

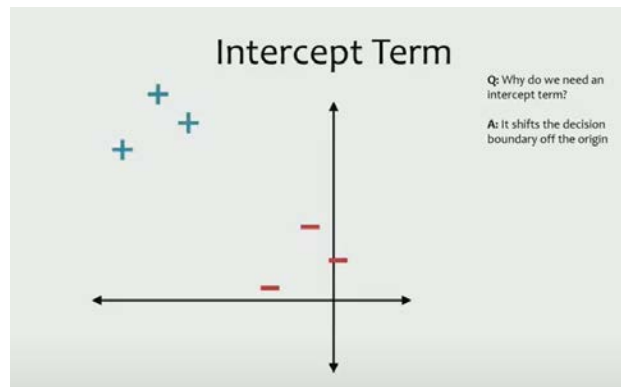
بسم الله الرحمن الرحيم

أهلاً بيكم يا شباب .. انهارد هنفكر في ال linear regression .. وقبل ما نوصل هناك .. هنلف كدا لفه سريعه علي ال perceptrons و هنكتشف هو ليه ال perceptron بيشتغل يعني .. ليه بنتكلم علي ال linear regression ... اول واحد هو اننا نتكلم علي regression algorithm و الفرق ما بين ال regression and classification ان الأوتبوت هيبقا real number بدل ماكان integer .. بالاضافه للسبب الاولاني .. انو بيبقي مثال علي اننا نبدأ نفكر في ال learning problems بطريقه مختلفه و في framework مختلفه .. بحيث هنا انك بتعمل recasting our learning problem كأنها optimization problem .. الفكره ديه اللي هي "Casting learning as an optimization problem is going to carry through a bunch of different methods and models we will think about"

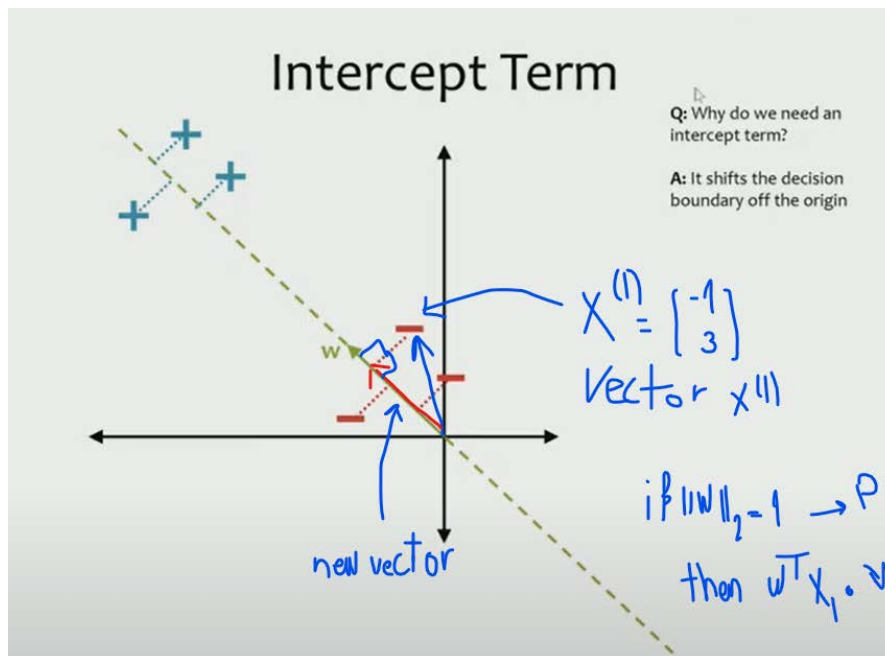
أسيمنت 3 نزل يا جميل .. يلا أبدأ فيه .. هو كلو نظري مفهوش برمجته

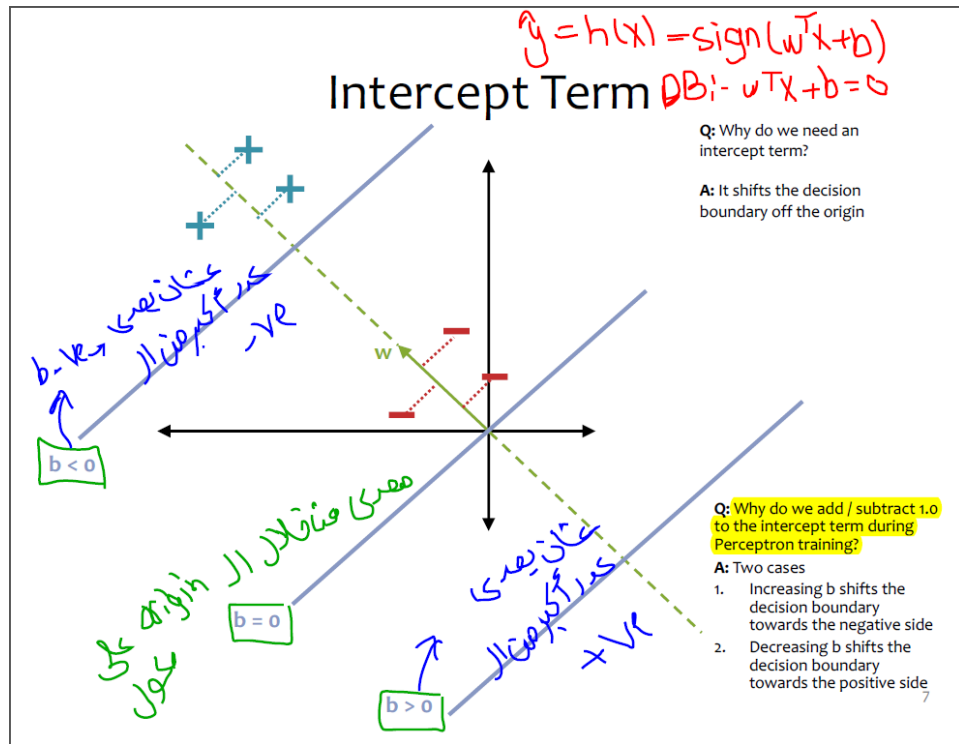
يلا نبدأ

احنا هنبدأ من اننا نفكر ازاى ال intercept term بيلعب دور في ال perceptron algorithm .. اول سؤال .. انت ليه محتاج intercept term اصلاً ... عشان هو بي shift ال decision boundary بعيد عن ال origin .. لو انت مكنتش عندك intercept term .. في الحاله ديه كل ال boundaries هتعتدي بال origin ... تعال نفكر في داتا سيت ..



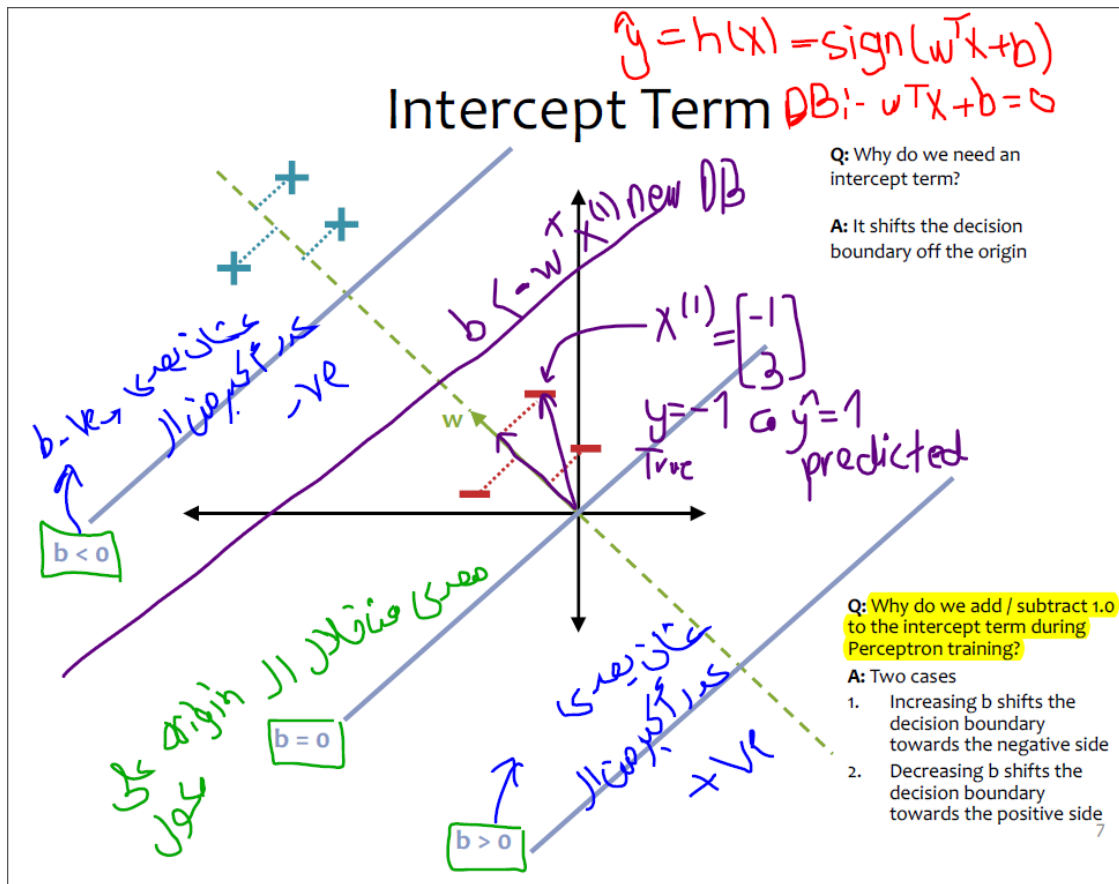
3 و - 3 .. تخيل ان عندك فيكتور اسمو w .. اختارناه ... تخيل انك ترسم خط في نفس الدايركشن بتاع ال w vector رايح للمالانهايه .. أول ما ترسم الخط ده .. عاوزين نفكر هيبقا الشكل عامل ازاى بعد ما نعمل projection لكل النقط اللي عندنا علي الخط الأخضر في الصوره ده ...





الدكتور سأل سؤال ... دلوقت ال b بتساوي صفر .. وچالك إكس 1 .. وال $w^T x_1$ طلعت بوزيتيف .. فدلوقت هي بقت misclassified .. ايه هي القيمة الجديدة اللي هتخلي ال $w^T x_1$ بقت negative لما اجيلها تاني اللفه الجايه في ال iterations بتاعت التريننج ..؟؟

فدلوقت انت بتقول ان $w^T x_1 + b$ بتطلع بوزيتيف .. بيقا انا محتاج اعدل علي ال bias انو بيقا نيجاتيف عشان اقدر احتوي النقطة الجديده اللي جايالي .. ففي الحاله ديه انت عندك معادلة الباوندري كوندشن هي اصلا $w^T x + b = 0$.. فبالتالي انت لوقلت ان $b = -w^T x_1$.. في الحاله ديه الخط الجديد هيعدي علي طول بالنقطه اللي اسمها x_1 ... انما عشان يقدر يحتويها ... فهبقا اقل من ال $-w^T x_1$..



فلما هتبقا قيمة ال bias سالب أكثر .. بيقا هتقدر تحتوي الرقم السالب اللي عندك اللي انت عملتلو misclassification .. طيب اللي تقدر نسألو هنا .. هو ليه ال perceptron مش بيشتغل علي طول .. كل ما يجيلك نقطه و تغلط فيها تقوم حاطط ال $b = +- w^T \cdot x$.. الرد ان اننا مش عاجزين نبعد اوي كذا .. احنا محتاجين نروح بس 1 في الناحيه اللي احنا محتاجينها .. فيا تجمع 1 او تطرح 1

طيب ال inductive bias أول حاجه ان ال decision boundary لازم تبقي linear .. ثاني حاجه ان ال most recent mistakes هي ال most important فيالتالي تخيل عندك a perceptron و عندك 100 مثال .. وبعدين يقوم عامل كلو صح .. بس هوبا دابل كيك .. تجيلو نقطه يعملها غلط .. ايه اللي هيحصل .. يعني احنا علي مدار 1000 مثال .. مغيرناش البارامترز بتاعتنا .. هنجي علي البطيخه ديه ونقوم مغيرنها .. هو فعلاً ده اللي هيحصل .. واحتمال كمان لما يغير تلاقي ال 1000 مثال الي فرحان بيهم دول .. كلهم بقو في الناحية الغلط ... فاحنا محتاجين ن prioritize ال most recent mistakes يعني ..

الدكتور عاجز يستعرض notation trick

Background: Hyperplanes $w^T x + b = 0$

Notation Trick: fold the bias b and the weights w into a single vector θ by prepending a constant to x and increasing dimensionality by one to get x' !

Hyperplane (Definition 1):
 $\mathcal{H} = \{x : w^T x = b\}$

Hyperplane (Definition 2):
 $\mathcal{H} = \{x' : \theta^T x' = 0 \text{ and } x'_0 = 1\}$

$\theta = [b, w_1, \dots, w_M]^T$

Half-spaces:

$\mathcal{H}^+ = \{x : \theta^T x > 0 \text{ and } x_0 = 1\}$

$\mathcal{H}^- = \{x : \theta^T x < 0 \text{ and } x_0 = 1\}$

دلوقت نقدر ن redefine our perceptron algorithm باستخدام التعريف اللي هو ..

(Online) Perceptron Algorithm

Data: Inputs are continuous vectors of length M . Outputs are discrete.

$$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots$$

where $x \in \mathbb{R}^M$ and $y \in \{+1, -1\}$

Prediction: Output determined by hyperplane.

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

$$\text{sign}(a) = \begin{cases} 1, & \text{if } a \geq 0 \\ -1, & \text{otherwise} \end{cases}$$

Assume $\theta = [b, w_1, \dots, w_M]^T$ and $x_0 = 1$

Learning: Iterative procedure:

- **initialize** parameters to vector of all zeroes
- **while** not converged
 - **receive** next example $(x^{(i)}, y^{(i)})$
 - **predict** $y' = h(x^{(i)})$
 - **if** positive mistake: **add** $x^{(i)}$ to parameters
 - **if** negative mistake: **subtract** $x^{(i)}$ from parameters

احنا عندنا شوية إكسات .. وعندنا شوية إيات .. الواي ديسكريت يا واد هه .. واحنا هنا هنفترض ان ال output متحدد بهايير بليين اللي هو $\text{sign}(\theta^T x)$.. هنا بقا اول حاجه حظ البارمترز في الاول بصفر .. وبعدين ابدأ لف .. وانت بتلف بتبص علي البريدكشن بتاعك .. لو كان بوزيتيف . اجمع ال x_i مع الثيتا .. لو نيجاتيف اطرحه من الثيتا .. لإن دلوقت لما هنتيجي تجمع الإكسات .. انت حرفيا بتجمع 1 علي ال b .. لإن أول entry من الإكس اللي عندك هو 1 .. ولما هنتيجي تطرح فانت فعلياً هتطرح 1 من ال b ..

(Online) Perceptron Algorithm

Data: Inputs are continuous vectors of length M . Outputs are discrete.

$$(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$$

$$\text{where } \mathbf{x} \in \mathbb{R}^M \text{ and } y \in \{+1, -1\}$$

Prediction: Output determined by

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

$$\text{Assume } \theta = [b, w_1, \dots, w_M]$$

Learning:

Algorithm 1 Perceptron Learning Algorithm

```

1: procedure PERCEPTRON( $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}$ )
2:    $\theta \leftarrow 0$ 
3:   for  $i \in \{1, 2, \dots\}$  do
4:      $\hat{y} \leftarrow \text{sign}(\theta^T \mathbf{x}^{(i)})$ 
5:     if  $\hat{y} \neq y^{(i)}$  then
6:        $\theta \leftarrow \theta + y^{(i)} \mathbf{x}^{(i)}$ 
7:   return  $\theta$ 

```

Implementation Trick: same behavior as our “add on positive mistake and subtract on negative mistake” version, because $y^{(i)}$ takes care of the sign

▷ Initialize parameters
 ▷ For each example
 ▷ Predict
 ▷ If mistake
 ▷ Update parameters

1 or -1

والصوره اللي فانت ديه انضف طريقه ت implement بيها ال algorithm ده ..

دلوقت هنخش في ال batch perceptron algorithm .. وهنسيب ورا ضهرنا ال online learning setting بتاع ال perceptron .. وهنرجع لل setting بتاع لما كان عندنا \mathcal{D} fixed training set .. ففي ال training time ... ابعت كل الداتا سبت بتاعتك للي هي بتتكون من N training example .. نفس الالجورزم المستخدم هو هو بس انت هتفضل تعيد وتزيد وتعيد وتزيد لحد ما ت converge ..

(Batch) Perceptron Algorithm

Learning for Perceptron also works if we have a fixed training dataset, \mathcal{D} . We call this the “batch” setting in contrast to the “online” setting that we’ve discussed so far.

Algorithm 1 Perceptron Learning Algorithm (Batch)

```

1: procedure PERCEPTRON( $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$ )
2:    $\theta \leftarrow 0$ 
3:   while not converged do
4:     for  $i \in \{1, 2, \dots, N\}$  do
5:        $\hat{y} \leftarrow \text{sign}(\theta^T \mathbf{x}^{(i)})$ 
6:       if  $\hat{y} \neq y^{(i)}$  then
7:          $\theta \leftarrow \theta + y^{(i)} \mathbf{x}^{(i)}$ 
8:   return  $\theta$ 

```

بس الدكتور هنا قال .. لازم تكون سألت نفسك سوالين .. يعني ايه الالجورزم ي converge .. و هل الالجورزم ده هي converge اصلا .. بس قبل ما نجابو هنعوز نعرف شوية حاجات و بعدين نرجع ثاني ..

الدكتور راح اتكلم سريعاً علي ان ال Perceptron algorithm ممكن تعمله derivation عن طريق طريقتين ... اول واحده انك ت extend the linear hinge loss علي ال online perceptron algorithm to the batch setting .. ثاني طريقه انك ت apply SGD .. عشان تقلل ال separator في extensions of perceptron ..

Extensions of Perceptron

- **Voted Perceptron**
 - generalizes better than (standard) perceptron
 - memory intensive (keeps around every weight vector seen during training, so each one can vote)
- **Averaged Perceptron**
 - empirically similar performance to voted perceptron
 - can be implemented in a memory efficient way (running averages are efficient)
- **Kernel Perceptron**
 - Choose a kernel $K(x', x)$
 - Apply the **kernel trick** to Perceptron
 - Resulting algorithm is **still very simple**
- **Structured Perceptron**
 - Basic idea can also be applied when y ranges over an exponentially large set
 - Mistake bound **does not** depend on the size of that set

Linear & non-linear DB

14

Perceptron Exercises

Question:

The parameter vector w learned by the Perceptron algorithm can be **written as a linear combination** of the feature vectors $x^{(1)}, x^{(2)}, \dots, x^{(N)}$.

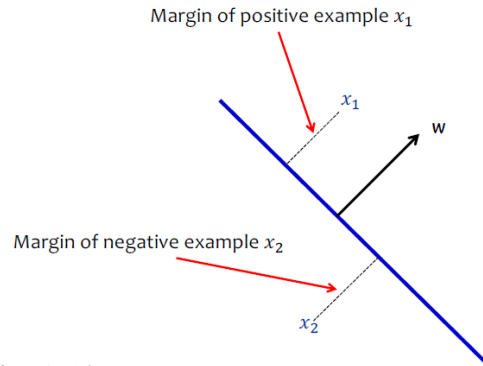
- $w = \alpha_1 x^1 + \alpha_2 x^2 + \dots + \alpha_N x^N$
- True, if you replace "linear" with "polynomial" above
 - True, for all datasets
 - False, for all datasets
 - True, but only for certain datasets
 - False, but only for certain datasets

15

دلوقت عاوزين نشوف ازاى الحاجات ديه بتشتغل .. فهحتاج نخط شوية تعريفات كدا ... أول واحد هو ال geometric margin .. و هنا انت هتعرف 3 مارجنز .. أول واحد بيقول .. "The margin of an example with respect to a linear separator W is the distance from X to the plane ... plain" .. بمعنى ان ال decision boundary هو ال plane .. والمسافه ديه بتعتمد علي ال W ..

Geometric Margin

Definition: The **margin** of example x w.r.t. a linear sep. w is the distance from x to the plane $w \cdot x = 0$ (or the negative if on wrong side)



Slide from Nina Balcan

تاني تعريف هو ان ال margin gamma w لشوية set of examples هنسميهم S with respect to a linear separator W .. انت عندك W اللي بتديك الديسيجن باوندرى .. دلوقت تخيل انك بت رسم 2 parallel lines .. ف فكر فيها كالاتي .. انت عندك W اللي بتديك الديسيجن باوندرى .. دلوقت تخيل انك بت رسم 2 parallel lines علي الديسيجن باوندرى بتاعك .. لحد ما يخط في نقطه من النقط اللي عندك .. واقف هنا واكتب المسافه اللي انت خدتها .. وهي ديه المارجن بتاعتك ..

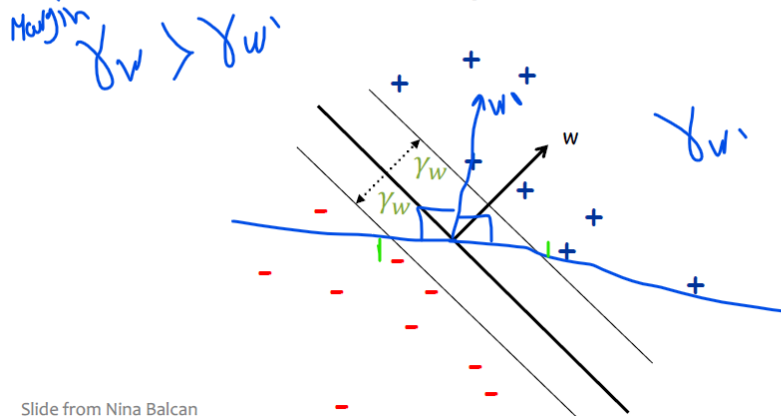
تالت تعريف .. المارجن جاما ل set of examples بالنسبه لل possible linear separators .. المارجن جاما W علي كل ال possible linear separators

تخيل عندك 2 linear separator

Geometric Margin

Definition: The **margin** of example x w.r.t. a linear sep. w is the distance from x to the plane $w \cdot x = 0$ (or the negative if on wrong side)

Definition: The **margin** γ_w of a set of examples S wrt a linear separator w is the smallest margin over points $x \in S$.



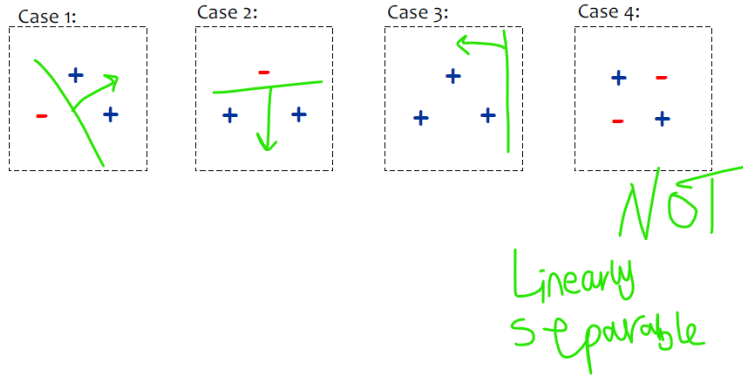
Slide from Nina Balcan

أكبر جاب ما بين ال + و ال - .. خد نص المسافه بتاعت المساحه الفاضيه ديه .. وهي ديه المارجن ..

هنحط تعريف كمان: في البايئري كلاسيكيشن بروبلم .. احنا بنقول linearly separable ... بييقوا decision boundary لو في يقدر يفصل النقط اللي عندنا ... مثلاً::

Linear Separability

Def: For a **binary classification** problem, a set of examples S is **linearly separable** if there exists a linear decision boundary that can separate the points



20

طيب تعال نتكلم علي ال perceptron mistake bound .. ده guarantee ا بيقول ان لو الداتا ليها margin gamma .. وكل النقط واقعه جوا كوره نص قطرها R ... والكوره متسنتره عند ال origin .. في الحاله ديه ال perceptron هيعمل اقل من أو بيساوي R/γ^2 غلطات .. ده بيقول ايه الكلام ده .. بيقول ان في خاصيه مهمه للألجورزم ده ... فأخيراً احنا هنا هنقدر نعرف ال converge هنا .. هنقول ان ال patch perceptron algorithm خلاص converged لما يبطل يعمل غلطات علي التريننج داتا .. perfectly classify the training data ... طب ايه اللي بيضمنهولك ال perceptron mistake bound .. بيضمنلك ان ل linearly separable dataset .. في الآخر ال perceptron algorithm هيبطل يعمل غلطات .. ففي الحاله ديه مش هي update the parameters again .. فهنا هنقول اننا converged خلاص ..

Analysis: Perceptron

Perceptron Mistake Bound

Guarantee: If data has margin γ and all points inside a ball of radius R , then Perceptron makes $\leq (R/\gamma)^2$ mistakes.

(Normalized margin: multiplying all points by 100, or dividing all points by 100, doesn't change the number of mistakes; algo is invariant to scaling.)

Def: We say that the (batch) perceptron algorithm has **converged** if it stops making mistakes on the training data (perfectly classifies the training data).

Main Takeaway: For **linearly separable** data, if the perceptron algorithm cycles repeatedly through the data, it will **converge** in a finite # of steps.

Slide adapted from Nina Balcan

22

تعال نكتب الكلام ده كرياضه ...

Analysis: Perceptron

Perceptron Mistake Bound

Theorem 0.1 (Block (1962), Novikoff (1962)).

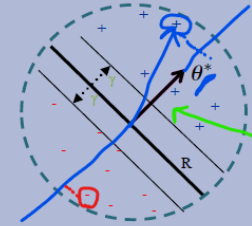
Given dataset: $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$.

Suppose:

1. Finite size inputs: $\|\mathbf{x}^{(i)}\| \leq R$
2. Linearly separable data: $\exists \theta^*$ s.t. $\|\theta^*\|_2 = 1$ and $y^{(i)}(\theta^* \cdot \mathbf{x}^{(i)}) \geq \gamma, \forall i$

Then: The number of mistakes made by the Perceptron algorithm on this dataset is

$$k \leq (R/\gamma)^2$$



$$\text{len} = \theta^* \cdot x$$

$\gamma > 0 \rightarrow \text{margin for } \theta^*$

Figure from Nina Balcan

23

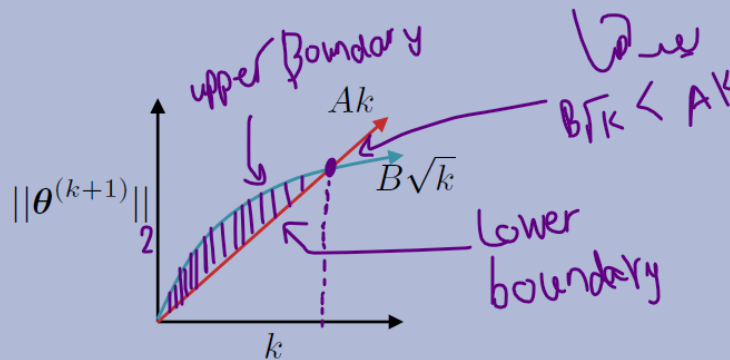
تعال نبص علي الاثبات .. بتاع ان ال perceptron mistake bound ... ال $\theta^{(k+1)}$ البارمترز بعد ال kth mistake ... احنا هنجيب upper bound and lower bound .. ال L2 norm of the parameters after the kth mistake .. ال upper bound هيقا قيمة B مضروب في ال \sqrt{k} .. ايه بقا هم الكيرفين اللي عندك ..

Analysis: Perceptron

Proof of Perceptron Mistake Bound:

We will show that there exist constants A and B s.t.

$$Ak \leq \|\theta^{(k+1)}\| \leq B\sqrt{k}$$



25

A هي السلوب بتاع الخط .. منعرفش ايه هي ال B .. بس هو هيقا كيرف وخلص يعني طالع اسرع من ال A .. بس بسبب ال sqrt هيبطأ و ي cross .. ال statement true بس في قيمه عند ال k هيحصل الكروس .. عندها ال statement isn't true .. فبالنالي مينفعش يكون في قيمه ال k تعدي الكروس بتاع الكيرفين .. فأكبر رقم ال k .. هو هيقا الرقم اللي بيتقاطعو عندو ..

يلا علي الريجريشن بسرعه عشان زهقت .. الهدف واضح وبسيط .. عندك شوية داتا .. و الإكس هو فيكتور و الواي هو سكيلر .. اتعلم فانكشن اللي هي كيرف .. اللي هتبقا احسن حاجه ... خد بالك ان الواي هنا هي scalar .. مش integer .. قبل كدا الواي كانت في الكلاسيفيكيشن ديسكريت ... الواي دلوقت هي رقم حقيقي ... هنا ال $h(x)$ مش هتبقا كلاسيفير عادي .. ده هيديك كيرف .. امثله كتيره الحقيقه .. انا نمت في آخر المحاضره .. فابقا اسمعها ثاني . هه