

انهارده هنتكلم علي الباك بروب .. احنا فجاء في مكان هنعسب ال gradients بالنسبه لأي differentiable function .. ولما هنعمل كذا .. ال paradigm بتاع ال learning as function approximation هبيان هو عبارته عن ايه .... لأن دلوقت انت تقدر تكتب ال objective و تكتب complicated parameterized function وتستخدم optimization عشان تتعلم parameters كويسه لل objective function ... تعال نشوف لما كنا بنتكلم علي اختيار ال architecture بتاع النورال نتورك كان قدامنا ايه ..

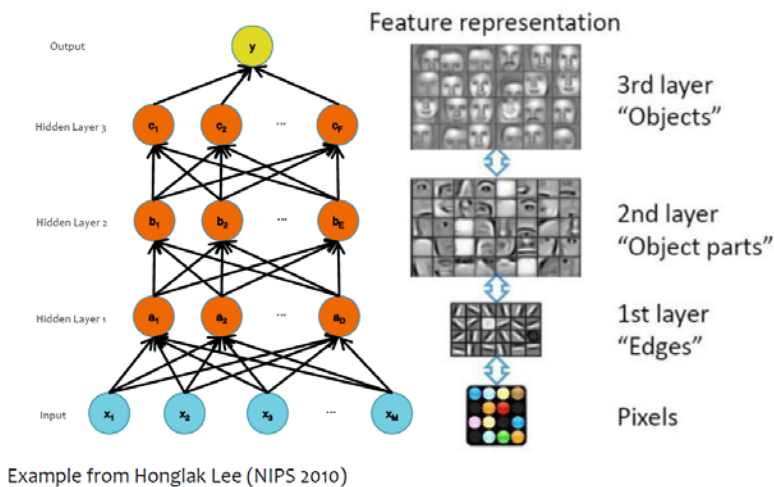
## Neural Network Architectures

Even for a basic Neural Network, there are many design decisions to make:

1. # of hidden layers (depth)
2. # of units per hidden layer (width)
3. Type of activation function (nonlinearity)
4. Form of objective function

الدكتور بيقول ان في سبب ممكن يخليك انك تفكر في layers اكتر .. وهو انك فعلا عارف ان النتورك بتاعك هتعمل حاجه احسن من اللي انت بتقدر تعملو عشان تقدر تجيب level of abstraction معين .. فافتكر في ال feature engineering ... قصة ان ال computer vision بيقدّر يطلعك features of image .. بس ايه اللي هيحصل لو هو اللي عملوده كان غلط ومش شغال ... وكان الطريقه الصح هو انك بيقا عندك initial function تقدر تطلع ال edges وبعدها تاخذ ال edges وتحطهم مع بعض ... فالدكتور رايح ناحيه اننا نستخدم ال NN في الحوار ده

## Different Levels of Abstraction



دلوقت الدكتور بيتكلم علي اختيار الأكتيفيشن فانكشنز .. راح يتكلم علي السجمويد فيقول ان السلوب بتاع السجمويد تقريباً صفر طول منت بتبعد عن الصفر .. وده مشكله عشان انت هتاخذ ال gradients .. فلو ال gradients كانت صغيره .. فانت مش هبيقا عندك progress كبير .. الدكتور بيتكلم علي السلايد الجايه و بيوضح يعني ان ال tanh فعلا بتبقا احسن من ال sigmoid .. الدكتور بيتكلم علي الكيرف الأسود اللي فوق خالص .. بيقولك ان ديه مشكله ال vanishing gradients .... اعتقد هو مسميها sigmoid level 5 ... يمكن المقصود هنا هو sigmoid(sigmoid(sigmoid(sigmoid(sigmoid))))

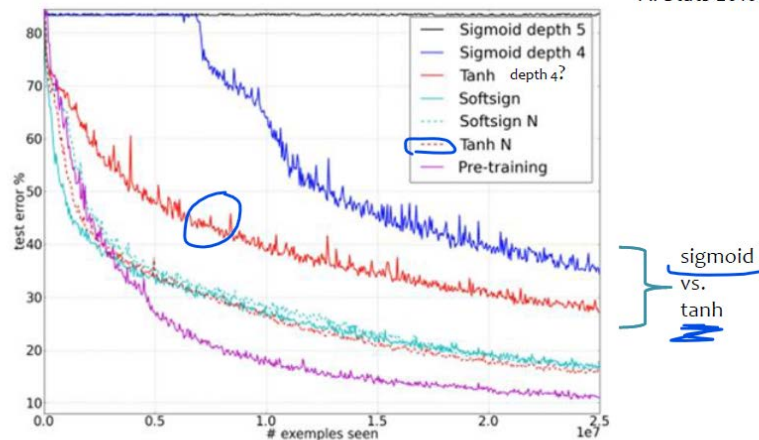
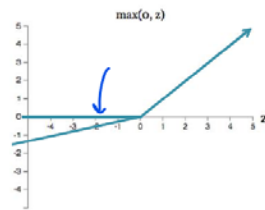


Figure from Glorot & Bentio (2010)

## Activation Functions

- A new change: modifying the nonlinearity
  - ReLU often used in vision tasks



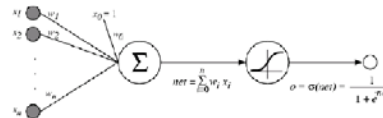
$$\max(0, w \cdot x + b).$$

Slide from William Cohen

Alternate 2: rectified linear unit

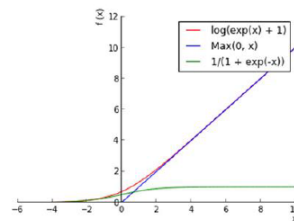
Linear with a cutoff at zero

(Implementation: clip the gradient when you pass zero)



## Activation Functions

- A new change: modifying the nonlinearity
  - ReLU often used in vision tasks

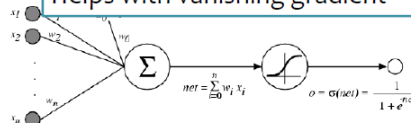


Slide from William Cohen

Alternate 2: rectified linear unit

Soft version:  $\log(\exp(x)+1)$

Doesn't saturate (at one end)  
Sparsifies outputs  
Helps with vanishing gradient



# Objective Functions for NNs

1. Quadratic Loss:
  - the same objective as Linear Regression
  - i.e. mean squared error
2. Cross-Entropy:
  - the same objective as Logistic Regression
  - i.e. negative log likelihood
  - This requires probabilities, so we add an additional “softmax” layer at the end of our network

	Forward	Backward
Quadratic	$J = \frac{1}{2}(y - y^*)^2$	$\frac{dJ}{dy} = y - y^*$
Cross Entropy	$J = y^* \log(y) + (1 - y^*) \log(1 - y)$	$\frac{dJ}{dy} = y^* \frac{1}{y} + (1 - y^*) \frac{1}{y - 1}$

الأوبجكتيف فانكشنز بتاعت ال NNs .. ليهم اسامي مختلفه شوية هنا في ال NN ... الهدف من السلايد ديه هو ال terminology يعني .. حاجه زي ال Quadratic loss .. اللي هي بتبقا  $(\text{true} - \text{predicted})^2$  .. ال  $y^*$  في السلايد هي ال predicted .... ديه هي ال Quadratic loss الي استخدمناها في ال mean squared error في ال linear regression .. عندك cross entropy هو نفس ال objective اللي بنستخدمو عشان ال binary logistic regression و ال multinomial .. هو بس بيبقا ال NLL بتاع الداتا .. حاجه من اللي تلاحظها في ال literature هو ان الناس بتستخدم ال NLL اللي هو ال cross entropy بشكل كتير في ال probabilistic models .. بس انت مش لازم تعمل كدا يعني عادي .. بس في اكثر الأوقات بتلاقى ال cross-entropy بيبقا احسن ...

## Objective Functions for NNs

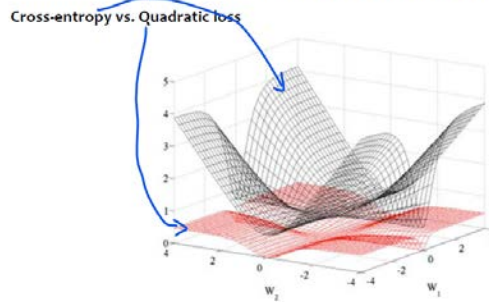
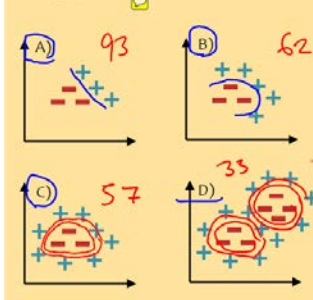


Figure 5: Cross entropy (black, surface on top) and quadratic (red, bottom surface) cost as a function of two weights (one at each layer) of a network with two layers,  $W_1$  respectively on the first layer and  $W_2$  on the second, output layer.

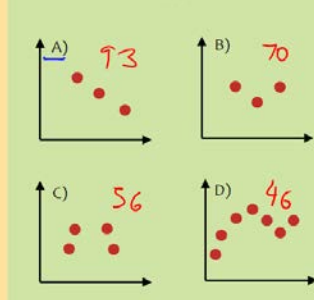
Figure from Glorot & Bentio (2010)

## Neural Network Errors

**Question A:** On which of the datasets below could a one-hidden layer neural network achieve zero classification error? **Select all that apply.**



**Question B:** On which of the datasets below could a one-hidden layer neural network for regression achieve nearly zero MSE? **Select all that apply.**



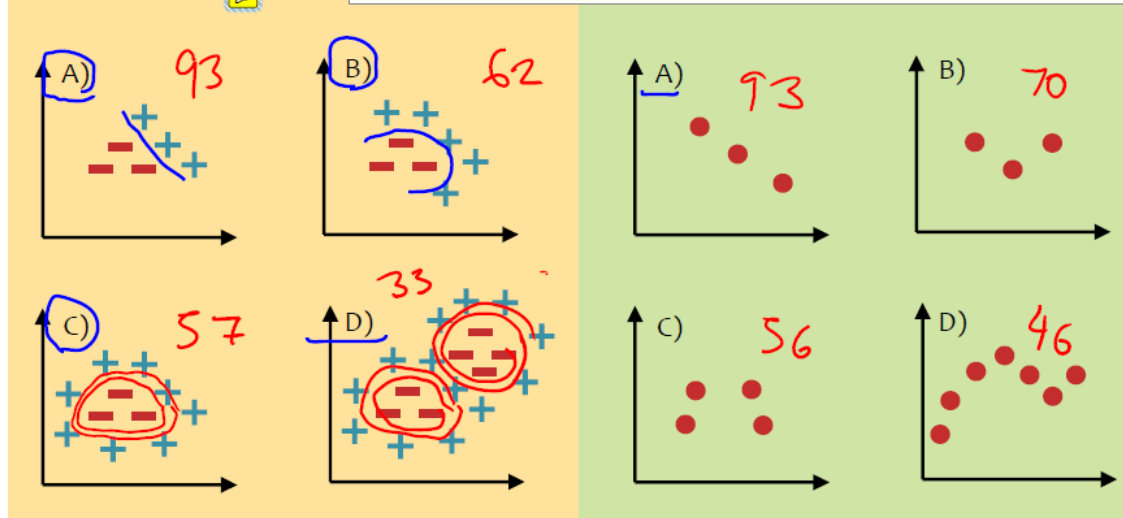
E = calamity

Neu

**Question A:** On which of the ...  
could a one-hidden layer ne ...  
achieve zero classification e ...  
that apply.

newname  
10/9/2021 2:43:08 PM

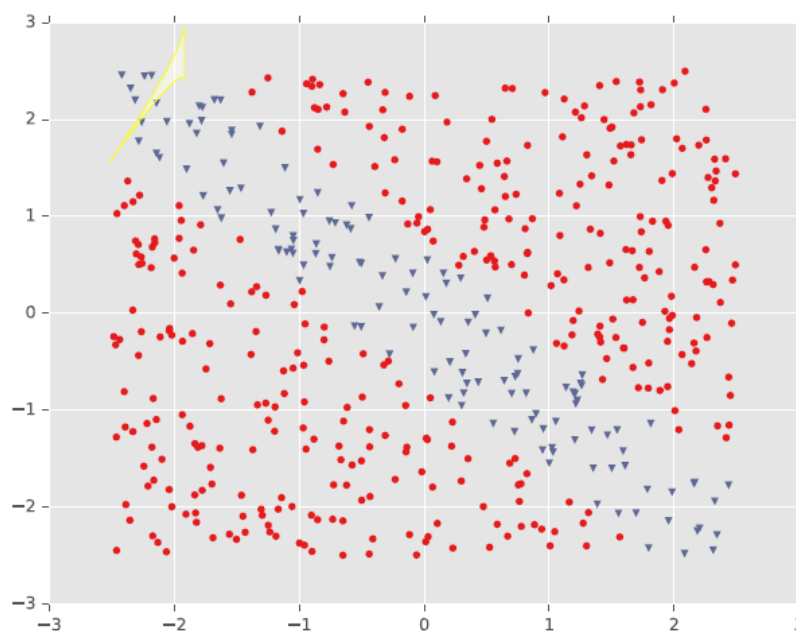
الاجابه هو كل اللي قدامك فيه ... تعال نرجع لل universal function approximator theorem .. تعال نبص علي ال regression ... ال  
Universal function approximator theorem يقولك ان  
"There exists a 1-Hidden layer neural network that can achieve arbitrarily close approximation to any function out there"  
.. approximation ده .. انت هتعرف تجيب ال probability distribution هنا بثبوتية approximate .. اللي هتطلعلك decision boundary اللي  
طب بالنسبة للكلاسيكيشن .. انت بتحاول ت approximate ال probability distribution اللي بتحاول تجيب .. فكر فيها .. كإنك بتحاول تجيب  
بتجيبك صفر ايزو ... لو هتفكر في الفانكشن اللي بتحاول تعملها أبوكسيميشن .. فكر فيها .. كإنك بتحاول تجيب probability distribution  
وانك بتحاول ت approximate ال probability distribution اللي هيطيلك الديسشن باوندرى اللي له صفر ايزو ... فانت هتعرف تجيب  
contour lines بالطريقه ديه



34

طبيب هنتكلم علي ايه اللي بيحصل علي ارض الواقع باستخدام ال neural network .. الدكتور خذ اول مثال فيقول ام عملناه باللوجستيك ريجريشن ..

## Example #1: Diagonal Band

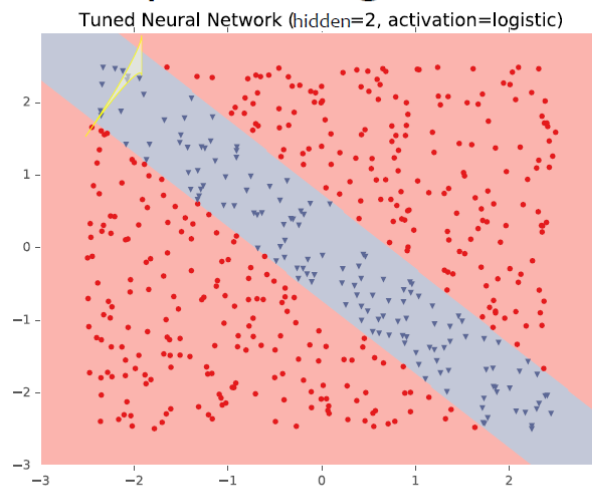


طلع كلو احمر .. زي الصورة اللي تحت ديه ..

## Example #1: Diagonal Band



## Example #1: Diagonal Band

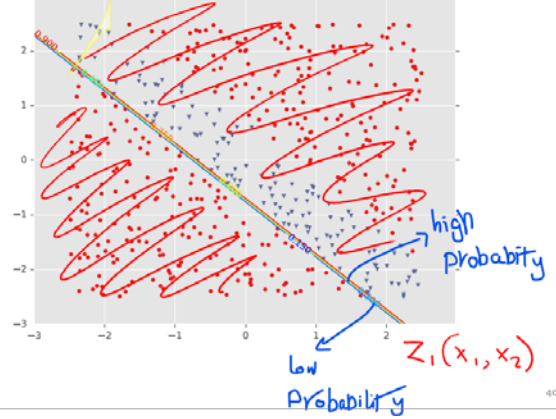


39

لما جينا عملناها بالنويرال نتورك .. الدنيا اشتغلت عادي ...

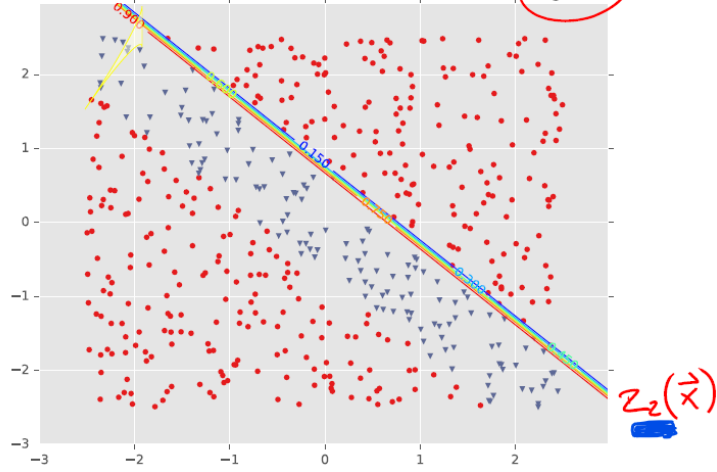
### Example #1: Diagonal Band

LR1 for Tuned Neural Network (hidden=2, activation=logistic)



### Example #1: Diagonal Band

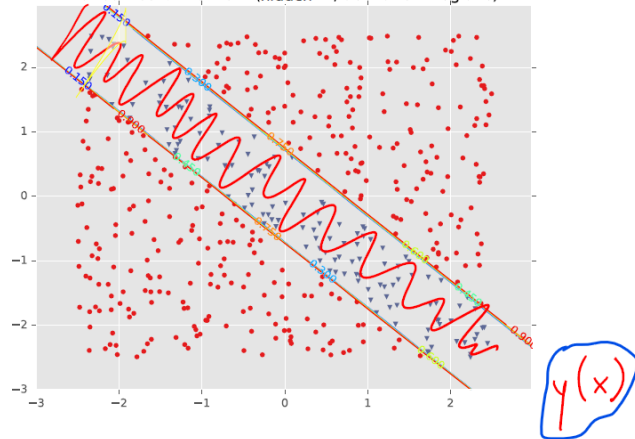
LR2 for Tuned Neural Network (hidden=2, activation=logistic)



41

### Example #1: Diagonal Band

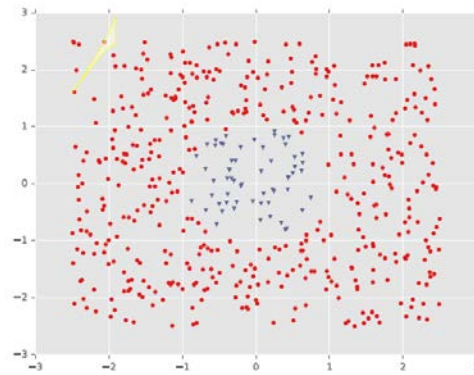
Tuned Neural Network (hidden=2, activation=logistic)



42

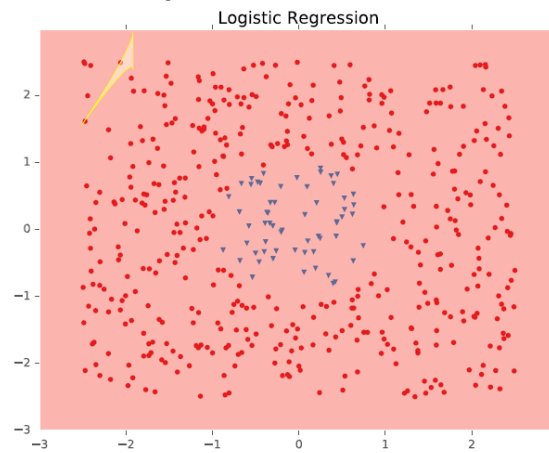
في المثال اللي فوق ده في 3 رسومات عشان الدكتور حبيب يوضح ان أول صوره هي أول layer و ثاني صوره هي ثاني layer و آخر صوره هو الإثنيتين مع بعض ....

### Example #2: One Pocket



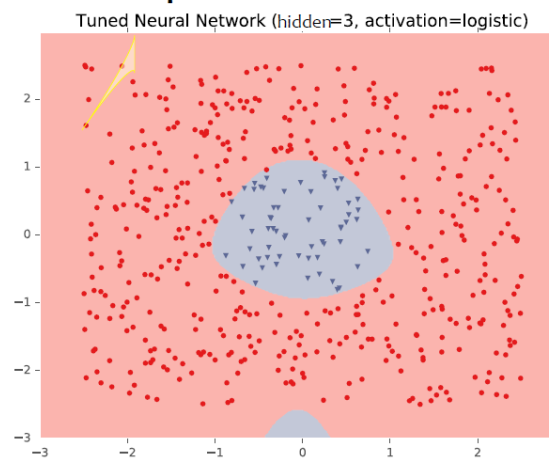
44

### Example #2: One Pocket



45

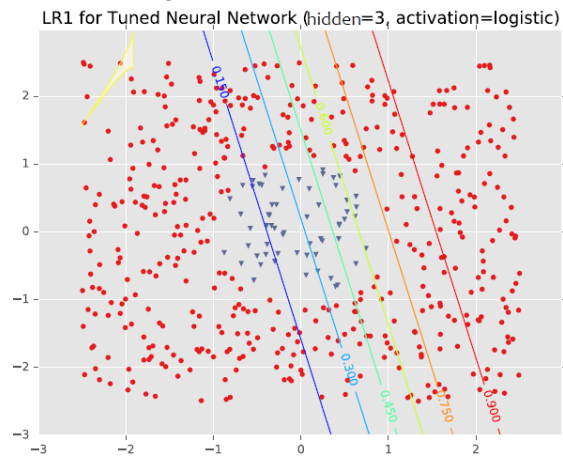
### Example #2: One Pocket



46

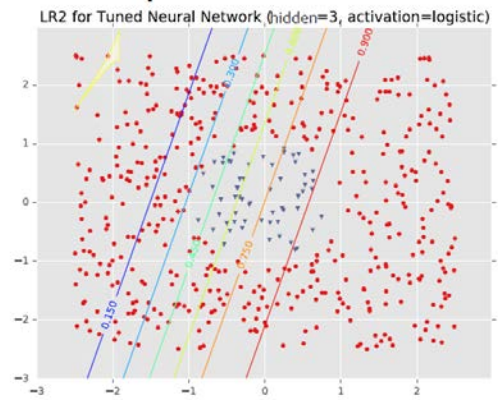


## Example #2: One Pocket



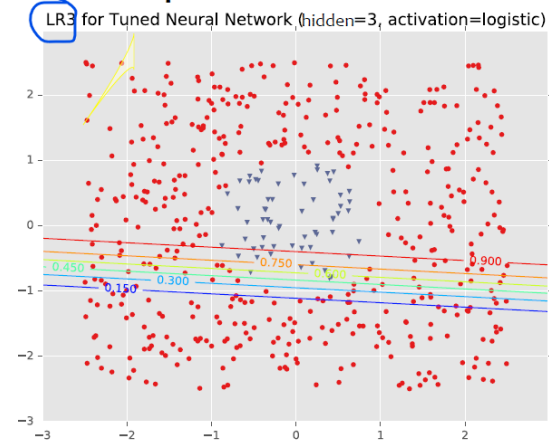
47

## Example #2: One Pocket



48

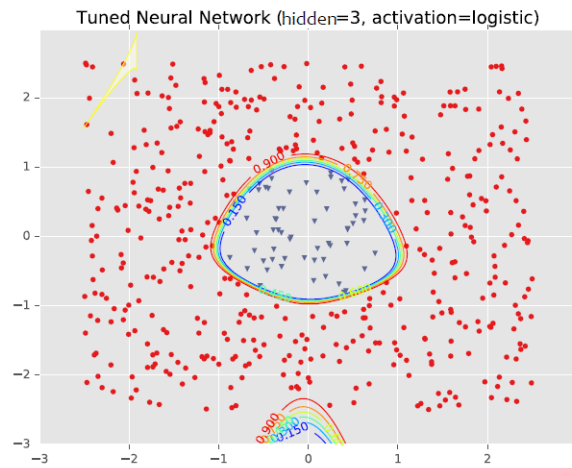
## Example #2: One Pocket



49



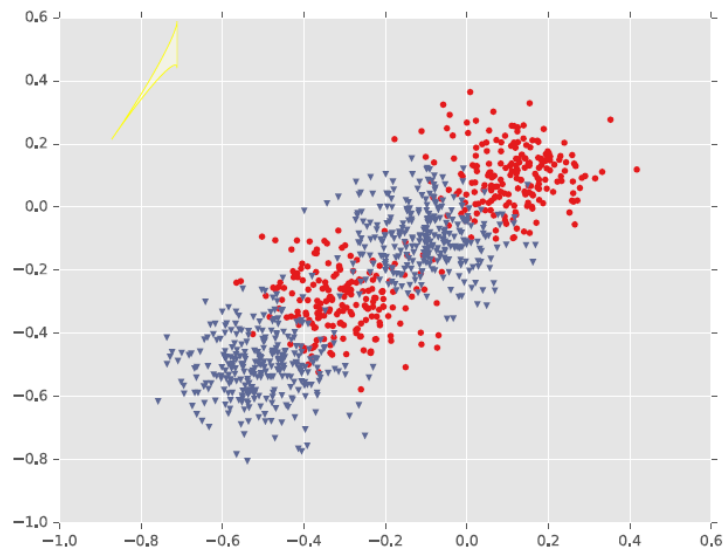
## Example #2: One Pocket



50

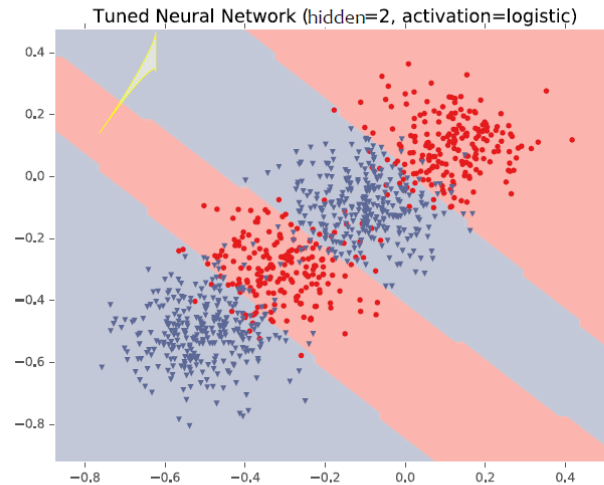
تعال نشوف مثال ثالث .. الدكتور بيقول هنا الهدف من المثال ده انو يشوف ال neural network وهي بت .. fail ..

## Example #3: Four Gaussians



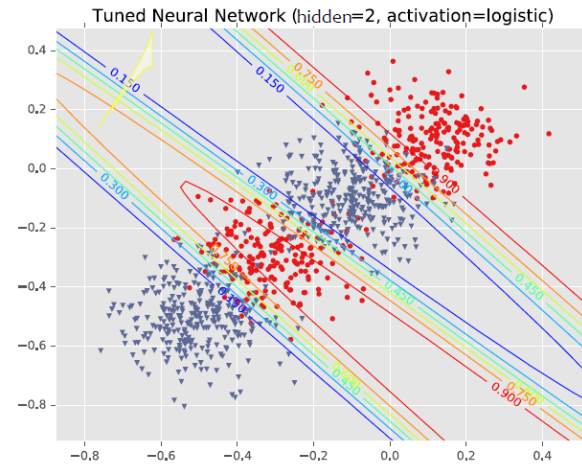
52

## Example #3: Four Gaussians



55

## Example #3: Four Gaussians



58

المهم الدكتور بيتكلم علي المثال الرابع .. هو استخدم hidden units 3 ف مره و ف مره ثانيه استخدم hidden units 4 .. لقي ان ال 4 مبقاش احسن من ال 3 ... ده عشان ال objective function كان non convex والطريقه اللي الدكتور استخدمها عشان يعمل ال training هنا كان انو عمل random initialization .. فبالنالي ال local minima في حالة ال 3 غير حالة ال 4 .....

طيب هنتكلم دلوقت علي ال ازاوي هنعسب ال gradients .. الدكتور حط سوالين .. اول واحد امتي هنعقد نحسب ال gradient بتاع أي نويرال نتورك و ثاني سؤال هو امتي هنعقد نخلي ال gradient computation efficient .. طلع في طرق كتيره .. الدكتور بيقولك عاوزنا نفكر في طرق more general هو ان عندك فانكشن اسمها  $f$  بتاخد فيكتور طولها  $A$  وبترجع فيكتور طولها  $B$  .. فلما يبقا عندك فانكشن زي ديه .. انت بتحتاج تحسب  $A*B$  different partial derivatives .... اللي هو partial derivative of each output with respect to each input .. طيب فال setting بتاع التريننج للنويرال نتورك .. مختلف شويه عن الماشين ليرنج .. عشان اللي عندنا هو بارمترز كتيره عشان ال  $A$  هتبقا كبيره .. بس  $B$  هتبقا بتساوي 1 ... الدكتور سأل سؤال ... ايه هي ال value that is the output that we are trying to optimize ال 1 .. هتبقا ال loss function اللي هو ال output of Loss( $\hat{y}$ ,  $y$ ) .. او اللي هو تقدر تقول انو ال  $J$  اللي هو هيبقا الفانكشن اللي بتاخد الأفرديج بتاع كل ال  $L_s$  .. في ال setting بتاعنا  $B = 1$  ..

# Training

# Approaches to Differentiation

## 1. Finite Difference Method

- Pro: Great for testing implementations of backpropagation
- Con: Slow for high dimensional inputs / outputs
- Required: Ability to call the function  $f(\mathbf{x})$  on any input  $\mathbf{x}$

## 2. Symbolic Differentiation

- Note: The method you learned in high-school
- Note: Used by Mathematica / Wolfram Alpha / Maple
- Pro: Yields easily interpretable derivatives
- Con: Leads to exponential computation time if not carefully implemented
- Required: Mathematical expression that defines  $f(\mathbf{x})$

## 3. Automatic Differentiation - Reverse Mode

- Note: Called Backpropagation when applied to Neural Nets
- Pro: Computes partial derivatives of one output  $f(\mathbf{x})$ , with respect to all inputs  $\mathbf{x}_i$  in time proportional to computation of  $f(\mathbf{x})$
- Con: Slow for high dimensional outputs (e.g. vector-valued functions)  $\rightarrow$  large  $B$
- Required: Algorithm for computing  $f(\mathbf{x})$

## 4. Automatic Differentiation - Forward Mode

- Note: Easy to implement. Uses dual numbers.
- Pro: Computes partial derivatives of all outputs  $f(\mathbf{x})$ , with respect to one input  $\mathbf{x}_i$  in time proportional to computation of  $f(\mathbf{x})$
- Con: Slow for high dimensional inputs (e.g. vector-valued  $\mathbf{x}$ )  $\rightarrow$  large  $A$
- Required: Algorithm for computing  $f(\mathbf{x})$

Given  $f : \mathbb{R}^A \rightarrow \mathbb{R}^B, f(\mathbf{x})$

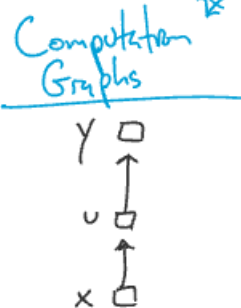
Compute  $\frac{\partial f(\mathbf{x})_i}{\partial x_j} \forall i, j$

## Backprop

Wednesday, February 26, 2020 11:58 AM

### Chain Rule

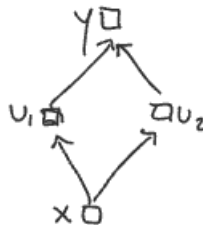
Def #1:  $y = f(u)$   
 $u = g(x)$



### DeFs

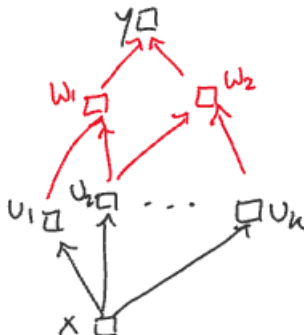
$$\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dx}$$

Def #2:  $y = f(u_1, u_2)$   
 $u_2 = g_2(x)$   
 $u_1 = g_1(x)$



$$\frac{dy}{dx} = \frac{dy}{du_1} \frac{du_1}{dx} + \frac{dy}{du_2} \frac{du_2}{dx}$$

Def #3:  $y = f(\vec{u})$   
 $\vec{u} = g(x)$



$$\frac{dy}{dx} = \sum_{k=1}^K \frac{dy}{du_k} \frac{du_k}{dx}$$

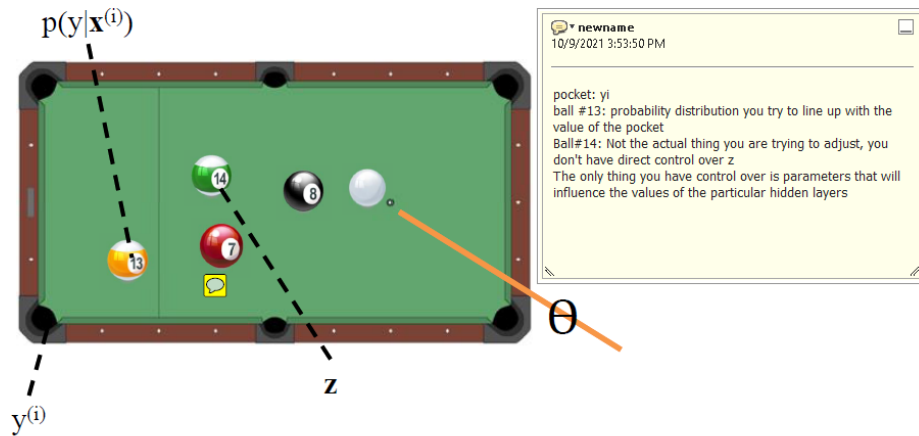
★ Holds for any intermediate quantities  $\vec{u}$

في التعريف رقم 3 .. الدكتور أكد علي حته انك حتي لما هتضيف ال  $w1$  and  $w2$  في النص .. هي مش هتفرق فحاجه لإنك في الآخر بتقول  $y = f(u)$  .. فهو انت متعرفش ال  $f$  اصلا بتبقا عامله ازاى ...

الدكتور رجع علي السلايدز وقال ان الباك بروبجيشن ما هو إلا repeated application of the chain rule فالتناس اللي بتفضل تصيح .. لو سمحت كفايه صياح .. هو بس بببقا شكلو رخم شويه لما بببقا عندك complicated computation graph .. المهم خد بالك برضو ان ال computation graph غير ال Neural Network ..

تعال ناخذ شوية intuition عن اللي بيحصل في الباك بروبجيشن ..

## Error Back-Propagation



Slide from (Stoyanov & Eisner, 2012)

87

بص علي مثال البلياردو اللي في المحاضره ... الدكتور دلوقت بيتكلم علي ازاى ان الباك بروب بيحصل نفس الحاجه عشان يآبدت البارمترز بتاعت النويرال نتورك ...

### Backprop Ex #1

$$y = f(x, z) = \exp(xz) + \frac{xz}{\log(x)} + \frac{\sin(\log(x))}{xz}$$

#### Forward Computation

Given  $x=2, z=3$

$$a = xz$$

$$b = \log(x)$$

$$c = \sin(b)$$

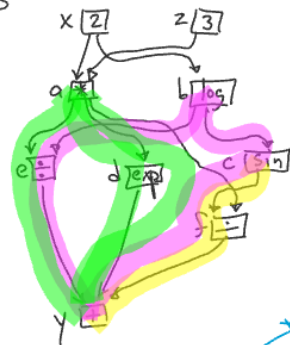
$$d = \exp(a)$$

$$e = a/b$$

$$f = c/a$$

$$y = d + e + f$$

#### Computation Graph



#### Backward Computation

$$g_y = \frac{\partial y}{\partial y} = 1$$

$$g_f = \frac{\partial y}{\partial f} = 1, g_e = \frac{\partial y}{\partial e} = 1, g_d = \frac{\partial y}{\partial d} = 1$$

$$g_c = \frac{\partial y}{\partial c} = \frac{\partial y}{\partial f} \frac{\partial f}{\partial c} = (g_f) \left( \frac{1}{a} \right)$$

$$g_b = \frac{\partial y}{\partial b} = \frac{\partial y}{\partial e} \frac{\partial e}{\partial b} + \frac{\partial y}{\partial f} \frac{\partial f}{\partial b} = (g_e) \left( \frac{-a}{b^2} \right) + (g_f) (\cos(b))$$

$$g_a = \frac{\partial y}{\partial a} = \frac{\partial y}{\partial d} \frac{\partial d}{\partial a} + \frac{\partial y}{\partial e} \frac{\partial e}{\partial a} + \frac{\partial y}{\partial f} \frac{\partial f}{\partial a} = (g_d) \left( \frac{1}{b} \right) + (g_e) (\exp(a)) + (g_f) \left( \frac{-c}{a^2} \right)$$

$$g_x = (g_a)(z) + (g_b) \left( \frac{1}{x} \right)$$

$$g_z = (g_a)(x)$$

#### Updates for Backprop:

$$g_x = \frac{\partial y}{\partial x} = \sum_{k=1}^K \frac{\partial y}{\partial u_k} \frac{\partial u_k}{\partial x}$$

$$\star = \sum_{k=1}^K (g_{u_k}) \left( \frac{\partial u_k}{\partial x} \right)$$

# Neural Network Training

- Consider a 2-hidden layer NN
- Params are  $\theta = [\alpha^{(1)}, \alpha^{(2)}, \beta]$
- SGD Training

Iterate until convergence:

① Sample  $i \in \{1, \dots, N\}$

② Compute gradient by backprop

$$g_{\alpha^{(1)}} = \nabla_{\alpha^{(1)}} J^{(i)}(\theta)$$

$$g_{\alpha^{(2)}} = \nabla_{\alpha^{(2)}} J^{(i)}(\theta)$$

$$g_{\beta} = \nabla_{\beta} J^{(i)}(\theta)$$

③ Step opposite the gradient

$$\alpha^{(1)} \leftarrow \alpha^{(1)} - \eta g_{\alpha^{(1)}}$$

$$\alpha^{(2)} \leftarrow \alpha^{(2)} - \eta g_{\alpha^{(2)}}$$

$$\beta \leftarrow \beta - \eta g_{\beta}$$

$$\theta \leftarrow \theta - \eta g_{\theta}$$

Recall:

$$\nabla J(\vec{a}, \vec{b}) = \nabla_{\vec{a}, \vec{b}} J(\vec{a}, \vec{b})$$

$$\nabla_{\vec{a}} J(\vec{a}, \vec{b}) = \begin{bmatrix} \partial J / \partial a_1 \\ \partial J / \partial a_2 \\ \vdots \\ \partial J / \partial a_n \end{bmatrix}$$

$$J^{(i)}(\theta) = \ell(h_{\theta}(\vec{x}^{(i)}), y^{(i)})$$

left out terms

$$\beta \leftarrow \beta - \delta \tilde{g}_\beta$$

Backprop Ex #2: for NN

Given: ① Dec. Fn.  $\hat{y} = h_\theta(\vec{x}) = \sigma\left(\underbrace{(\alpha^{(3)})^T}_{\tilde{z}^{(2)}} \sigma\left(\underbrace{(\alpha^{(2)})^T}_{\tilde{z}^{(1)}} \sigma\left(\underbrace{(\alpha^{(1)})^T}_{\tilde{z}^{(0)}} \vec{x}\right)\right)\right)$

left out intercept terms

② loss Fn.  $J = \ell(\hat{y}, y^*) = y^* \log(\hat{y}) + (1 - y^*) \log(1 - \hat{y})$  ← cross entropy

①

Forward

Given  $\vec{x}, \alpha^{(1)}, \alpha^{(2)}, \alpha^{(3)}, y^*$

$$\tilde{z}^{(0)} = \vec{x}$$

for  $i = 1, 2, 3$ :

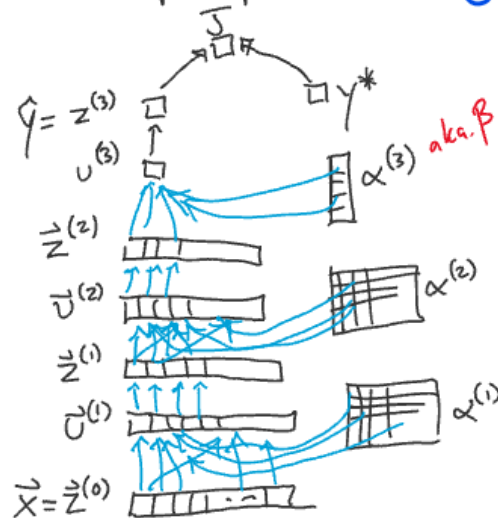
$$\vec{u}^{(i)} = (\alpha^{(i)})^T \tilde{z}^{(i-1)}$$

$$\tilde{z}^{(i)} = \sigma(\vec{u}^{(i)})$$

$$\hat{y} = z^{(3)}$$

$$J = \ell(\hat{y}, y^*)$$

② Comp Graph



③ Backward

$$g_J = [1]$$

$$g_{\hat{y}} = \frac{y^*}{\hat{y}} - \frac{(1 - y^*)}{(1 - \hat{y})}$$

Loss

for  $i = 3, 2, 1$ :

$$\begin{aligned} g_{\vec{u}^{(i)}} &= \dots \\ g_{z^{(i-1)}} &= \dots \\ g_{\alpha^{(i)}} &= \dots \end{aligned}$$

$$g_x = g_{z^{(0)}}$$