

أول حاجة هنرسم صوره كدا بتوضح احنا بنعمل ايه في ال linear regression ... عندك داتا سبت .. شوية نقط .. و بتحاول تشوف أحسن خط يحاول يجيب أحسن فيت للنقط ديه .. احنا هنا هندور علي خط optimal .. عاوزين ن formalize اللي بيحصل الاول .. الداتا نفسها بتتكون من شوية tuples .. إنبوت فيتشرز و تارجت أوبتوت .. ال x لما بتيقا bold ده عبارته عن فيكتور أوف فيتشرز .. زي مثلاً عدد أوض النوم عدد المطايخ عدد الحمامات ... التارجت هبيقا numerical value .. يعني مثلاً زي ال house prices .. أي حاجة bold هي فيكتور ... أي حاجة مش bold هي سكيلر عادي ... التارجت هنا اسمو تارجت عشان هو ال ground truth

احنا هدفنا اننا نلاقي hypothesis h اللي تقدر ت map x ل t ..

Problem

- Data: $\{(x_1, t_1), (x_2, t_2), \dots, (x_N, t_N)\}$
 - $x = \langle x_1, x_2, \dots, x_M \rangle$: input vector
 - t : target (continuous value)
- Problem: find hypothesis h that maps x to t
 - Assume that h is linear:

$$y(x, w) = w_0 + w_1 x_1 + \dots + w_M x_M = w^T \begin{pmatrix} 1 \\ x \end{pmatrix}$$

column vector

- Objective: minimize some loss function

– Euclidean loss: $L_2(w) = \frac{1}{2} \sum_{n=1}^N (y(x_n, w) - t_n)^2$

prediction

دلوقت احنا عاوزين نحل المشكله عن طريق optimization .. هو بيتكلم علي convex optimization .. ايه هي ال convex .. هي اللي بتحدد هل المشكله سهله ولا لا .. non convex ديه بتبقا صعبه ...

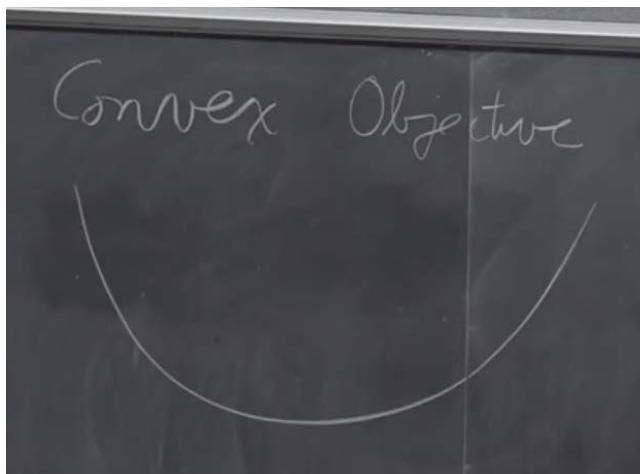
Optimization

- Find best w that minimizes Euclidean loss

$$w^* = \operatorname{argmin}_w \frac{1}{2} \sum_{n=1}^N \left(t_n - w^T \begin{pmatrix} 1 \\ x_n \end{pmatrix} \right)^2$$

- Convex optimization problem
 - \Rightarrow unique optimum (global)

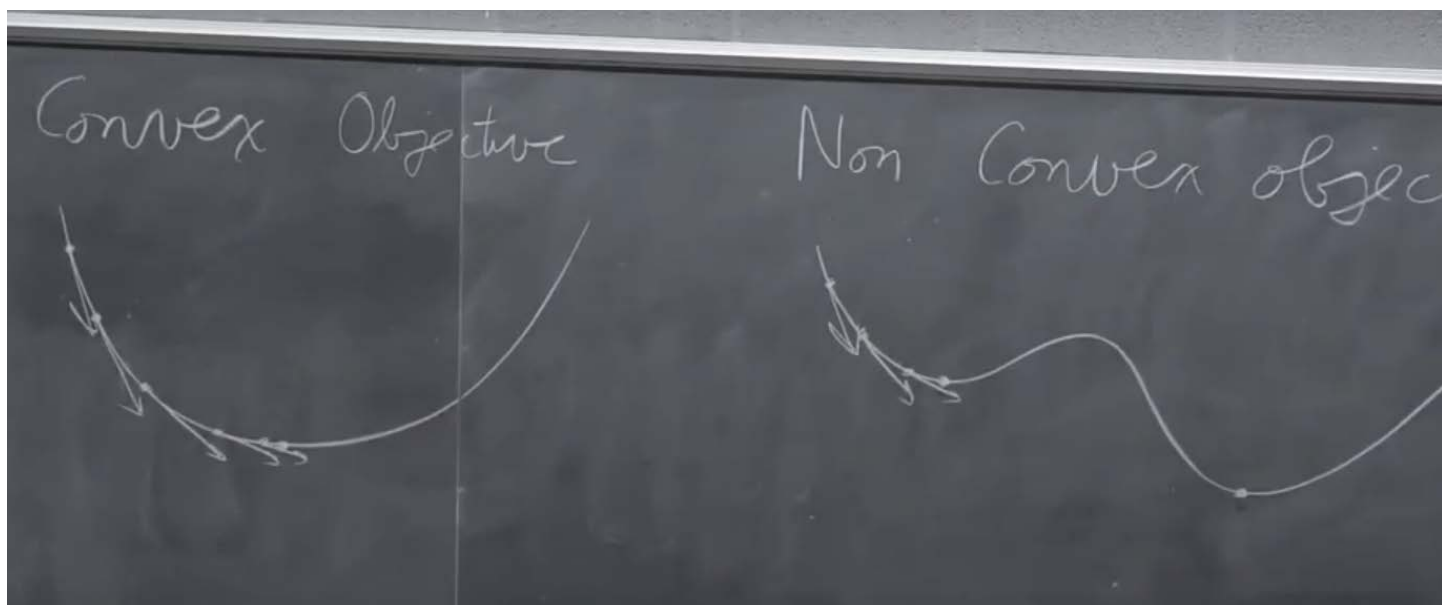
الدكتور رسم علي الصبوره بالطباشير .. هه . رسم convex objective .. وهو عاوز يجيب ال minimum of this function .. وبيقولك الموضوع سهل عشان انت عندك نقطه واحده منيم .. و منين ما تبدأ في الكيرف بتاعك هتعرف تنزل تحت عكس ال gradient .. و توصل للمنيم .. جامد ..



حاجه مثلاً gradient descent .. بتنزل علي الكيرف لحد ما توصل ..

هنيجي نبص علي Non convex objective

أدي كيرف عندو 2 منيم ... واحده لوكال منيمام و الثانيه جلوبال منيمام .. المشكله انك هتوصل للوكال منيمام تمام .. بس ال $\text{gradient} = 0$ وانا مش عارف اتحرك خلاص ... فبالتالي gradient descent مش هتعرف ت converge علي الجلوبال منيمام .. الا يعني لو بدأت من نقطه معينه



لل non convex optimization problems معندناش reliable techniques عشان نقدر ن converge بسرعه يعني ..

المشكله اللي احنا بنعتبرها في المحاضره هي convex .. و الهدف بتاعنا هو quadratic objective .. زي الصوره اللي علي اليمين فوق ... نرجع تاني للسلايدز ...

الدكتور قال اننا ممكن نلاقي ال gradient و ننزل خطوات وهكذا .. في حل احسن .. ان تقدر تحسب ال gradient عادي .. بس في شرط ان ال $\text{gradient} = 0$ انت وقفت خلاص غير كذا افضل اتحرك .. وهي ديه النقطه الوحيده اللي هي جلوبال منيمام ... احنا عاوزين نلاقي ال W اللي بت satisfy the equation .. اللي فيها ال Summation .. خد بالك ان ال W هي اللي مش عارفنها .. فهي ال variable .. احنا هنعوز نفجيب بقا ال W .. فقالك هتجيبها ازاي ..

$$Aw = b$$

Solution

- Let $\bar{\mathbf{x}} = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$ then $\min_{\mathbf{w}} \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \bar{\mathbf{x}}_n)^2$
- Find \mathbf{w}^* by setting the derivative to 0
$$\frac{\partial L_2}{\partial w_j} = \sum_{n=1}^N (t_n - \mathbf{w}^T \bar{\mathbf{x}}_n) \bar{x}_{nj} = 0 \quad \forall j$$
$$\Rightarrow \sum_{n=1}^N (t_n - \mathbf{w}^T \bar{\mathbf{x}}_n) \bar{\mathbf{x}}_n = 0$$
- This is a linear system in \mathbf{w} , therefore we rewrite it as $\mathbf{A}\mathbf{w} = \mathbf{b}$
where $\mathbf{A} = \sum_{n=1}^N \bar{\mathbf{x}}_n \bar{\mathbf{x}}_n^T$ and $\mathbf{b} = \sum_{n=1}^N t_n \bar{\mathbf{x}}_n$

Solution

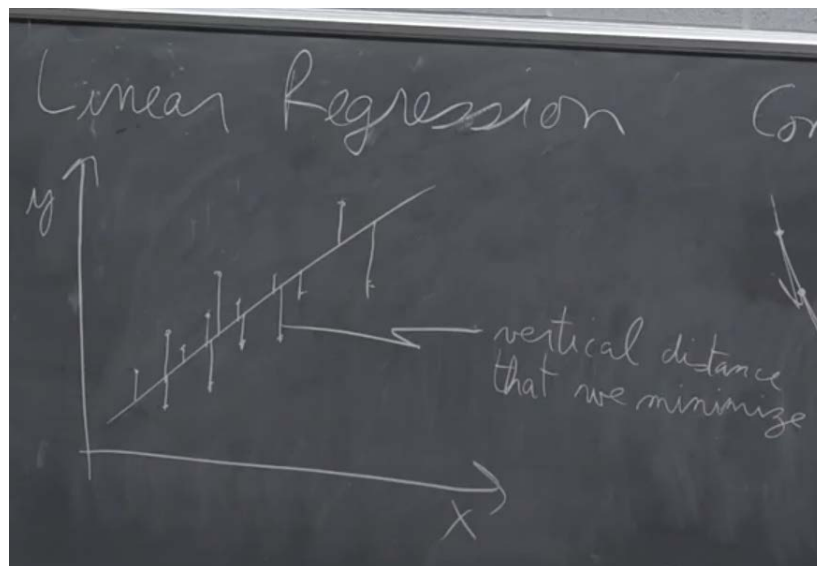
- If training instances span \mathbb{R}^{M+1} then \mathbf{A} is invertible:

$$\mathbf{w} = \mathbf{A}^{-1} \mathbf{b}$$

- In practice it is faster to solve the linear system $\mathbf{A}\mathbf{w} = \mathbf{b}$ directly instead of inverting \mathbf{A}
 - Gaussian elimination
 - Conjugate gradient
 - Iterative methods

الدكتور هيرسم ع الصبوره :

احنا مثلاً لاقينا ال line بتاع ال linear regression اللي محتاجينو .. لما تلاقي الخط اللي انت محتاجو .. احنا هدفنا نقلل ال Euclidean distance with respect to y .. فاحنا هدفنا نقلل ال vertical distance ما بين ال lines .. المهم خد بالك ان كل نقطه هتعوّز تشد الخط في ناحيه .. ففي الآخر كلو هي stabilize عند ال minimum شد في الخط



لما بنلاقي الحل بناءً علي ال least square .. الحل مش بيبقي robust .. عشان ممكن يحصل perturbation للانبوت .. فده هيغير ال line هيبقا فين .. وانت عاوز تشوف لو حصل تغيير بسيط .. الخط هيتغير تغيير بسيط .. بس انت ممكن تلاقي ان لو حصل تغيير في الانبوت .. هيجصل تغيير كبير جداً في الخط .. فعشان كذا هو unstable .. فده كذا overfitting .. ومشكلة ال least square انو too sensitive to noise ... هنعمل مثال بسيط:

Ex:1 هستخدم ال solution سلايد اللي في السلايدز بناعت الدكتور

$$\begin{aligned} \text{Ex:1} \\ A &= \begin{pmatrix} 1 & 1 \\ 0 & \epsilon \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 0 & \epsilon \end{pmatrix}^T = \begin{pmatrix} 2 & \epsilon \\ \epsilon & \epsilon^2 \end{pmatrix} \\ b &= 1 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + 1 \begin{pmatrix} 1 \\ \epsilon \end{pmatrix} = \begin{pmatrix} 2 \\ \epsilon \end{pmatrix} \\ A^{-1} &= \begin{pmatrix} 1 & -1/\epsilon \\ -1/\epsilon & 2/\epsilon^2 \end{pmatrix} \quad W = A^{-1}b = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \end{aligned}$$

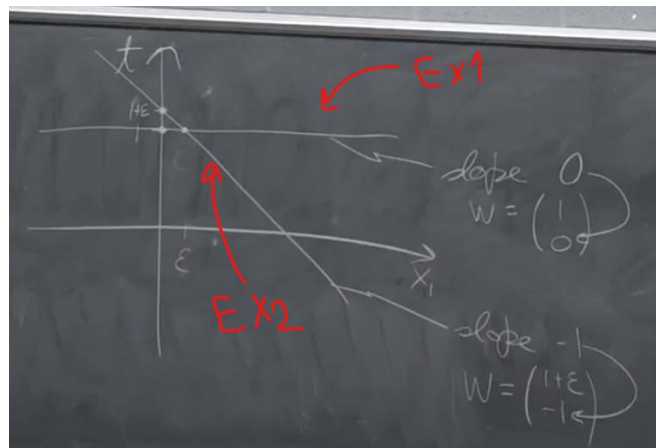
تعال نشوف ثاني مثال .. الفرق هنا اننا حطينا تارجت $1+\epsilon$ بدل ما كان 1 .. تغيير بسيط .. تعال نشوف التأثير عامل ازاي ..

$$\begin{aligned} \text{Ex:2} \\ A &= \begin{pmatrix} 2 & \epsilon \\ \epsilon & \epsilon^2 \end{pmatrix} \quad A^{-1} = \begin{pmatrix} 1 & -1/\epsilon \\ -1/\epsilon & 2/\epsilon^2 \end{pmatrix} \\ b &= (1+\epsilon) \begin{pmatrix} 1 \\ 0 \end{pmatrix} + 1 \begin{pmatrix} 1 \\ \epsilon \end{pmatrix} = \begin{pmatrix} 2+\epsilon \\ \epsilon \end{pmatrix} \\ W &= A^{-1}b = \begin{pmatrix} 1+\epsilon \\ -1 \end{pmatrix} \end{aligned}$$

قيم ال W اتغيرت .. اول رقم بدل 1 بقي $1+\epsilon$.. انما ثاني رقم راح من 0 ل -1 .. تغيير بسيط .. هوب قام ضاربك الماجنتيود ... فديه مشكله مينفعش انها تحصل .. الدكتور هيرسم صوره ثانيه توضح ليه ده بيحصل .. اول حاجه هنرسم x_1 علي ال x axis و ال y axis هو التارجت .. خذ ال ϵ مع ال 1 و الصفر مع ال $1+\epsilon$.. فهنرسم خط ال solution ... في الحل الي اترسم ده انت عندك slope قيمتو -1 .. فبالتالي الحل

$$W = (1+\epsilon)$$

-1)



ازاي نتعامل مع المشكله ديه .. هو ال regularization .. هنعط penalty term هيجاول ي constrain نوع الحل اللي احنا ناخدو في الاعتبار .. الهدف الاساسي ان كنا ن minimize the Euclidean distance ما بين الأوتبوت وبين التارجت .. بس دلوقت كمان احنا هنقل ال magnitude بتاع ال weight ... ماهو لو انت عندك weights كبيره مقارنة بأرقام تانيه .. كذا انت رايح ف داهيه .. فكرة ال regularization هتوضح المحاضره الجايه... فمتستعجلش

Regularization

- Idea: favor smaller values
- Tikhonov regularization: add $\|w\|_2^2$ as a penalty term
- Ridge regression:

$$w^* = \underset{w}{\operatorname{argmin}} \frac{1}{2} \sum_{n=1}^N (t_n - w^T \bar{x}_n)^2 + \frac{\lambda}{2} \|w\|_2^2$$

where λ is a weight to adjust the importance of the penalty

$\|w\|_2^2$.. السكوير بتاع ال Euclidian norm .. واللمضه ديه هايبر بارمتر تعملو أوبتيمائيزيشن .. ايه اللي هيحصل لما نضيف ال penalty term ... السستم لسه linear .. بس ال solution اتغير ..

Regularization

- Solution: $(\lambda I + A)w = b$
- Notes
 - Without regularization: eigenvalues of linear system may be arbitrarily close to 0 and the inverse may have arbitrarily large eigenvalues.
 - With Tikhonov regularization, eigenvalues of linear system are $\geq \lambda$ and therefore bounded away from 0. Similarly, eigenvalues of inverse are bounded above by $1/\lambda$.

الدكتور رسم ع الصبوره بالطباشير ثاني ...

$$\text{Let } \lambda = 0.05 \text{ \& } \varepsilon = 0.1$$

$$(\lambda I + A) = \begin{pmatrix} 2+\lambda & \varepsilon \\ \varepsilon & \varepsilon^2 + \lambda \end{pmatrix} = \begin{pmatrix} 2.05 & 0.1 \\ 0.1 & 0.06 \end{pmatrix}$$

$$(\lambda I + A)^{-1} = \begin{pmatrix} 0.531 & -0.885 \\ -0.885 & 18.1416 \end{pmatrix}$$

لو رجعنا تاني لمثال رقم 1 و 2

<p>Ex: 1</p> $b = \begin{pmatrix} 2 \\ \varepsilon \end{pmatrix} = \begin{pmatrix} 2 \\ 0.1 \end{pmatrix}$ $W = \begin{pmatrix} 0.9735 \\ 0.0442 \end{pmatrix}$	<p>Ex: 2</p> $b = \begin{pmatrix} 2+\varepsilon \\ \varepsilon \end{pmatrix} = \begin{pmatrix} 2.1 \\ 0.1 \end{pmatrix}$ $W = \begin{pmatrix} 1.0265 \\ -0.0442 \end{pmatrix}$
---	--

كدا الارقام بقت قريبيه من بعضها