

يلا علي الريجريشن .. اللي هنععمل انهارده .. اننا هنفكر ثاني في اللينير ريجريشن .. وديه خطوه جميله الحقيقه لاننا هنفكر فجه كإننا بنحل أوبتيمائيزيشن بروبليم .. اللي هنعملو انهارده اننا هنقعد نترنح ما بين اننا ندور علي طريقه ال optimization و بعدين نشوف ازاى نطبقها علي ال linear regression .. وهنعمل الترنح ده حوالي 3 او 4 مرات هنفضل نترنح كدا لحد ما المحاضره تخلص .. طب افرض المحاضره مخلصتش .. هنفضل نترنح برضو .. يلا بينا .. هومورك 3 نزل يا حلاوه .. هه

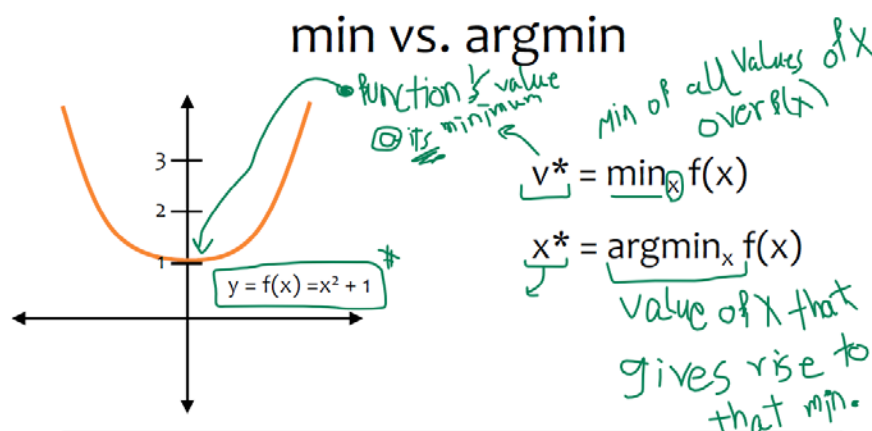
آخر مره اتكلمنا علي التعريف بتاع ال linear function و ازاى ال linear regression بيقلل ال residuals ما بين البريديشن و الأوبزرفيشن .. تعال بس الاول نقول كام حاجه .. الاوبتيمائيزيشن ليه سيتنتج معين .. ان الفانكشن اللي انت عاوز تعملها اوبتيمائيزيشن مش شرط تبقى هي اللي احنا مهتمين بيه ... زي مثلاً اللايكللي هوود فانكشن بتاعت اللوجستيك ريجريشن .. بس الي انت مهتم بيه اصلاً هو الايرور بتاع ال test data عشان الجينيرالايزيشن .. ثاني حاجه مهمه ف الستنج بتاع الاوبتيمائيزيشن هو ان الداتا اصلاً ممكن تبقى نويزي و في الحاله ديه انك تجيب optimal percentage مش هتبقا اهم حاجه يعني .. بس في شوية حاجات تانيه ممكن تخلي ال precision مهمه يعني .. ثالث حاجه هو ال early stopping وده مهم جداً الحقيقه عشان يبساعد ال generalization ..

## Optimization for ML

Not quite the same setting as other fields...

- Function we are optimizing might not be the true goal  
(e.g. likelihood vs generalization error)
- Precision might not matter  
(e.g. data is noisy, so optimal up to  $1e-16$  might not help)
- Stopping early can help generalization error  
(i.e. “early stopping” is a technique for regularization – discussed more next time)

الدكتور بيقول لما بيتكتب ال argmin .. هل انت بتبقا فاهم ايه المقصود منها .. طبعاً انا حمار فبفترض انه عاوز الحاجه اللي تخلي الحاجه اقل حاجه .. هه .. بطيخه .. فالدكتور هيقول ايه اللي بيبقا مقصود



1. Q: What is  $v^*$ ?

$v^* = 1$ , the minimum value of the function

2. Q: What is  $x^*$ ?

$x^* = 0$ , the argument that yields the minimum value

قال  $\text{agmin } x$  هي قيمة الإكس اللي إبت قيمة الفانكشن المنيم ... احنا هنستخدم الحاجات ديه عشان نقدر نفكر في الليبير ريجريشن كفافكشن أبروكسميشن .. هنا الداتا هتبقى شوية  $x$  and  $y$  pairs و  $x$  هتبقا فيكتور طولها  $M$  و الواي هتبقا رقم حقيقي .. احنا هنا هنفترض ان الداتا  $D$  كانت Generated باستخدام البروسيس ديه: البروسيس بتقول ان إكس كانت  $\text{sampled from some probability distribution } p$  star .. واحنا مش عارفين ايه هي ال  $P$  star اصلا هي unknown prob. Distribution .. وبعدين هنجيب ال  $y$  عن طريق اننا نباضي الإكسات لل unknown function  $h$  star .. فهنا الفرضيه علي ان ال  $p$  star و ال  $h$  star مش عارفينهم ... قدام شوية هنعدل علي الفرضيات ديه و نخطط كمان شوية نويز .. يعني المصدر بتاع ال  $y$  مش هيبقا deterministic ....

تعال نختار hypothesis space .. وهنا ال hypothesis space هيبقا هو ال all linear functions in  $M$  dimensional space .. فده هيبقا زي ال perceptron hypothesis space ماعدا اننا معندناش sign حوالين ال  $\theta^T x$  .. وبعدين هنختار an objective function هتبقا ال mean squared error وده اللي اتكلمنا عليه المره اللي فاتت اننا هنقل ال Sum of residuals .. احنا هناخد كل residual و نربعو وبعدين نجمعهم كلهم .. فين ال  $w$  و ال  $b$  .. في الثيتا .. فهنا انت تقدر تفترض ان أول فيتشر في الإكس هي دايماً بواحد .. و ده هيبقا ال bias بتاعك ..

## Linear Regression as Function Approximation

$\mathcal{D} = \{x^{(i)}, y^{(i)}\}_{i=1}^N$   
where  $x \in \mathbb{R}^M$  and  $y \in \mathbb{R}$

1. Assume  $\mathcal{D}$  generated as:

$$x^{(i)} \sim p^*(\cdot) \quad \text{unknown}$$

$$y^{(i)} = h^*(x^{(i)}) + \epsilon$$

2. Choose hypothesis space,  $\mathcal{H}$ :  
all linear functions in  $M$ -dimensional space

$$\mathcal{H} = \{h_\theta : h_\theta(x) = \theta^T x, \theta \in \mathbb{R}^M\}$$

3. Choose an objective function:  
mean squared error (MSE)

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N e_i^2$$

$$= \frac{1}{N} \sum_{i=1}^N (y^{(i)} - h_\theta(x^{(i)}))^2$$

$$= \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \theta^T x^{(i)})^2$$

4. Solve the unconstrained optimization problem via favorite method:

- gradient descent
- closed form
- stochastic gradient descent
- ...

look for value of vector  $\theta$  that minimizes that function  $J(\theta) \rightarrow \hat{\theta}$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta)$$

5. Test time: given a new  $x$ , make prediction  $\hat{y}$

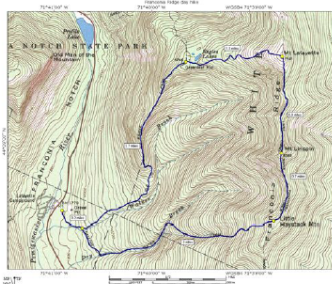
$$\hat{y} = h_{\hat{\theta}}(x) = \hat{\theta}^T x$$

10

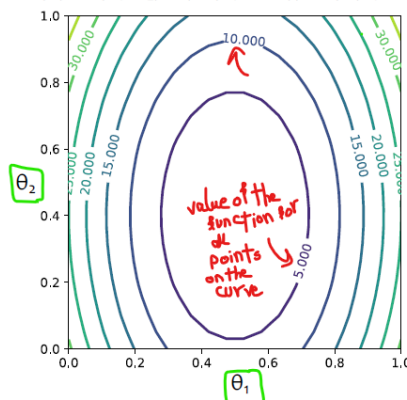
## Contour Plots

### Contour Plots

1. Each level curve labeled with value
2. Value label indicates the value of the function for all points lying on that level curve
3. Just like a topographical map, but for a function



$$J(\theta) = J(\theta_1, \theta_2) = (10(\theta_1 - 0.5))^2 + (6(\theta_2 - 0.4))^2$$



12

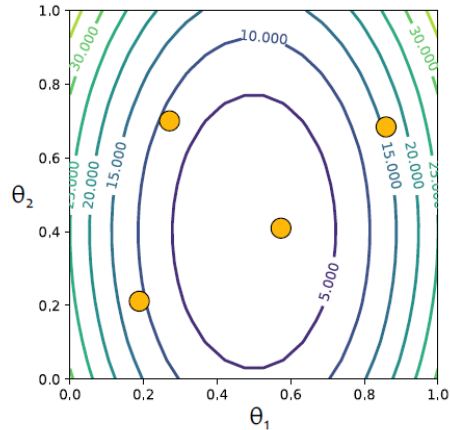
هنا بقي أول أوبيتيمائيزيشن الجورزم .. اسمو random guessing ... لو هنجي نختار البارمترز بتاعت الموديل بتاعنا بناءً علي الألجورزم ده .. تعال نشوف ايه اللي هيجصل ..

## Optimization by Random Guessing

### Optimization Method #0: Random Guessing

1. Pick a random  $\theta$
2. Evaluate  $J(\theta)$
3. Repeat steps 1 and 2 many times
4. Return  $\theta$  that gives smallest  $J(\theta)$

$$J(\theta) = J(\theta_1, \theta_2) = (10(\theta_1 - 0.5))^2 + (6(\theta_2 - 0.4))^2$$



t	$\theta_1$	$\theta_2$	$J(\theta_1, \theta_2)$
1	0.2	0.2	10.4
2	0.3	0.7	7.2
3	0.6	0.4	1.0
4	0.9	0.7	19.2

13

لو جيت تبص علي اللي بيحصل هتلاقي انك بتجيب أرقام راندم ما بين التيتا 1 و 2 ... وبعدين تقوم معوض وتجيب الناتج بتاع ال  $J(\theta_1, \theta_2)$  .... هتلاقي ان الأرقام بتاعت ال  $J$  كبرت فجأة .. فهتسأل نفسك سؤال .. هو ليه بيعمل كذا .. ما الكونتور عندو اهو ما يروح عليه و يجيب ال minimum .. الدكتور اعترض ع النقطة ديه وقال طب لو انت عندك 10000 تيتا.. هتبص ع الكونتور ازاي ... وهتعلم الأوبيتيمائيزيشن ازاي ..

تعال نبص علي الحوار ده من سكوب ال linear regression

## Optimization by Random Guessing

### Optimization Method #0: Random Guessing

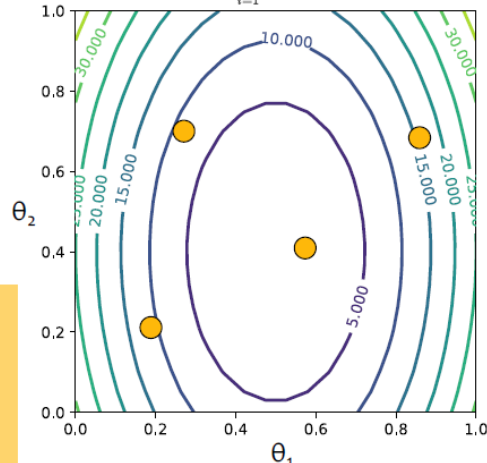
1. Pick a random  $\theta$
2. Evaluate  $J(\theta)$
3. Repeat steps 1 and 2 many times
4. Return  $\theta$  that gives smallest  $J(\theta)$

### For Linear Regression:

- **objective function** is Mean Squared Error (MSE)
- $MSE = J(w, b)$   

$$= J(\theta_1, \theta_2) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \theta^T x^{(i)})^2$$
- contour plot: each line labeled with MSE – **lower means a better fit**
- **minimum** corresponds to parameters  $(w, b) = (\theta_1, \theta_2)$  that **best fit** some training dataset

$$J(\theta) = J(\theta_1, \theta_2) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \theta^T x^{(i)})^2$$



t	$\theta_1$	$\theta_2$	$J(\theta_1, \theta_2)$
1	0.2	0.2	10.4
2	0.3	0.7	7.2
3	0.6	0.4	1.0
4	0.9	0.7	19.2

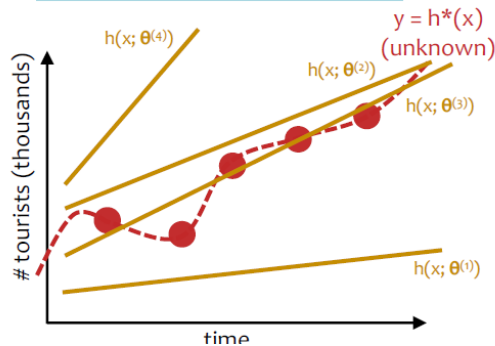
14

طب تعال نرن ال random guessing علي ال linear regression .. فهحتاج لكل hypothesis حاجتين slope and intersect ..

# Linear Regression by Rand. Guessing

## Optimization Method #0: Random Guessing

1. Pick a random  $\theta$
2. Evaluate  $J(\theta)$
3. Repeat steps 1 and 2 many times
4. Return  $\theta$  that gives smallest  $J(\theta)$



## For Linear Regression:

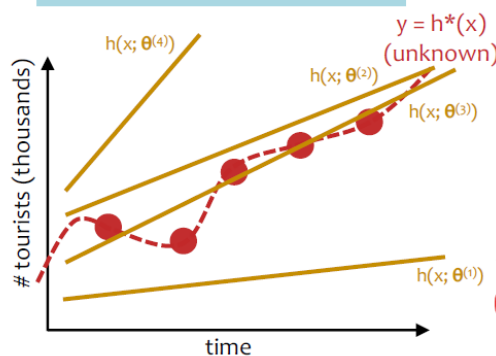
- target function  $h^*(x)$  is **unknown**
- only have access to  $h^*(x)$  through **training examples**  $(x^{(i)}, y^{(i)})$
- want  $h(x; \theta^{(t)})$  that **best approximates**  $h^*(x)$
- **enable generalization** w/inductive bias that restricts hypothesis class to **linear functions**

الخطوط فيه هي الراندم guessing اللي حصل .. تعال نخط الصورتين مع بعض .. اللي هو الجراف والكتنور مع بعض ..

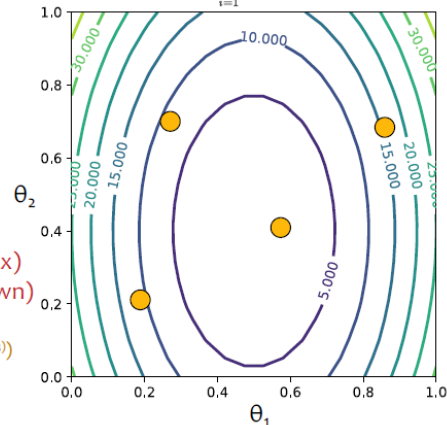
# Linear Regression by Rand. Guessing

## Optimization Method #0: Random Guessing

1. Pick a random  $\theta$
2. Evaluate  $J(\theta)$
3. Repeat steps 1 and 2 many times
4. Return  $\theta$  that gives smallest  $J(\theta)$



$$J(\theta) = J(\theta_1, \theta_2) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \theta^T x^{(i)})^2$$



t	$\theta_1$	$\theta_2$	$J(\theta_1, \theta_2)$
1	0.2	0.2	10.4
2	0.3	0.7	7.2
3	0.6	0.4	1.0
4	0.9	0.7	19.2

Return  $\theta_3$

16

طيب تعال نفكر في ال optimization method #1 .. عشان نعمل كدا هحتاج شوية ديفينيشنز .. أول واحد هو unconstrained optimization .. فلما بنتكلم علي unconstrained optimization .. انت بتقول ان عندك فانكشن اسمها  $J(\theta)$  و ال  $\theta$  فيه بت map من M dimensional space ل رقم واحد بس .. وهدفنا اننا نلاقي ال  $\hat{\theta}$  اللي هي ال argmin of all possible  $\theta$  of  $J(\theta)$  .. للماشين ليرننج .. ال  $\theta$  هنا هي أوبجكتيف فانكشن ... والثيتا هتبقا البارامترز بتاعتك .. لو هنجي نشغل ال gradient descent .. لازم نفهم ايه هو ال gradient ... فدلوقت هنعوز نفهم يعني ايه ال derivatives .. هنفكر كدا بس ايه هو ال interpretations بتاعت ال derivatives .. بص ع الصورة ..

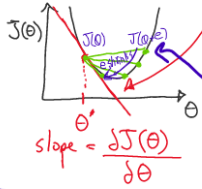
## Unconstrained Optimization

Given fn.  $J(\theta)$   $J: \mathbb{R}^M \rightarrow \mathbb{R}$

Goal is  $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta)$   
 ← parameters  
 ← For ML  $J$ 's an obj. fn.

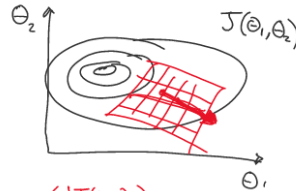
## Derivatives

① Deriv. as Slope of Tangent



$$\text{slope} = \frac{dJ(\theta)}{d\theta}$$

③ Deriv as Tangent Plane



② Deriv as a Limit

$$\frac{dJ(\theta)}{d\theta} = \lim_{e \rightarrow 0} \frac{J(\theta+e) - J(\theta)}{e}$$

limit of secants is tangent

$$\left\langle \frac{dJ(\theta_1, \theta_2)}{d\theta_1}, \frac{dJ(\theta_1, \theta_2)}{d\theta_2} \right\rangle$$

← e shrinks → secant line approaches the tangent line

## Gradient

The gradient of  $J$

$$\nabla J(\theta) = \begin{bmatrix} \frac{dJ(\theta)}{d\theta_1} \\ \frac{dJ(\theta)}{d\theta_2} \\ \vdots \\ \frac{dJ(\theta)}{d\theta_M} \end{bmatrix}$$

← first order partial derivatives

## Gradient Descent

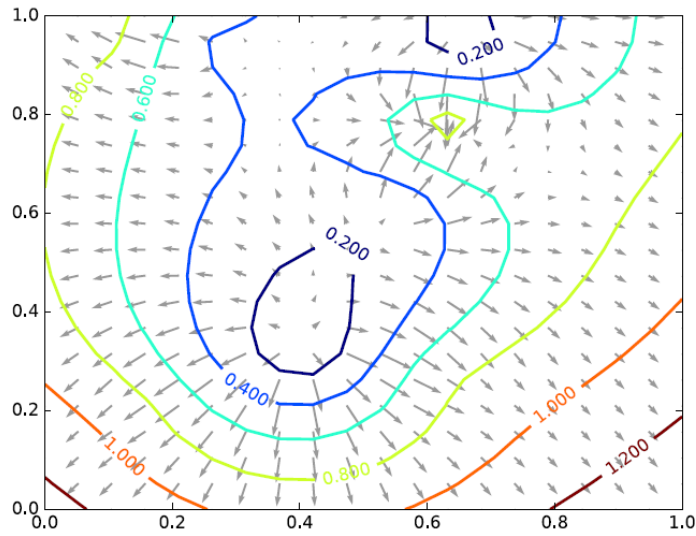
Algorithm:

- ① Choose an initial point  $\hat{\theta}$
- ② Repeat For  $t=1, 2, 3, \dots$ 
  - a) Compute gradient  $\vec{g} = \nabla J(\theta)$
  - b) Choose a step size  $\gamma_t > 0$
  - c) Update  $\hat{\theta} \leftarrow \hat{\theta} - \gamma_t \vec{g}$
- ③ Return  $\hat{\theta}$  when stopping criterion is met

الدكتور شرح مثال عليه هو و مراته ... الحقيقه المثال حلو فاسمعو من المحاضره مش هكتبو . المهم انو في الاخر بيتارجت حته انك عارف طول منت بتنزل من ع الجبل انت عارف انك مسيرك تلاقى العربيه بتاعتك .. الالجورزم ده هو ال gradient descent .. تعال نشوف الكونتور بلوت بتاع الفانكشن

انت هتحتاج شوية معلومات من الجريدينيت .. عشان انت لما هتيجي ت evaluate ال gradient علي نقطه معينه في الفانكشن ... هيديك فيكتور دايماً بيبشاور ناحية ال steepest ascend اللي هو الاتجاه اللي لو انت مشيتو .. فتهزود قيمة الفانكشن .. فبالتالي الرسمه اللي تحت هي grid of gradients ولكل gradient عندك فيكتور بتاع التانجنت بليين ..

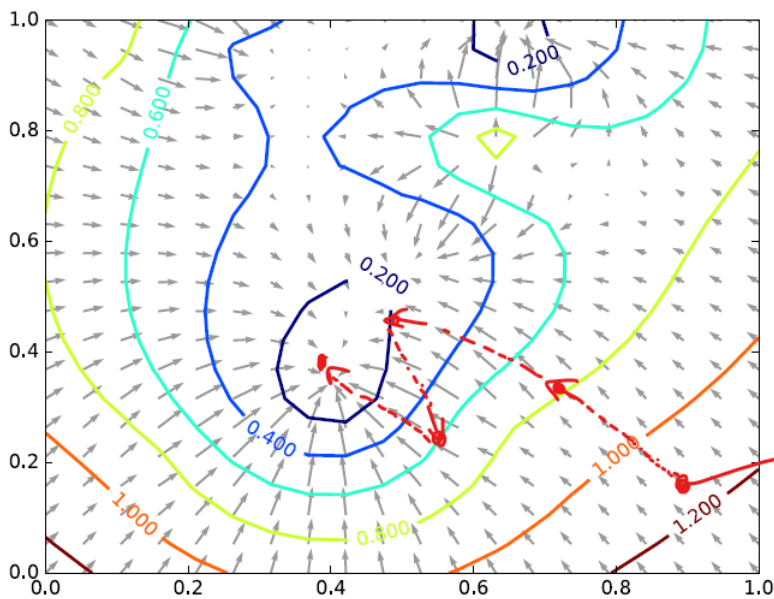
# Gradients



These are the **gradients** that Gradient **Ascent** would follow.

لو جينا نبص علي النيجاتيف جريدينتس .. هيشاورو ناحية ال steepest descent

## (Negative) Gradients

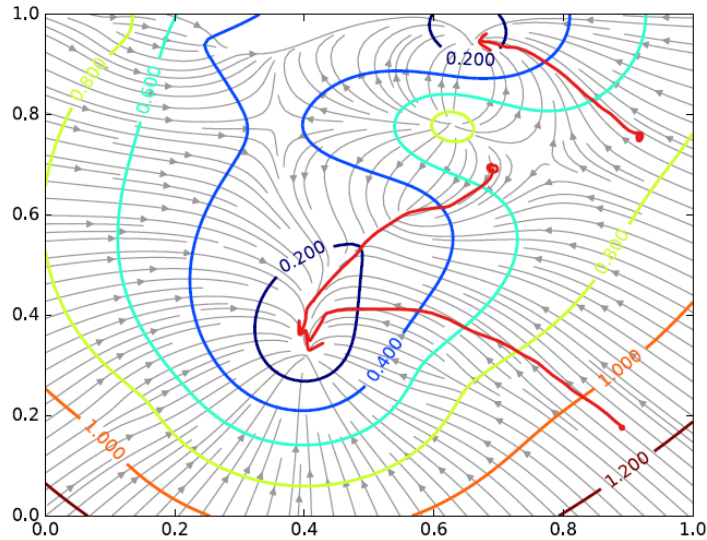


These are the **negative** gradients that Gradient **Descent** would follow.

لو جينا نبص علي الطريق اللي هنمشي عليه من أي starting point ..



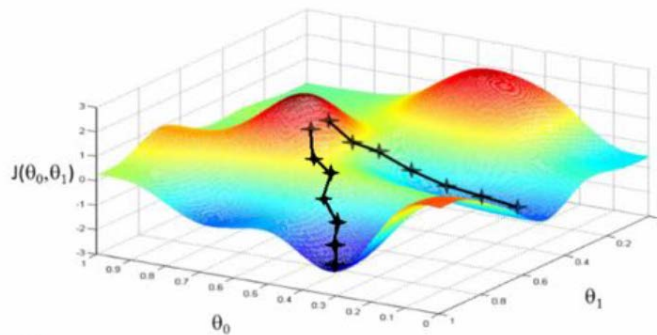
# (Negative) Gradient Paths



Shown are the **paths** that Gradient Descent would follow if it were making **infinitesimally small steps**.

## Pros and cons of gradient descent

- Simple and often quite effective on ML tasks
- Often very scalable
- Only applies to smooth functions (differentiable)
- Might find a local minimum, rather than a global one



e courtesy of William Cohen

25

تعال بقي نشوف الـ gradient descent كـ الجورنم . هيمشي كالاتي :

1. اختيار نقطة البداية بتاعتك .. اللي هي الـ Theta vector

2. Repeat the following steps:

a. Compute gradient  $g = \text{gradient of } J(\theta)$

b. Choose step size, which is some value  $\text{Gamma} > 0$  and it's a real number. It says how much we extend or

shrink that gradient vector to decide how far to step

c. Update Our paramtere  $\text{Theta} \leftarrow \text{Theta} - \text{Gamma} * g$

3. Return theta when Some stopping criterion is met

وده يشاباب هو الجريدنيت ديسينت .. في شوية أسئلة .. ازاي نختار الـ starting point .. أول طريقه هي اللي عملناها للـ perceptron .. اللي هو حط كلو اصفار .. ثاني طريقه اختارهم random ..

طب بالنسبه لل stopping criterion .. انت محتاج ت stop taking steps .. أياً كان ال steps اللي هتاخدها هي هتقل .. ليه هتقل .. عشان لو جيت تبص ع الصوره هتلاقي ان الجريدينتس اللي عند الحتت اللي فعلاً steep .. بتلاقي الماجنتيود بتاعها كبير .. ماجنتيود جامد .. بس عند المنيم الفيكتور بيبقا رفيع أوي وصغير كذا ... ليه .. لأن الفيكتور بيبرزنت ان السلوب بتاع البلين هنا هيبقا بصفر .. وفعلاً عند الإكزاكت منيم .. هو فعلاً هيبقا سلوب صفر .. فالستوبنج كرايترين هتبقا ان

Whenever the gradient of  $J(\theta)$  and taken its L2 norm is  $<$  some epsilon (small number i.e.  $10^{-3}$ ), then in this case the length of that gradient vector must be really short and so it's almost 0 and we are probably close to the minimum

آخر حاجه هو ال step size ... اول اختيار هو الفيكسيد فاليو .. جاما ..  $\Gamma = 0.1$  .. وفي أویشن ثاني هو ال exact line search .. وفي أویشن ثالث اللي هو ال backtracking line search .. فكرة اللالين سيرش هو .. لو افترضنا ان عندنا فانكشن و المنيم بتاعها هو في نص الأوضه اللي احنا فيها .. والجرادينت ببشاور بشكل مباشر ناحية المنيم ده .. هيبقا حاجه جامده جداً لو عرفت انت محتاج كام ستيب .. بس انت مش هتعرف كام ستيب .. فهتأخذ ستيب واحد ... وبعدين تشوف هل انا احسن .. وبعدين تقول انا اخذ 2 ستيبس .. هل انا احسن .. هتقول اه فعلاً انا كذا احسن .. وبعدين تأخذ 4 ستيبس وتشوف هل كذا احسن .. تقول انا ايه اللي جابني ف غريبه احنا ف عجيبه .. ارجع خطوتين لورا "Backtrack" .. هل انت احسن .. اه .. خذ خطوه قدام .. هل انت احسن .. اه وتقريباً كذا انا وصلت عجيبه .. هه ..

في تكنيك رابع وهو انك تستخدم  $\text{schedule } \Gamma_t = \Gamma_0 / ((t-1) * \Gamma_0 + 1)$

### Gradient Descent

#### Algorithm:

- ① Choose an initial point  $\tilde{\theta}$
- ② Repeat For  $t=1,2,3,\dots$ 
  - a) Compute gradient  $\tilde{g} = \nabla J(\tilde{\theta})$
  - b) Choose a step size  $\gamma_t > 0$
  - c) Update  $\tilde{\theta} \leftarrow \tilde{\theta} - \gamma_t \tilde{g}$
- ③ Return  $\tilde{\theta}$  when stopping criterion is met

#### Remarks

##### Starting Point:

- ①  $\tilde{\theta} = 0$
- ②  $\tilde{\theta}$  randomly

##### Stopping Criterion:

$$\|\nabla J(\tilde{\theta})\|_2 < \epsilon$$

$\epsilon = 10^{-8}$

##### Step Sizes:

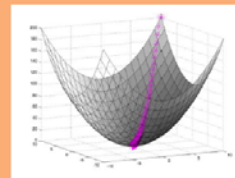
- Fixed value  $\gamma = 0.1$
- exact line search
- backtracking line search
- Schedule  $\gamma_t = \frac{\gamma_0}{(t-1)\gamma_0 + 1}$

### Gradient for Linear Regression

## Gradient Descent

### Algorithm 1 Gradient Descent

- 1: procedure  $\text{GD}(\mathcal{D}, \theta^{(0)})$
- 2:  $\theta \leftarrow \theta^{(0)}$
- 3: while not converged do
- 4:  $\theta \leftarrow \theta - \gamma \nabla_{\theta} J(\theta)$
- 5: return  $\theta$



In order to apply GD to Linear Regression all we need is the **gradient** of the objective function (i.e. vector of partial derivatives).

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{d}{d\theta_1} J(\theta) \\ \frac{d}{d\theta_2} J(\theta) \\ \vdots \\ \frac{d}{d\theta_M} J(\theta) \end{bmatrix}$$

27

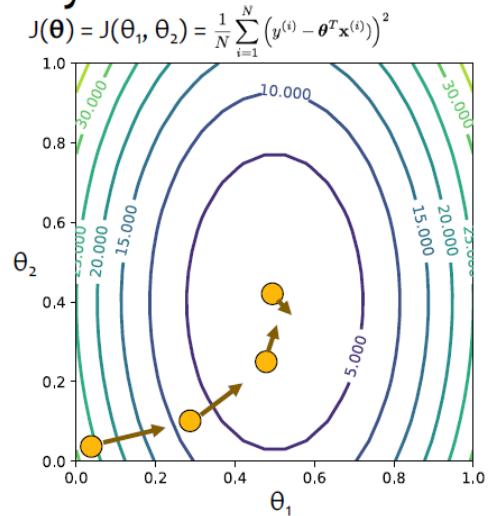
تعال نشوف بقا ايه هو ال gradient في ال linear regression ... لما كنا بنبص من فوق كذا . احنا كنا علوزين ناخذ ال mean squared error function و نعملها minimization فاحنا هنعوز نحسبها ال gradient بتاع ال mean squared error .. وده هيديلنا طريقه اننا نحسب ال .. gradient descent



# Linear Regression by Gradient Desc.

## Optimization Method #1: Gradient Descent

1. Pick a random  $\theta$
2. Repeat:
  - a. Evaluate gradient  $\nabla J(\theta)$
  - b. Step opposite gradient
3. Return  $\theta$  that gives smallest  $J(\theta)$



t	$\theta_1$	$\theta_2$	$J(\theta_1, \theta_2)$
1	0.01	0.02	25.2
2	0.30	0.12	8.7
3	0.51	0.30	1.5
4	0.59	0.43	0.2

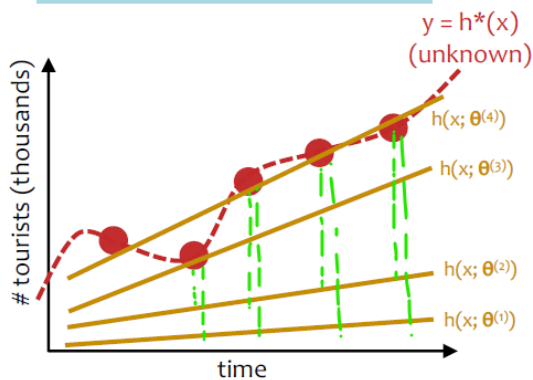
31

المهم تعال نرجع للمثال بتاعنا اللي هو ال linear function اللي بنحاول ن fit on the data .. الإكس أكسيس هو التايم و عدد السياح علي الواي أكسيس  
و عندك التروو فانكشن اللي انت مش عارفها .. و هنعوز نبص علي السلوب والإنترسييت ... هنفضل ناخذ خطوات .. هتلاحظ ان ال Mean squared  
error بقا احسن .. اللي هو قيمته قلت .. ليه .. عشان المسافات ما بين التروو و البريديكشن قلت

# Linear Regression by Gradient Desc.

## Optimization Method #1: Gradient Descent

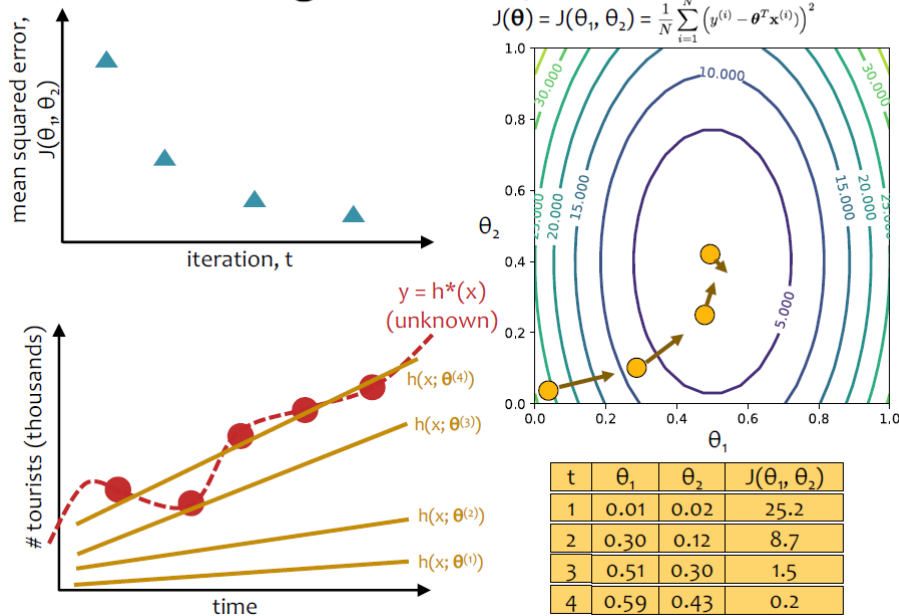
1. Pick a random  $\theta$
2. Repeat:
  - a. Evaluate gradient  $\nabla J(\theta)$
  - b. Step opposite gradient
3. Return  $\theta$  that gives smallest  $J(\theta)$



t	$\theta_1$	$\theta_2$	$J(\theta_1, \theta_2)$
1	0.01	0.02	25.2
2	0.30	0.12	8.7
3	0.51	0.30	1.5
4	0.59	0.43	0.2

32

# Linear Regression by Gradient Desc.



35

احنا اللي احنا محتاجينو دلوقت هو ال actual gradient for linear regression... هنشوف دلوقت عظمه في الدرفيجن لل gradient for linear regression

هنبدأ بان عندنا معادلة ال mean squared errors ..

Gradient for Linear Regression

MSE:  $J(\theta) = \frac{1}{N} \sum_{i=1}^N J^{(i)}(\theta)$  where  $J^{(i)}(\theta) = \frac{1}{2} (y^{(i)} - \theta^T x^{(i)})^2$

$\frac{\partial J^{(i)}(\theta)}{\partial \theta_j} = \frac{1}{2} (y^{(i)} - \theta^T x^{(i)}) \frac{\partial}{\partial \theta_j} (y^{(i)} - \theta^T x^{(i)})$

$= (y^{(i)} - \theta^T x^{(i)}) \frac{\partial}{\partial \theta_j} (y^{(i)} - \sum_{m=1}^M \theta_m x_m^{(i)})$

$= -(y^{(i)} - \theta^T x^{(i)}) x_j^{(i)}$

$\nabla J^{(i)}(\theta) = \begin{bmatrix} \frac{\partial J^{(i)}(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J^{(i)}(\theta)}{\partial \theta_M} \end{bmatrix}$

$= -(y^{(i)} - \theta^T x^{(i)}) \vec{x}^{(i)}$

$\nabla J(\theta) = \nabla \left( \frac{1}{N} \sum_{i=1}^N J^{(i)}(\theta) \right) = \frac{1}{N} \sum_{i=1}^N \nabla J^{(i)}(\theta)$

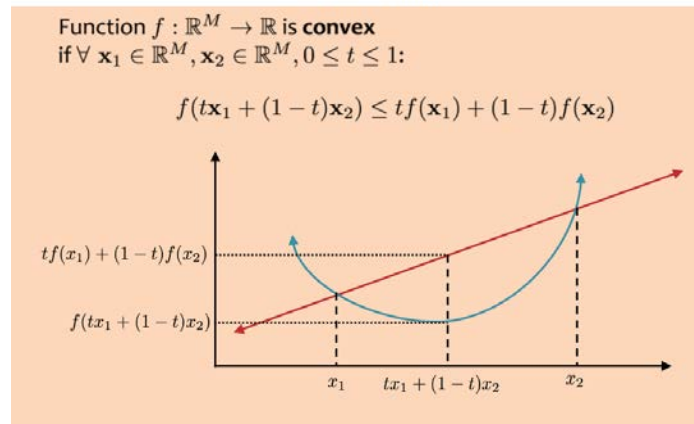
$= \frac{1}{N} \sum_{i=1}^N -(y^{(i)} - \theta^T x^{(i)}) \vec{x}^{(i)}$

Handwritten notes: "عشانها تهتق 2 ال", "doesn't affect argmin", "Take weighted avg", "Scalar", "vector".

طيب دلوقت هنتكلم علي ال convexity ... احنا نقول فانكشن تبعنا كونفيكس لو لكل ال pairs of points  $x_1$  and  $x_2$  و لكل قيم ال  $t$  ما بين الصفر و الواحد .. انت تقدر تاخذ function of the average of the 2 points و ده بيقا اقل من ال average function evaluated at those 2 points .. افترض ان عندك نقطة  $x_1$  و نقطة  $x_2$  .. وانت عاوز تتأكد ان لأي قيمة لل  $t$  ما بين الصفر و الواحد .. هيديك نقطة علي ال red sequin اللي بيوصل ما بين الإكس 1 و 2 ... اللي هو الخط الأحمر ... و convexity بتقول ان أيأ كان القيمة اللي علي الخط الأحمر ده .. هتبقا أكبر من القيمة الحقيقية اللي هتجيبها اللي هي  $tx_1 + (1-t)x_2$  .. وده لكل ال possible pairs  $x_1$  and  $x_2$  .. فلو انت عندك convex function ديه ليها شوية خصائص .. بس

تعال نحت شوية تعريفات .. لو عندك جينيريك فانكشن اسمها  $f(x)$  القيمة  $x^*$  هي جلوبال منيمم .. if and only if ال  $f(x^*)$  اقل من ال  $f(x)$  لكل القيم المحتمله لل  $x$  .. هنقول انها لو كالم منيمم ده لو موجود some epsilon some window اللي هو .. هتلاقى ان عندك local minimum .. فلو عندك convex function .. single local minimum هي global minimum .. لو لاقيت أي منيمم .. أمسك فيه ... " أخيراً لاقيتك .. رايح فين" ... انما لو NonConvex ... مش شرط

## Convexity



## Convexity

Suppose we have a function  $f(x): \mathcal{X} \rightarrow \mathcal{Y}$ .

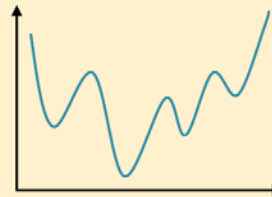
- The value  $x^*$  is a **global minimum** of  $f$  iff  $f(x^*) \leq f(x), \forall x \in \mathcal{X}$ .
- The value  $x^*$  is a **local minimum** of  $f$  iff  $\exists \epsilon$  s.t.  $f(x^*) \leq f(x), \forall x \in [x^* - \epsilon, x^* + \epsilon]$ .

### Convex Function



- Each **local minimum** is a **global minimum**

### Nonconvex Function

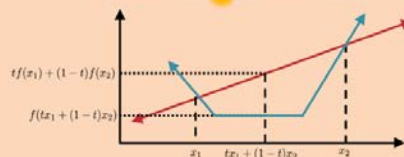


- A nonconvex function is **not convex**
- Each **local minimum** is **not necessarily a global minimum**

## Convexity

Function  $f: \mathbb{R}^M \rightarrow \mathbb{R}$  is **convex**  
if  $\forall \mathbf{x}_1 \in \mathbb{R}^M, \mathbf{x}_2 \in \mathbb{R}^M, 0 \leq t \leq 1$ :

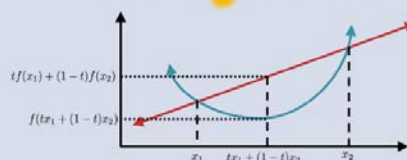
$$f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \leq tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$



Each **local minimum** of a **convex** function is also a **global minimum**.

Function  $f: \mathbb{R}^M \rightarrow \mathbb{R}$  is **strictly convex**  
if  $\forall \mathbf{x}_1 \in \mathbb{R}^M, \mathbf{x}_2 \in \mathbb{R}^M, 0 \leq t \leq 1$ :

$$f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) < tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$



A **strictly convex** function has a **unique global minimum**.

# Calculus and Optimization

## In-Class Exercise

Plot three functions:

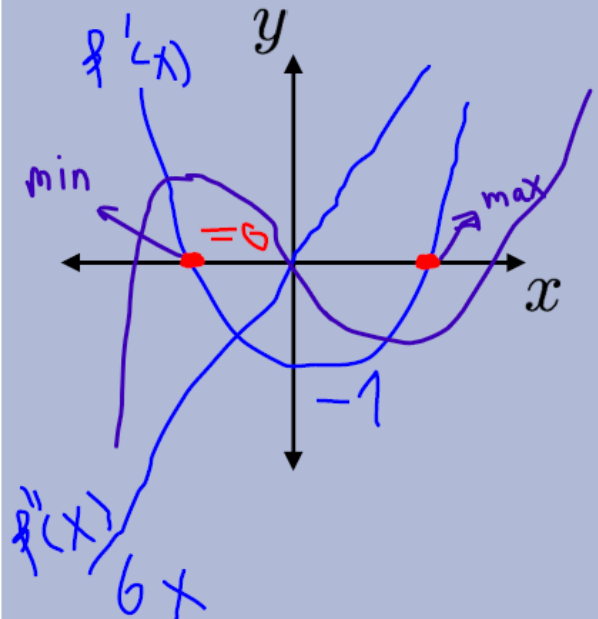
$$1. f(x) = x^3 - x$$

$$2. f'(x) = \frac{\partial y}{\partial x}$$

$$3. f''(x) = \frac{\partial^2 y}{\partial x^2}$$

$\nabla = 0$   
 or  
 min or max  
 $-ve$   
 min  
 $+ve$   
 max

## Answer Here:



50

فبالتالي لما تيجي تحل في الكلوزد فورم ... تقدر ت jump علي طول لل local minimum .. لو انت بس عرفت حاجه عن ال gradient بتاع الفانكشن .. فالحل في الكلوزد فورم .. انك مثلاً عندك M dimensional function  $J(\theta)$  and  $\theta$  is M dim. Function

1. هنحل المعادله بتاعت ال  $\text{Gradient}(J(\theta)) = 0$  for  $\theta$

2. Test for min, max or saddle point using second dervative

فالدكتور بيقول .. هو مش هيبقا حاجه جامده جداً لو بدل ما كنا بنأخذ شوية STEPS صغيره بجريدينت ديسسنت .. احنا نحل علي طول المين سكويرد إيرور ونجيب قيمة للثيتا ليها الجريدينت بصفر .. لو لاقينا القيمه ديه .. والفانكشن اصلاً كونفيكس .. بيقا خلاص يا ماما ان بيقا انت وصلت للجلوبال منيم ...

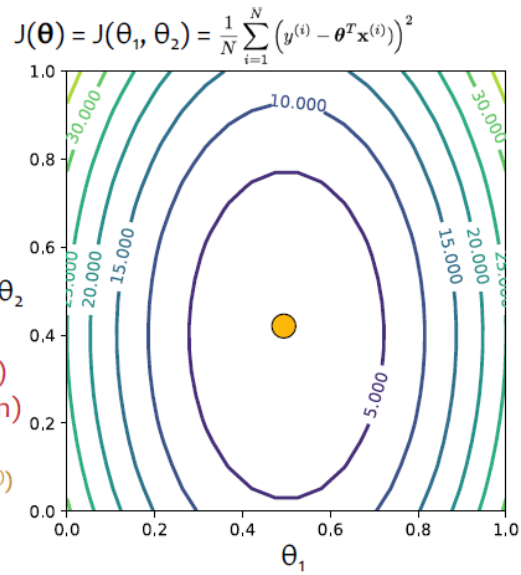
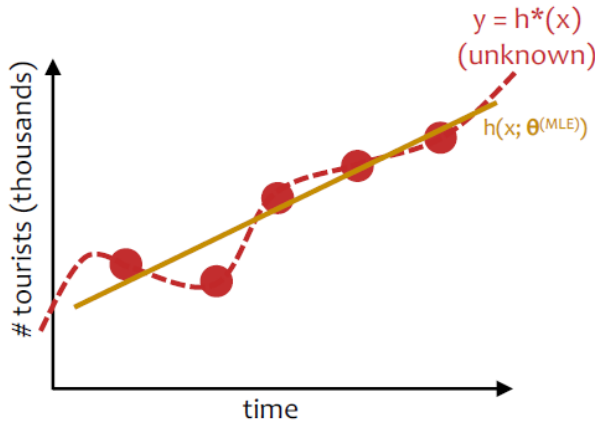
# Linear Regression: Closed Form

## Optimization Method #2: Closed Form

1. Evaluate

$$\theta^{\text{MLE}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

2. Return  $\theta^{\text{MLE}}$



t	$\theta_1$	$\theta_2$	$J(\theta_1, \theta_2)$
MLE	0.59	0.43	0.2

55

هي خطوه واحده وخلصنا فعلاً ... ده الكلود فورم سلوشن .. فالكلوزد فورم سلوشن قول مثلاً عندك فيكتور اسمو  $\mathbf{y}$  جاواه كل الأوتبوت فالبيوز .. وبعدين هيبقا عندك ماتركس اسمها  $\mathbf{X}$  وديه بنتكون من كل الترئينج إكزامبلز بتاعتنا .. من أول فيتشر لحد آخر فيتشر .. الماتركس ديه الدايمنشز بتاعتها ... Design Matrix اسمها  $\mathbf{X}$  وبعدين نكتب ال  $J(\theta)$

## Closed Form Solution

Ex:  $M$ -dim fn.  $J(\theta)$ ,  $\theta \in \mathbb{R}^M$

① Solve  $\nabla J(\theta) = 0$  for  $\theta$

② test for min/max, or saddle point  
using second derivatives

## Closed Form Solution for Lin Reg.

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^{(1)} & \dots & x_M^{(1)} \\ \vdots & & \vdots \\ x_1^{(N)} & \dots & x_M^{(N)} \end{bmatrix}$$

$\vec{x}^{(1)}$   
 $\vec{x}^{(N)}$   
 Design Matrix

① Write  $J(\theta)$

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (y^{(i)} - \theta^T \vec{x}^{(i)})^2$$

Ordinary Least Squares (OLS)

$$= \frac{1}{N} \cdot \frac{1}{2} (\mathbf{X} \theta - \vec{y})^T (\mathbf{X} \theta - \vec{y})$$

② "Normal Equations" + gradient set to zero

$$\nabla J(\theta) = \mathbf{X}^T \mathbf{X} \theta - \mathbf{X}^T \vec{y} = 0$$

① Write  $J(\theta)$

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (y^{(i)} - \theta^T x^{(i)})^2 \quad \text{Ordinary Least Squares (OLS)}$$

$$= \frac{1}{N} \cdot \frac{1}{2} (\mathbf{X} \vec{\theta} - \vec{y})^T (\mathbf{X} \vec{\theta} - \vec{y})$$

② "Normal Equations" ← gradient set to zero

$$\nabla J(\theta) = \mathbf{X}^T \mathbf{X} \vec{\theta} - \mathbf{X}^T \vec{y} = 0$$

③ Solve for  $\vec{\theta}$

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

$$= \underset{\theta}{\operatorname{argmin}} J(\vec{\theta})$$

وده ال closed form solution

## Computational Complexity of OLS

To solve the Ordinary Least Squares problem we compute:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (y^{(i)} - (\theta^T x^{(i)}))^2$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{Y})$$

The resulting shape of the matrices:

$$\underbrace{\begin{pmatrix} \mathbf{X}^T & \mathbf{X} \end{pmatrix}}_{M \times M}^{-1} \underbrace{\begin{pmatrix} \mathbf{X}^T & \mathbf{Y} \end{pmatrix}}_{M \times 1}$$

**Background: Matrix Multiplication** Given matrices **A** and **B**

- If **A** is  $q \times r$  and **B** is  $r \times s$ , computing **AB** takes  $O(qrs)$
- If **A** and **B** are  $q \times q$ , computing **AB** takes  $O(q^{2.373})$
- If **A** is  $q \times q$ , computing  $A^{-1}$  takes  $O(q^{2.373})$ .

**Computational Complexity of OLS:**

$\mathbf{X}^T \mathbf{X}$	$O(M^2 N)$
$(\quad)^{-1}$	$O(M^{2.373})$
$\mathbf{X}^T \mathbf{Y}$	$O(MN)$
$(\quad)^{-1}(\quad)$	$O(M^2)$
total	$O(M^2 N + M^{2.373})$

Linear in # of examples, N  
Polynomial in # of features, M

57

طبعاً دلوقت فرحت .. هيبه هوب دابل كيك نعمل الطريقه ديه بقا لكل الماشين ليرننج ولأبي حاجه فيها لينير ريجريشن .. طبعاً انت دافع تمن الحل ده ك computational cost و التمن ده مش سهل .. كل ما عدد الفيتشرز زاد و عدد ال examples زادو .. كل ما هتدفع أكثر .. فده مش حل كويس الحقيقه غير لما يكون ال N and M small