times (n+1)$

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \mathbb{R}^m$$

$$\min_{\theta} J(\theta)$$

$$\vec{\theta} = (X^T X)^{-1} \cdot X^T y \text{ without regularization}$$

for regularization;

$$\vec{\theta} = \left(X^T X + \lambda \begin{bmatrix} 0 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \\ & & & & 0 \end{bmatrix} \right)^{-1} \cdot X^T y$$

zeros zeros
(m+1) x (m+1)

Example $m=2$

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

non-invertible/singular.

what about non-invertibility now? (Advanced)

Suppose $m \leq n$
exam # features

$$\vec{\theta} = (X^T X)^{-1} \cdot X^T y$$

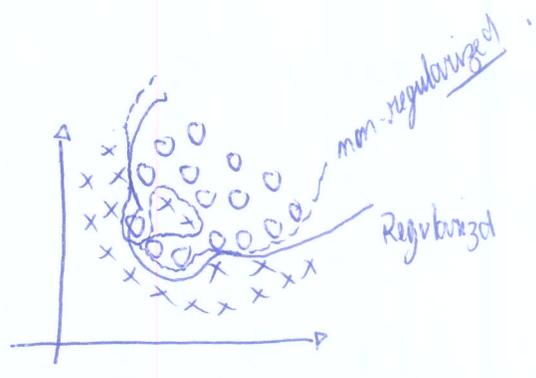
but

$$\text{if } \lambda > 0, \vec{\theta} = \left(X^T X + \lambda \begin{bmatrix} 0 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix} \right)^{-1} \cdot X^T y$$

we'll have no singularity anymore. why some features will decrease magnitude and simplify model.

Regularized Logistic Regression

RegL can ~~be~~ also be prone to overfitting:



$$\hat{f}_\theta(x) = g \left(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 \cdot x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \dots \right)$$

Cost function

$$J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m x^{(i)} (\log \hat{f}(x^{(i)})) + (1 - y^{(i)}) \cdot \log (1 - \hat{f}(x^{(i)})) \right] + A$$

introducing Regularization

+ dimensions

We just add

$$+ \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$

for $\theta_1, \dots, \theta_m$.

A

How to implement this?

Gradient Descent

Repeat {

$$\theta_0 \leftarrow \theta_0 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (f(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)}$$

$$\theta_j \leftarrow \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (f(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} + \frac{\lambda}{m} \theta_j \right]$$

$j \in \{1, \dots, m\}$.

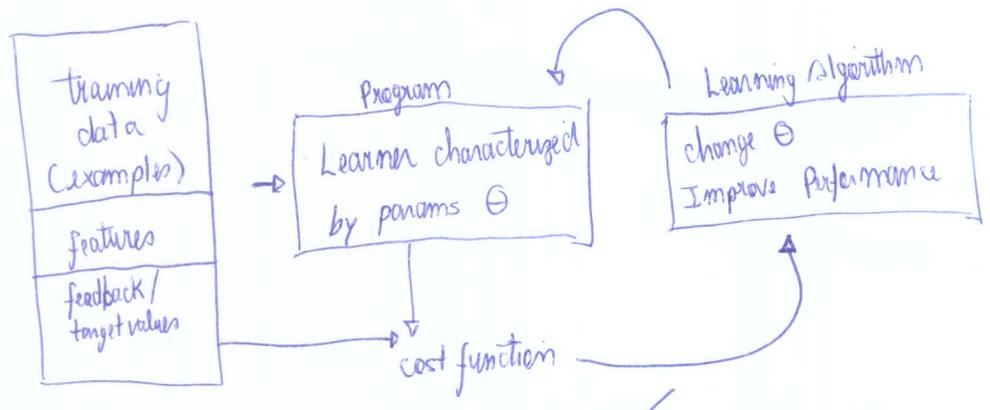
}

$\frac{\partial J(\theta)}{\partial \theta_j}$

remember $f_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$

Wrapping up what we had so far

Notation
 features x
 targets y
 predictions \hat{y}
 params θ



In the case of linear regression,



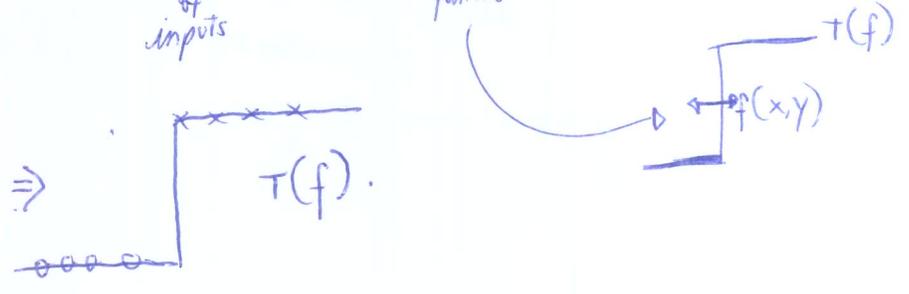
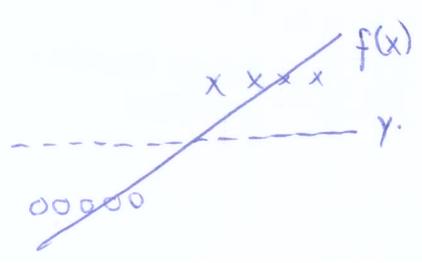
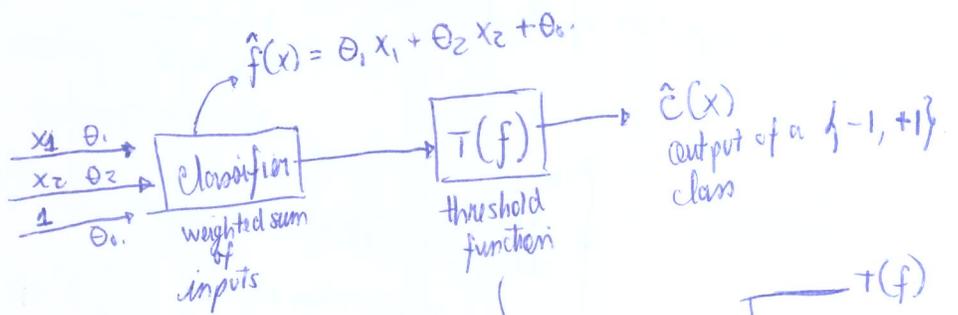
"Predictor"
 Evaluate line
 $\mu = \theta_0 + \theta_1 x_1$
 $\hat{y}(x) = \theta_0 + \theta_1 x_1$

$\mathcal{Y} \in \mathbb{R}$

while with classification, we seek to predict a discrete value/target \mathcal{Y} .

Now that we already saw linear regression and logistic regression (classifier) let's see a similar linear algorithm for classification ~~is~~ called perceptron.

Perceptron with 2 features



Done