

انهارده هنتكلم عن ال binary logistic regression .. هنتكلم بعدها ال multinomial logistic regression وده مختلف عن ال BLR ... آخر مره قلنا البرنسبل بتاع ال MLE .. ان عندنا داتا من شوية أمثله  $X_i$  واننا هنتختار البارمترز اللي بت maximize likelihood of data تحت ال iid assumption فده معناه اننا هنتختار البارمترز Theta اللي هت argmax ال  $L(\theta)$  اللي هي ال product of the probability of each of the examples .. ولو عاوز تفكر ف الكلام ده من ناحية الماشين ليرننج برسبيكتيف اللي هو ال classification .. هنا ال  $X_i$  هي ال general data examples انما في الماشين ليرننج ال  $X_i$  بتبقا pairs of feature vector  $(X_i, Y_i)$  ... واللي بنعملو لما بنحاول نحل مشكلة ال argmax .. ان عندنا بارمترز و عندنا فانكشن اسمها  $L(\theta)$  وبنحاول نلاقي الثيتا اللي بت maximize الفانكشن ديه .. الدكتور بيتكلم في السلايد و الرسمتين اللي تحت .. بيقول ان اللي علي الشمال ديه في ال 1D .. اللي علي اليمين في ال 2D .. في ال 2D عندك ثيتا 1 و ثيتا 2 .. والفانكشن بتاعتك في ال D2 .. واحنا بنحاول ن climb to the top of the nearest hill .. ده هو ال setting for finding the maximum likelihood estimate .. الي احنا هنعملو قدام هو اننا هنقول احنا ضيعنا كل الوقت ده في عشان نحط ال setting بتاع ال optimization algorithm .. والي هي اصلا بتعرف ازاي ت walk downhill .. بس هنا انت عاوز ت maximize .. فازاي تقدر تستخدم الجوردم بيعرف يتحرك ناحية ال nearest valley عشان يحل مشكله شكلها عامل زي ال hill .. المنطقي انك تاخذ الفانكشن بتاعتك اللي عاوز ت maximize .. وتقوم مديها اشاره سالب .. بيقا عملت ال flip اللي يخليك تتحرك downhill فانت برضو هنا بقي هدفك انك ت minimize the negative function

## MLE

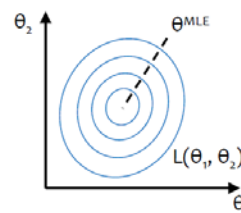
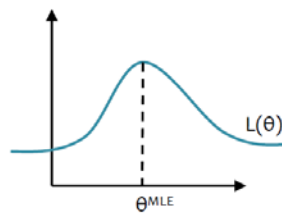
Suppose we have data  $\mathcal{D} = \{x^{(i)}\}_{i=1}^N$

### Principle of Maximum Likelihood Estimation:

Choose the parameters that maximize the likelihood of the data.

$$\theta^{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^N p(\mathbf{x}^{(i)} | \theta)$$

Maximum Likelihood Estimate (MLE) ( $\hat{x}^{(i)}, \hat{y}^{(i)}$ )



5

ال MLE بيحاول انو يحجز as much probability mass as possible للفكره اللي احنا شفناها.. الفكره كلها انك تحاول ت fit the data .. و MLE هيوديك ناحية الأوفر فينتنتج للداتا اللي انت لاحظتها ..

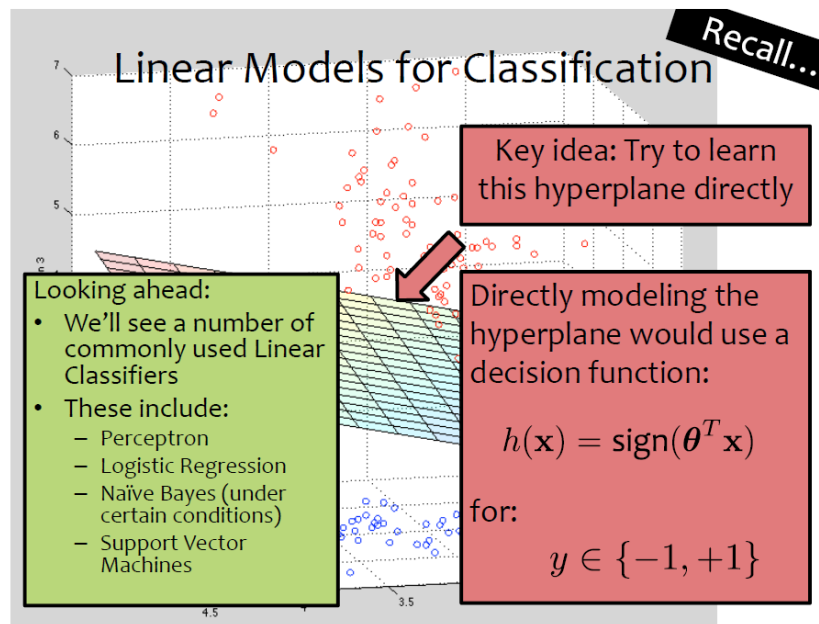
## MLE

What does maximizing likelihood accomplish?

- There is only a finite amount of probability mass (i.e. sum-to-one constraint)
- MLE tries to allocate **as much** probability mass **as possible** to the things we have observed...

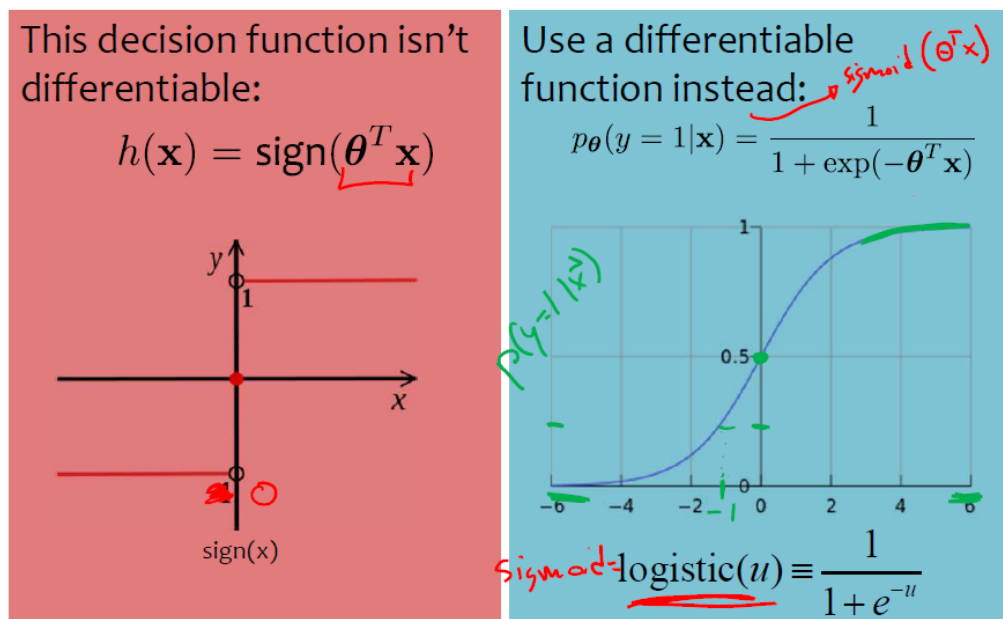
... at the expense of the things we have **not** observed

احنا هنستخدم ال MLE عشان نتعلم logistic regression model .. احنا هنفترض ان الأوتوتس هتبقا باينري .. ف في ال Logistic regression .. خد بالك اننا دلوقت برضو بنتكلم علي .. classification بعض النظر عن الاسم يعني اللي هو regression .. وهنا ال classification هيتعلم linear model زي البرسيبترون كذا ... و علي طول استخدم مثلا حاجه زي ال  $\text{sign}(\theta^T x)$  .. احنا انهارد هنعرف ال linear classifier اللي هو ال linear regression .. وبعدين هنعرف ال objective function .. وديه هتيجي من ال MLE وبعدين optimize using gradient descent عشان تتعلم البارامترز .. وبعدين ت predict the class with highest probability under the model

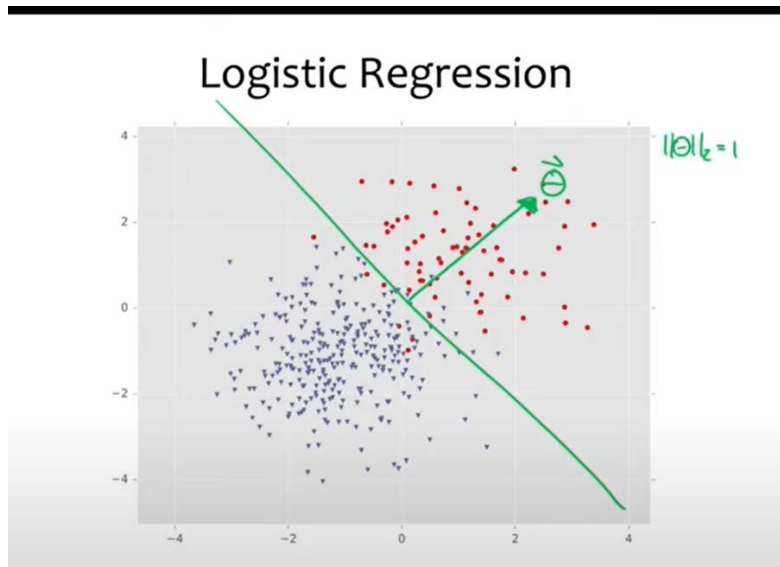


طيب في الطريق ده احنا هيقابلنا مشكله .. ال decision function not differentiable ... بس احنا قلنا اننا هنن ال gradient descent علي الفانكشن ديه .. فهحتاج نحسب ال partial deriv. .. بس هنا مش هنعرف نشق .. خلاص هنستخدم a differentiable function .. و هي ديه ال sigmoid .. او اللوجستيك فانكشن ..

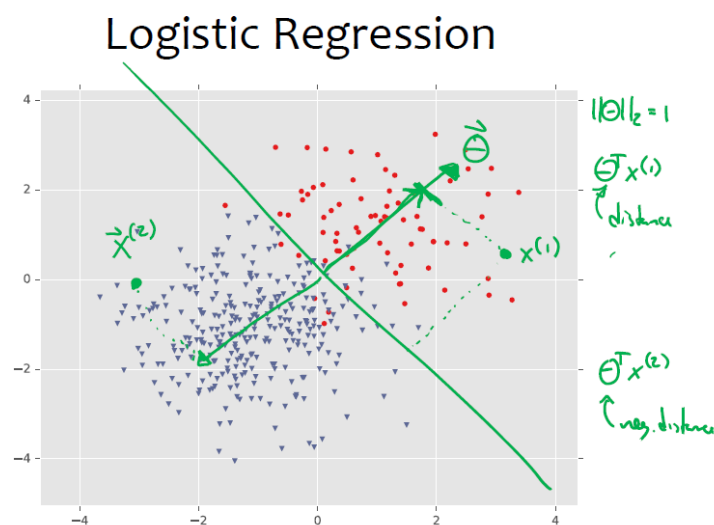
## Using gradient ascent for linear classifiers



تعال نبص علي الداتا سيت .. وتخيل ان عندنا DB لينير وعندك الفيكتور ثيتا .. في الدايركشن بتاع ال + إكزامبلز ..



افترض عندك نقطه اسمها  $x_1$  .. وانت عاوز تحسب قيمة الـ  $\theta^T x_1$  .. الدوت برودكت ده بيعبر عن الماجنتيود بتاع البروجكشن ... وده فالدوت برودكت هو المسافه ما بين الإكس 1 للثيتا فيكتور .. طيب تعال نخط نقطه ثانيه اسمها  $x_2$  .. تعال نوقعها علي الثيتا .. هيديلك فيكتور في الاتجاه الثاني .. في الحاله ديهِ الثيتا ترانسبوز إكس 2 هيبقا في النيجاتيف دايركشن ..



## Logistic Regression

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

$$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N \text{ where } \mathbf{x} \in \mathbb{R}^M \text{ and } y \in \{0, 1\}$$

**Model:** Logistic function applied to dot product of parameters with input vector.

$$p_{\theta}(y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})}$$

**Learning:** finds the parameters that minimize some objective function.  $\theta^* = \underset{\theta}{\operatorname{argmin}} J(\theta)$

**Prediction:** Output is the most probable class.

$$\hat{y} = \underset{y \in \{0,1\}}{\operatorname{argmax}} p_{\theta}(y|\mathbf{x})$$

1. احنا المره اللي فاتت قلنا مش هنستخدم  $x$  .. خلينا بقا نستخدم  $x$  المره ديه
2. الموديل : هنعرف الموديل هنا .. ال  $y$  يا صفر يا واحد ... فالتالي احتمالية ان الواي تبقا بواحد أو صفر هيبقا واخد. bernolli distr. اللي هو عندك بارمتر إسمو فاي .. يقلك  $y = 1$  how likely we are to get  $y = 1$  و  $1 - \phi$  هيبقا هو الاحتمالية ان الواي تبقا بصفر ... احنا المره اللي فاتت احنا اتعلمنا البارمتر فاي علي طول .. المره ديه هنعرف الفاي انها فانكشن في الثيتا و الإكس .. وهنشتغل علي السجمويد في الحاله ديه ..  $\text{sigmoid}(u) = 1 / (1 + \exp(-u))$  بحيث ان ال  $u = \text{Theta} \cdot X$  .... تعال بقا نحط  $p(y|x)$  هيبقا ليه قيمتين .. لو الواي بواحد هتبقا السجمويد .. والواي بصفر هتبقا واحد ماينص السجمويد ... ده كذا الموديل بتاعنا
3. الأوبجكتيف .. -ve conditional log likelihood .. هتكتب ال  $L(\text{theta})$  ... هتبقا زي الصوره اللي مخطوطه .. خد بالك ان  $p(y|x, \text{theta})$  هي بتتقال بروبابلتي أوف واي جيفن الإكس و البارمترز ثيتا .. الدكتور ساعات ممكن يكتبها زي الي باللون الاخضر اللي هو يحط الثيتا في ال subscript ده عشان يعني كإنك بتقول ان الثيتا هي البارمترز بتاعت البروبابلتي دستريبيوشن ..

## Binary Logistic Regression

① Assume: Use  $\vec{x}$

② Model:  $\phi = \sigma(\vec{\theta}^T \vec{x})$  where  $\sigma(u) = \frac{1}{1 + \exp(-u)}$   
 $y \sim \text{Bernoulli}(\phi)$

$$p(y|\vec{x}) = \begin{cases} \sigma(\vec{\theta}^T \vec{x}) & \text{if } y=1 \\ 1 - \sigma(\vec{\theta}^T \vec{x}) & \text{if } y=0 \end{cases}$$

③ Objective:

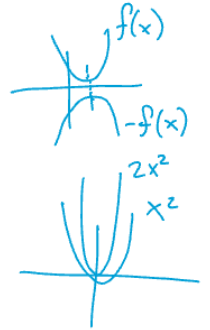
$$l(\vec{\theta}) = \sum_{i=1}^N \log p(y^{(i)} | x^{(i)}, \vec{\theta})$$

$$J(\vec{\theta}) = -\frac{1}{N} l(\vec{\theta}) = -\frac{1}{N} \sum_{i=1}^N -\log p(y^{(i)} | x^{(i)}, \vec{\theta})$$

$J(\vec{\theta})$  for SGD

★ negative average conditional log-likelihood for Log. Reg. is convex

$$\begin{aligned} \theta_{MLE} &= \underset{\theta}{\text{argmax}} l(\theta) \\ &= \underset{\theta}{\text{argmin}} -l(\theta) \\ &= \underset{\theta}{\text{argmin}} -\frac{1}{N} l(\theta) \end{aligned}$$

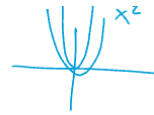


تعال نكتب الأوبجكتيف فانكشن  $J(\text{theta})$  .. هتبقا  $1/N$  مضروب في ال  $L(\text{theta})$  ... هنشتغل بال -ve average conditional log likelihood .....  
 وده هيبقا convex function .. الدكتور حط السالب جوا الساممشن .. وبعدين سأل ليه .. هنستفاد ايه .. عشان نقدر ن minimize ال  $J(\text{theta})$  . بس ليه حطيناه جوا الساممشن . ده عشان انت اللي عاوز تعملو انك تعرف small quantity لكل تريننج إكزامبل اللي نقدر نقلو .. السبب ف حاجه زي كذا .. هو ان عندنا single quantity نقدر نسميها  $J_i(\text{theta})$  و بعدين نقدر نستخدم ال  $J_i(\text{theta})$  لل SGD ... تعال نفكر في ازاي هنستخدم الفانكشن ديه .. انت هتحاول تلاقي البارمترز الثيتا .. اللي هي ال  $\text{argmax } L(\text{theta})$  ... لو خدت النيجاتيف بتاع الفانكشن ... فانت بتاخذ ال  $\text{argmin}$  ... المهم تعال ناخذ المشتقه ونحسب الجريدينس

$$J(\theta) = -\frac{1}{N} \ell(\theta) = -\frac{1}{N} \sum_{i=1}^N -\log p(y^{(i)} | x^{(i)}; \theta)$$

$J^{(i)}(\theta)$  for SGD

$$\begin{aligned} \theta_{MLE} &= \underset{\theta}{\operatorname{argmax}} \ell(\theta) \\ &= \underset{\theta}{\operatorname{argmin}} -\ell(\theta) \\ &= \underset{\theta}{\operatorname{argmin}} -\frac{1}{N} \ell(\theta) \end{aligned}$$



\* need derivatives for learning

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta_m} &= \frac{\partial}{\partial \theta_m} \left( -\log p(y^{(i)} | x^{(i)}; \theta) \right) \\ &= \begin{cases} \frac{\partial}{\partial \theta_m} -\log[\sigma(\theta^T x^{(i)})] & \text{if } y^{(i)} = 1 \\ \frac{\partial}{\partial \theta_m} -\log[1 - \sigma(\theta^T x^{(i)})] & \text{if } y^{(i)} = 0 \end{cases} \\ &= \dots \\ &= \dots \quad \leftarrow \text{recitation} \\ &= - \left( \underbrace{y^{(i)}}_{\text{truth}} - \underbrace{\sigma(\theta^T x^{(i)})}_{\text{prob of } y=1} \right) \underbrace{x_m^{(i)}}_{\text{feature}} \end{aligned}$$

$$\nabla J^{(i)}(\theta) = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} = - \left( y^{(i)} - \sigma(\theta^T x^{(i)}) \right) \vec{x}^{(i)}$$

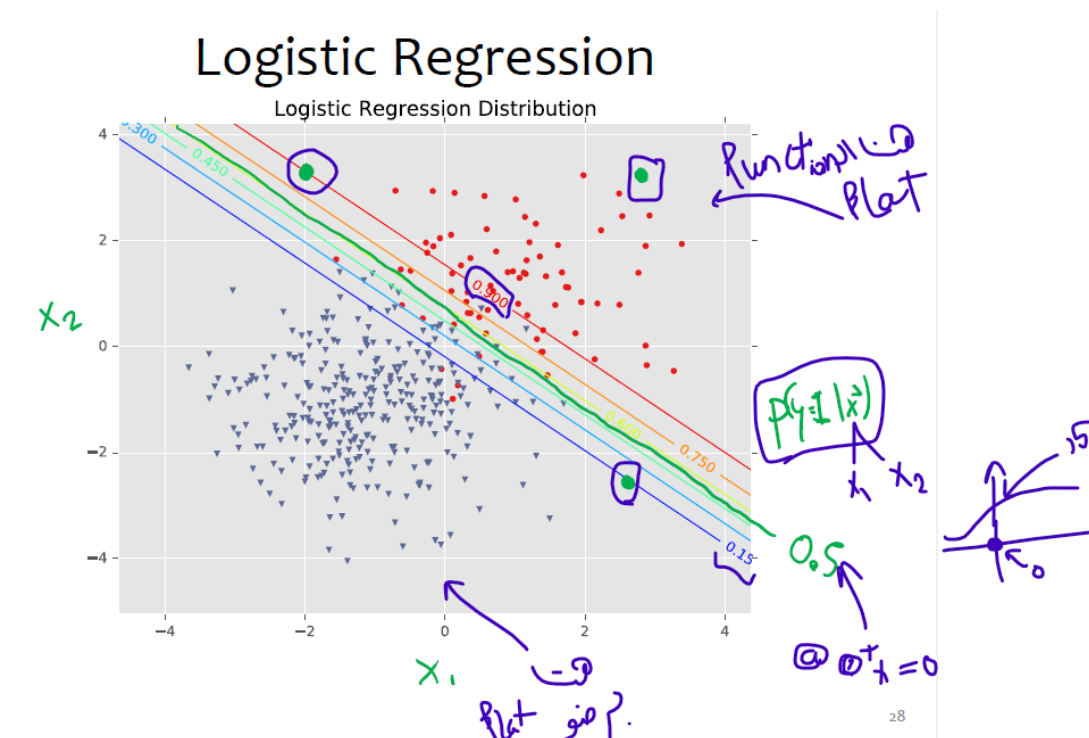
$$\nabla J(\theta) = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} = \frac{1}{N} \sum_{i=1}^N \nabla J^{(i)}(\theta)$$

④ Find  $\hat{\theta}$  by gradient descent or SGD

⑤ Predict

$$\begin{aligned} \hat{y} &= \underset{y \in \{0,1\}}{\operatorname{argmax}} p(y | \vec{x}) \\ &= \begin{cases} 1 & \text{if } p(y=1 | \vec{x}) \geq 0.5 \\ 0 & \text{otherwise} \end{cases} \\ &= \operatorname{sign}(\theta^T x) \in \{0,1\} \end{aligned}$$

الدكتور راح يشوف مثال علي اللوجستيك ريجريشن علي الدات سيت ..



فانت عشان تعمل بريدكشن في الموديل ده .. هو ال  $\text{argmax of } y$ , either 0 or 1 for which it maximizes  $(p(y|x))$  ده هنا غير المجاورتي  
 فووت ... القيمه الأعلى هنا بتعتمد علي الإكس .. عشان الفانكشن فاي ... بتعتمد علي الثيتا والإكس .. في طرق ثانيه تقدر تكتب بيها الإكسبرشن ده ..

④ Find  $\hat{\theta}$  by gradient descent or SGD

⑤ Predict

$$\begin{aligned}\hat{y} &= \underset{y \in \{0,1\}}{\text{argmax}} p(y|\vec{x}) \\ &= \begin{cases} 1 & \text{if } p(y=1|\vec{x}) \geq 0.5 \\ 0 & \text{otherwise} \end{cases} \\ &= \text{sign}(\theta^T x) \in \{0,1\}\end{aligned}$$

فده برضو لينير كلاسيفير ..

## Maximum Conditional Likelihood Estimation

**Learning:** finds the parameters that minimize some objective function.

$$\theta^* = \underset{\theta}{\text{argmin}} J(\theta)$$

We minimize the *negative log conditional likelihood*:

$$J(\theta) = -\log \prod_{i=1}^N p_{\theta}(y^{(i)} | \mathbf{x}^{(i)})$$

Why?

- ~~1. We can't maximize likelihood (as in Naïve Bayes) because we don't have a joint model  $p(x,y)$~~
2. It worked well for Linear Regression (least squares is MCLE)

31

## Maximum Conditional Likelihood Estimation

**Learning:** Four approaches to solving  $\theta^* = \underset{\theta}{\text{argmin}} J(\theta)$

**Approach 1:** Gradient Descent

(take larger – more certain – steps opposite the gradient)

**Approach 2:** Stochastic Gradient Descent (SGD)

(take many small steps opposite the gradient)

**Approach 3:** Newton's Method

(use second derivatives to better follow curvature)

~~**Approach 4:** Closed Form???~~

~~(set derivatives equal to zero and solve for parameters)~~

Logistic Regression does not have a closed form solution for MLE parameters.

33

## SGD for Logistic Regression

### Question:

Which of the following is a correct description of SGD for Logistic Regression?

### Answer:

At each step (i.e. iteration) of SGD for Logistic Regression we...

- 5 A. ☒ compute the gradient of the log-likelihood for all examples (2) update all the parameters using the gradient
- calm B. ☒ (1) ask Matt for a description of SGD for Logistic Regression, (2) write it down, (3) report that answer
- 6 C. ☒ compute the gradient of the log-likelihood for all examples (2) randomly pick an example (3) update only the parameters for that example
- 19 D. (1) randomly pick a parameter, (2) compute the partial derivative of the log-likelihood with respect to that parameter, (3) update that parameter for all examples
- 56 E. ☒ (1) randomly pick an example, (2) compute the gradient of the log-likelihood for that example, (3) update all the parameters using that gradient
- 16 F. (1) randomly pick a parameter and an example, (2) compute the gradient of the log-likelihood for that example with respect to that parameter, (3) update that parameter using that gradient

34

## Gradient Descent

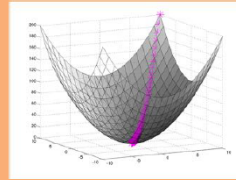
Recall...

### Algorithm 1 Gradient Descent

```

1: procedure 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 Logistic Regression all we need is the **gradient** of the objective function (i.e. vector of partial derivatives).

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

35

## Stochastic Gradient Descent (SGD)

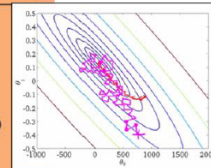
Recall...

### Algorithm 1 Stochastic Gradient Descent (SGD)

```

1: procedure SGD( $\mathcal{D}, \theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:     for  $i \in \text{shuffle}(\{1, 2, \dots, N\})$  do
5:        $\theta \leftarrow \theta - \gamma \nabla_{\theta} J^{(i)}(\theta)$ 
6:   return  $\theta$ 

```



We can also apply SGD to solve the MCLE problem for Logistic Regression.

We need a per-example objective:

$$\text{Let } J(\theta) = \sum_{i=1}^N J^{(i)}(\theta) \\ \text{where } J^{(i)}(\theta) = -\log p_{\theta}(y^i | \mathbf{x}^i).$$

36



# Logistic Regression vs. Perceptron

## Question:

**True or False:** Just like Perceptron, one step (i.e. iteration) of SGD for Logistic Regression will result in a change to the parameters only if the current example is incorrectly classified.

## Answer:

A = accuracy  
B = T  
C = F 100%



يلا نتكلم علي ال mini patch SGD .. في ال gradient descent احنا حسبنا ال gradient من كل الأمثلة .. بالنسبة لل SGD احنا قَرَبنا ال true gradient بالنسبة لل gradient of one randomly chosen example .. في ال mini-patch SGD احنا بنقَرَب ال true gradient عن طريق اتنا ناخذ الأفرديج بتاع ال k randomly chosen examples

## Mini-Batch SGD

- **Gradient Descent:**  
Compute true gradient exactly from all N examples
- **Stochastic Gradient Descent (SGD):**  
Approximate true gradient by the gradient of one randomly chosen example
- **Mini-Batch SGD:**  
Approximate true gradient by the average gradient of  $k$  randomly chosen examples

5

40

لكل الألبورزمز احنا عندنا نفس الفورمات .. أبديت الثيتا بالقيمة بتاعتها ماينص الجاما في ال  $g$  .. في الجريدينت ديسنت .. كانت ال  $g$  هي الفوللل جريدينت .. بالنسبة لل mini-batch SGD .. عندنا الأفرديج بتاع شوية أمثلة  $z$  ال  $i$  ديه هي ارقام 1 و 2 و 3 و 4 .. لحد  $N$  . وديه راندوملي تشوزن إكزامبلز .. وبعدين تاخذ الداتا ديه و تحسب منهم ال gradient ..

## Mini-Batch SGD

while not converged:  $\theta \leftarrow \theta - \alpha g$

## Three variants of first-order optimization:

$$\text{Gradient Descent: } g = \nabla J(\theta) = \frac{1}{N} \sum_{i=1}^N \nabla J^{(i)}(\theta)$$

$$\text{SGD: } g = \nabla J^{(i)}(\theta) \quad \text{where } i \text{ sampled uniformly}$$

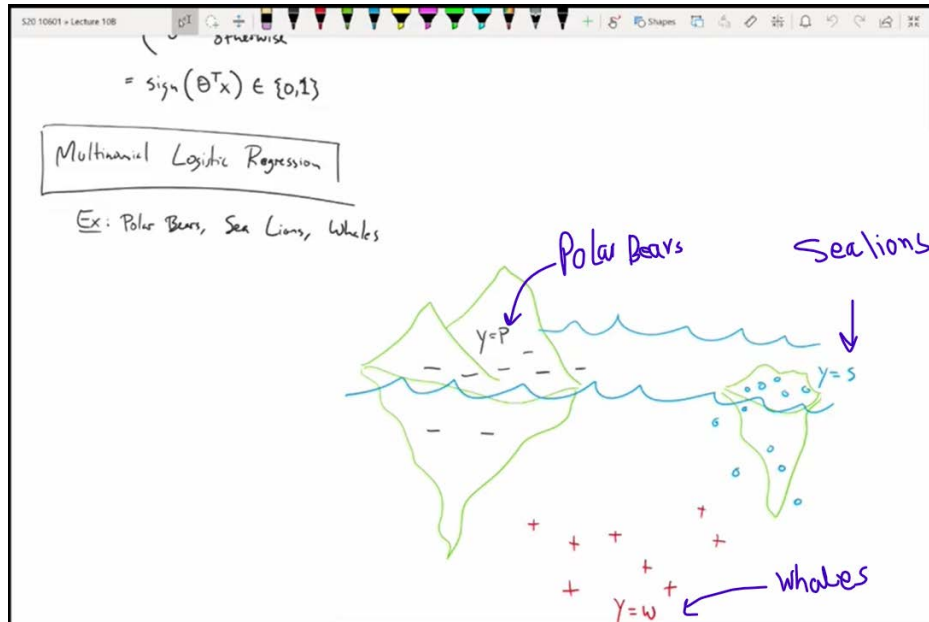
$$\text{Mini-batch SGD: } g = \frac{1}{S} \sum_{s=1}^S \nabla J^{(i_s)}(\theta) \quad \text{where } i_s \text{ sampled uniformly } \forall s$$

#samples  $\rightarrow N$

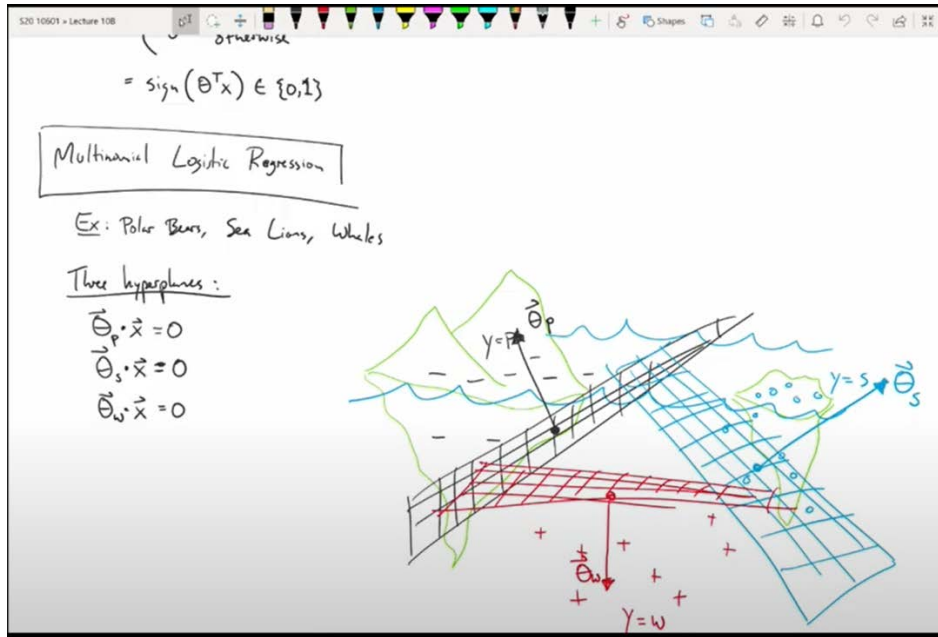


دلوقت هنعرف موديل جديد اسمو Multinomial logistic regression .. وعشان نعمل كذا ... تعال نوجه الشكر لعم dave و عم kevin عشان حرفياً لولا مجهودهم لا كان زمان المحاضرات اتسجلت ولا كان زمانى قعدت اذاكر المنهج بتاعهم .. ربنا بياركلهم ويوفقهم ..

الدكتور بيقول بقا لو عندنا درون و عاوزين نعمل classification ل 3 حاجات .. Polar bears, sea lions, whales ..



الدكتور حط هايير بلينز بتفصل كل واحد عن الآخرين .. فانت عندك 3 هايير بلينز ..



# Multinomial Logistic Regression

Ex: Polar Bears, Sea Lions, Whales

Three hyperplanes:

$$\vec{\theta}_p \cdot \vec{x} = 0$$

$$\vec{\theta}_s \cdot \vec{x} = 0$$

$$\vec{\theta}_w \cdot \vec{x} = 0$$

$$p(y=p|\vec{x}) = \exp(\vec{\theta}_p \cdot \vec{x}) / Z(\vec{x}, \theta)$$

$$p(y=s|\vec{x}) = \exp(\vec{\theta}_s \cdot \vec{x}) / Z(\vec{x}, \theta)$$

$$p(y=w|\vec{x}) = \exp(\vec{\theta}_w \cdot \vec{x}) / Z(\vec{x}, \theta)$$

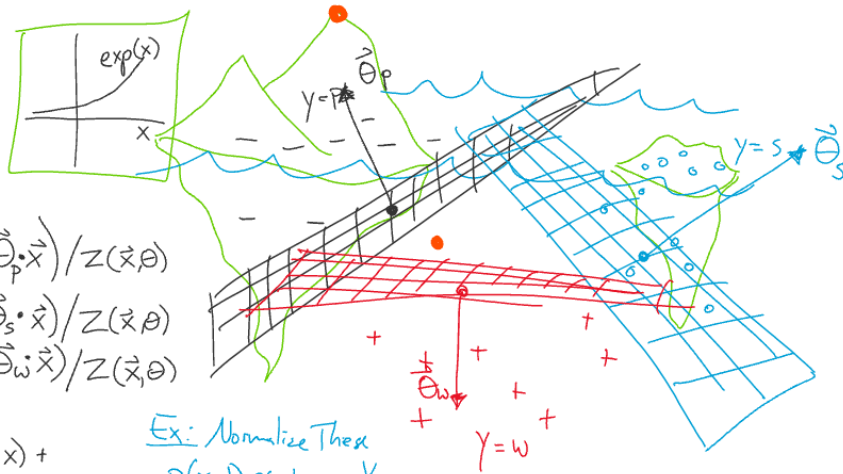
$$Z(\vec{x}, \theta) = \exp(\theta_p \cdot x) + \exp(\theta_s \cdot x) + \exp(\theta_w \cdot x)$$

Ex: Normalize These

$$p(x=1) \propto 1 = 1/5$$

$$p(x=2) \propto 3 = 3/5$$

$$p(x=3) \propto \frac{1}{3} = 1/5$$



Def: Multiclass Classification  $\vec{x} \in \mathbb{R}^M$ ,  $y \in \{1, 2, \dots, K\}$

② Model:  $p(y|\vec{x}) = \frac{\exp(\vec{\theta}_y \cdot \vec{x})}{\sum_{k \in \mathcal{Y}} \exp(\vec{\theta}_k \cdot \vec{x})}$

$\sum_{k \in \mathcal{Y}} \exp(\vec{\theta}_k \cdot \vec{x})$   $Z(\vec{x}, \theta)$  normalization constant

where  $\theta = K \times M$  matrix =  $\begin{bmatrix} \vec{\theta}_1 \\ \vec{\theta}_2 \\ \vdots \\ \vec{\theta}_K \end{bmatrix}$   
polar bear  $\vec{\theta}_1$   
sea lion  $\vec{\theta}_2$   
whale  $\vec{\theta}_K$

$$\phi = [\phi_1, \phi_2, \dots, \phi_K]^T$$

$$\phi_k = p(y=k|\vec{x})$$

$y \sim \text{Categorical}(\phi)$  same as Multinomial with 1 sample } weighted die roll

Def: Multiclass Classification  $\vec{x} \in \mathbb{R}^M$ ,  $y \in \{1, 2, \dots, K\}$

② Model:  $p(y|\vec{x}) = \frac{\exp(\vec{\theta}_y \cdot \vec{x})}{\sum_{k \in \mathcal{Y}} \exp(\vec{\theta}_k \cdot \vec{x})}$

where  $\Theta = K \times M$  matrix =  $\begin{bmatrix} \vec{\theta}_1 \\ \vec{\theta}_2 \\ \vdots \\ \vec{\theta}_K \end{bmatrix}$   
 pole  
 sen  
 while

$z(\vec{x}, \Theta)$   
 normalizes  
 constant

$$\Theta = [\theta_1, \theta_2, \dots, \theta_K]^T$$

$$\theta_k = p(y=k|\vec{x})$$

$y \sim \text{Categorical}(\Theta)$  ← same as Multinomial with 1 sample weighted die roll

③ Learning by MLE

Conditional Log-Likelihood:

$$\ell(\Theta) = \log \left[ \prod_{i=1}^N p(y^{(i)}|x^{(i)}, \Theta) \right] = \sum_{i=1}^N \log p(\dots)$$

$$J(\Theta) = -\frac{1}{N} \ell(\Theta) \quad \leftarrow \star \text{convex} \star$$

$$\hat{\Theta} = \underset{\Theta}{\text{argmin}} J(\Theta) \quad \text{use GD or SGD}$$

Gradient:

$$\begin{aligned} \frac{\partial J^{(i)}(\Theta)}{\partial \theta_{km}} &= \frac{\partial}{\partial \theta_{km}} (-\log p(y^{(i)}|x^{(i)}, \Theta)) \\ &= \dots \quad \leftarrow \text{HW4} \quad \begin{cases} 1 & \text{if } y^{(i)} = k \\ 0 & \text{otherwise} \end{cases} \\ &= -(\mathbb{1}(y^{(i)} = k) - p(y^{(i)}|x^{(i)}, \Theta)) x_m^{(i)} \end{aligned}$$

$$\frac{\partial J^{(i)}(\Theta)}{\partial \vec{\theta}_k} = \nabla_{\vec{\theta}_k} J^{(i)}(\Theta) = -(\mathbb{1}(y^{(i)} = k) - p(y^{(i)}|x^{(i)}, \Theta)) \vec{x}^{(i)}$$

$$= \begin{bmatrix} \frac{\partial J^{(i)}}{\partial \theta_{k1}} \\ \frac{\partial J^{(i)}}{\partial \theta_{k2}} \\ \vdots \\ \frac{\partial J^{(i)}}{\partial \theta_{kM}} \end{bmatrix} \quad \text{compute for each } k$$

④ SGD for learning/opt.

⑤ predict the most probable class.