

انهارده هنتكلم شويه عن الديب ليرننج و كمان هنفكر في ال PAC learning وده يعتبر أول theoretical justifications للي احنا بنعملو لحد دلوقت ... هنتكلم الاول علي high-level picture علي ال RNN لمدة نص ساعه .. والباقي هيخش في ال PAC .. يادوب .. احنا بدأنا نتكلم علي ال RNN وخذنا مثال علي ال time series dataset

## Dataset for Supervised Part-of-Speech (POS) Tagging

Data:  $\mathcal{D} = \{x^{(n)}, y^{(n)}\}_{n=1}^N$

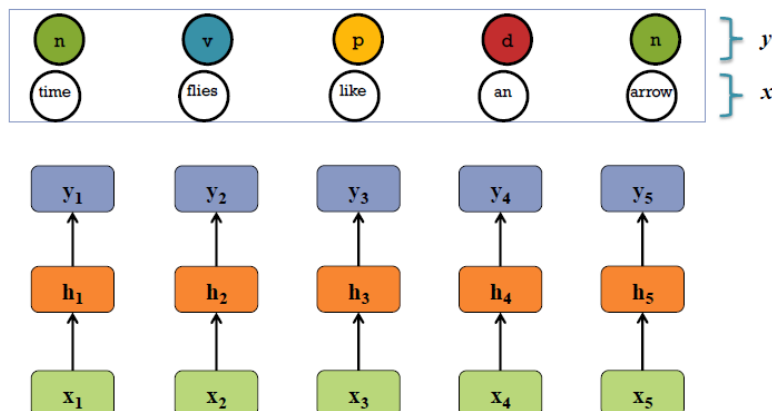
Sample 1:	<div>n</div>	<div>v</div>	<div>p</div>	<div>d</div>	<div>n</div>	$y^{(1)}$
	<div>time</div>	<div>flies</div>	<div>like</div>	<div>an</div>	<div>arrow</div>	$x^{(1)}$
Sample 2:	<div>n</div>	<div>n</div>	<div>v</div>	<div>d</div>	<div>n</div>	$y^{(2)}$
	<div>time</div>	<div>flies</div>	<div>like</div>	<div>an</div>	<div>arrow</div>	$x^{(2)}$
Sample 3:	<div>n</div>	<div>v</div>	<div>p</div>	<div>n</div>	<div>n</div>	$y^{(3)}$
	<div>flies</div>	<div>fly</div>	<div>with</div>	<div>their</div>	<div>wings</div>	$x^{(3)}$
Sample 4:	<div>p</div>	<div>n</div>	<div>n</div>	<div>v</div>	<div>v</div>	$y^{(4)}$
	<div>with</div>	<div>time</div>	<div>you</div>	<div>will</div>	<div>see</div>	$x^{(4)}$

5

فانت ممكن بيقا عندك speech مثلا او hand writtern digits ... وهكذا .. وكان الدكتور قال ان اللي عاوزين نعملو هو اننا نفكر في ازاي نقدر conceptualize باستخدام ال fixed size neural network .. عشان نشغل علي حاجه variable length input and output .. واللي عاوزين نفكر فيه هنا هو اننا نبني حاجه بتشتغل زي اللي كنا بنعملو في ال feature engineering discussion ... فهييqa عندنا مثلا نويرال نتورك هتاخذ كلمه مثلا و بيقا عندك hidden unit و predict a preobability distribution over a part of speech وبعدين نكرر الحوار ده لكل واحد من ال speech tags اللي عندنا .. فهنلاقي عندنا 5 نويرال نتوركس منفصلين عن بعض .. في مشكله في الحوار ده ..

## Time Series Data

**Question 1:** How could we apply the neural networks we've seen so far (which expect **fixed size input/output**) to a prediction task with **variable length input/output**?



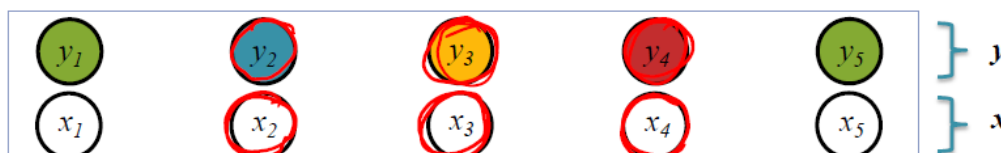
9

المشكلة اللي الدكتور بيقلها هو انك لما هتيجي تبص علي كلمة flies .. ازاي هتعرف هي noun or verb ؟ .. هل احنا بنتكلم علي الدبان .. ولا حاجه بتطير .. الطريقه اللي انت بتعرف بيها هي انه ي واحد فيهم .. انك تبص علي كل كلمه في الجملة اللي عندك .. فانت عارف ان في الحاله ديه ان كلمه flies وراها كلمة time فهي فعلياً verb .. بس الصوره اللي فانت من السلايد مش بتعامل حاجه زي كذا .. انت كل اللي عارفه هو الكلمه اللي عندك دلوقت ... في poll question ... ازاي ن incorporate context ..

newname  
10/14/2021 3:06:32 AM

في ناس كثيره اختارو الاختيار H .. وده معناه هنا ان اللي هيجعل ان احنا هن construct a feature function تبص علي ال 3 كلمات اللي عندنا .. وده لوحده ماسك 3 اختيارات ... هل احنا لازم ندخل في الاعتبار  $y_i$  .. الحقيقة علي حسب ازاي بتعمل ال formulation بتاع ال multi-nomial-logistic regression ... فالدكتور بيقل الطريقه اللي بتعمل بيها formulation ان مفيش مساحه تضيق فيتشرز لل output tag .. بس فعلياً انت تقدر بيضا عندك features of output value بس بتحتاج تعمل reformulation للطريقه اللي بتكتب بيها ال logistic regression .. مشكلة الاختيار رقم H هو انو معاه  $y(i+1)$

Question 2: How could we incorporate context (e.g. words to the left/right, or tags to the left/right) into our solution?



Multiple Choice:

Working left-to-right, use features of...

cs.kml

	$x_{i-1}$	$x_i$	$x_{i+1}$	$y_{i-1}$	$y_i$	$y_{i+1}$
A	✓					
B				✓		
C	✓			✓		
D	✓			✓	✓	✓
E	✓	✓		✓	✓	✓
F	✓	✓	✓	✓		
G	✓	✓	✓	✓	✓	
H	✓	✓	✓	✓	✓	✓

10

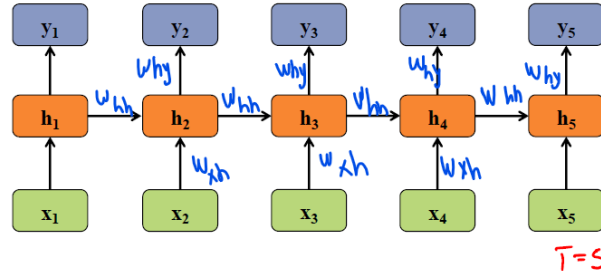
اقرا الكلام اللي مكتوب في السلايد . المهم الدكتور كان بيقل ان اخر اختيار هيقا فيه مشكلة ال  $y(i+1)$  في الحاله ديه المشكله هتبقا ان مثلاً يعني ازاي نبريدكت  $y_3$  احنا هنبص علي  $y_2$  and  $y_4$  .. في التست تايم انت معندكش أي تاج اصلاً .. ال  $y_2$  and  $y_4$  هيجوا منين .. الدكتور بيقل ان اصلاً predicted them .. فازاي عملت prediction for  $y_4$  انت استخدمت  $y_3$  and  $y_5$  ... فالدكتور قالك انت عشان تجيب  $y_4$  .. هتحتاج  $y_3$  .. وعشان تبريدكت  $y_3$  هتحتاج  $y_4$  .. في ال OS ده اسمو deadlock .. هنا اسمو circular dependency في الفيتشر فانكشن بتاعتنا .. بدل كذا اللي نقدر نعملو هو ان نعمل ال tagging problem عن طريق اننا نروح من الشمال لليمين .. نفس طريقه القرايه بتاعت الإنسان .. تعال نبص علي الاختيار بتاع F and G ... انت عمرك ما بتبص قدام هنا .. انت بس بتبص وراك ..

طيب ال RNN هتبقا نفس الفكره ديه .. هو انك ت incorporate contexts of  $x_i$  .. احنا هنفكر ازاي نحط left and right context .. all of the left and right context فتعال نشوف ازاي هنجيب كل الكونتيكست اللي علي الشمال .. الفكره هنبني RNN اول  $y_1$  هتبقا زي standard feat forward NN عادي .. انما تاني  $y_2$  prediction هيبني ال hidden layer  $h_2$  عن طريق انو يبص علي  $h_1$  و  $x_2$  .. و هفضل نتحرك بالطريقه ديه ... ازاي هنبني حاجه زي كذا وتبقا مفيده .. تعال بس نبص علي ال notation الاول ..

# Recurrent Neural Networks (RNNs)

inputs:  $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$   
hidden units:  $\mathbf{h} = (h_1, h_2, \dots, h_T), h_i \in \mathcal{R}^J$   
outputs:  $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$   
 nonlinearity:  $\mathcal{H}$

Definition of the RNN:  
 $h_t = \mathcal{H}(W_{xh}x_t + W_{hh}h_{t-1} + b_h)$   
 $y_t = W_{hy}h_t + b_y$

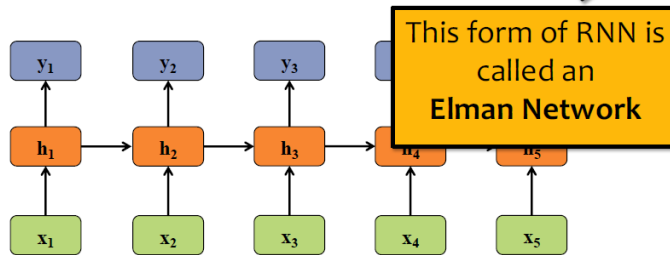


11

# Recurrent Neural Networks (RNNs)

inputs:  $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$   
 hidden units:  $\mathbf{h} = (h_1, h_2, \dots, h_T), h_i \in \mathcal{R}^J$   
 outputs:  $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$   
 nonlinearity:  $\mathcal{H}$

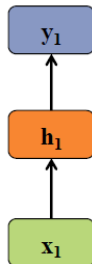
Definition of the RNN:  
 $h_t = \mathcal{H}(W_{xh}x_t + W_{hh}h_{t-1} + b_h)$   
 $y_t = W_{hy}h_t + b_y$



# Recurrent Neural Networks (RNNs)

inputs:  $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$   
 hidden units:  $\mathbf{h} = (h_1, h_2, \dots, h_T), h_i \in \mathcal{R}^J$   
 outputs:  $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$   
 nonlinearity:  $\mathcal{H}$

Definition of the RNN:  
 $h_t = \mathcal{H}(W_{xh}x_t + W_{hh}h_{t-1} + b_h)$   
 $y_t = W_{hy}h_t + b_y$



- If  $T=1$ , then we have a standard feed-forward **neural net with one hidden layer**
- All of the deep nets from last lecture required **fixed size inputs/outputs**

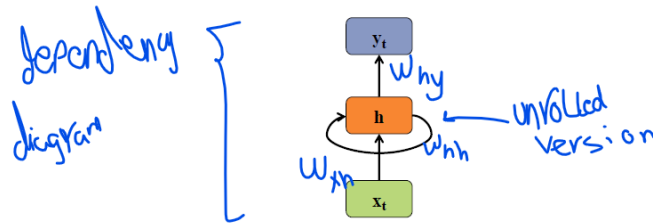
Background	A Recipe for Machine Learning
<p>1. Given training data:</p> <p><math>\{x_i, y_i\}_{i=1}^N</math></p> <p><i>Handwritten: <math>\in \mathbb{R}^T</math> (pointing to <math>x_i</math>), <math>\leftarrow</math> vector (pointing to <math>x_i</math>), <math>\rightarrow</math> bold = vector (pointing to <math>y_i</math>)</i></p>	<p>3. Define goal:</p> $\theta^* = \arg \min_{\theta} \sum_{i=1}^N \ell(f_{\theta}(x_i), y_i)$
<p>2. Choose each of these:</p> <ul style="list-style-type: none"> <li>Decision function</li> <li><math>\hat{y} = f_{\theta}(x_i)</math></li> <li>Loss function</li> <li><math>\ell(\hat{y}, y_i) \in \mathbb{R}</math></li> </ul>	<p>4. Train with SGD:</p> <p>(take small steps opposite the gradient)</p> $\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla \ell(f_{\theta}(x_i), y_i)$ <p><i>Handwritten: + time steps in it (pointing to the summation index i)</i></p>

14

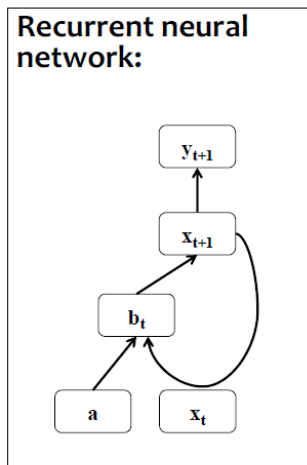
## Recurrent Neural Networks (RNNs)

inputs:  $x = (x_1, x_2, \dots, x_T), x_i \in \mathbb{R}^I$   
hidden units:  $h = (h_1, h_2, \dots, h_T), h_i \in \mathbb{R}^J$   
outputs:  $y = (y_1, y_2, \dots, y_T), y_i \in \mathbb{R}^K$   
nonlinearity:  $\mathcal{H}$

Definition of the RNN:  
 $h_t = \mathcal{H}(W_{xh}x_t + W_{hh}h_{t-1} + b_h)$   
 $y_t = W_{hy}h_t + b_y$

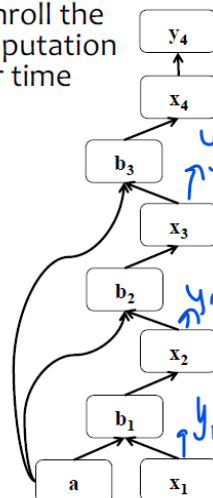


## Background: Backprop through time



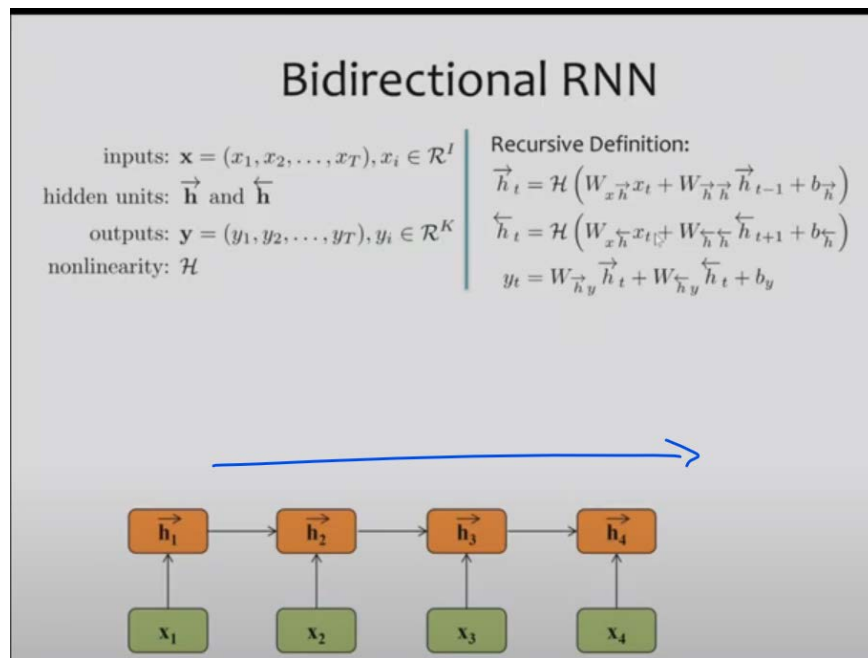
### BPTT:

1. Unroll the computation over time

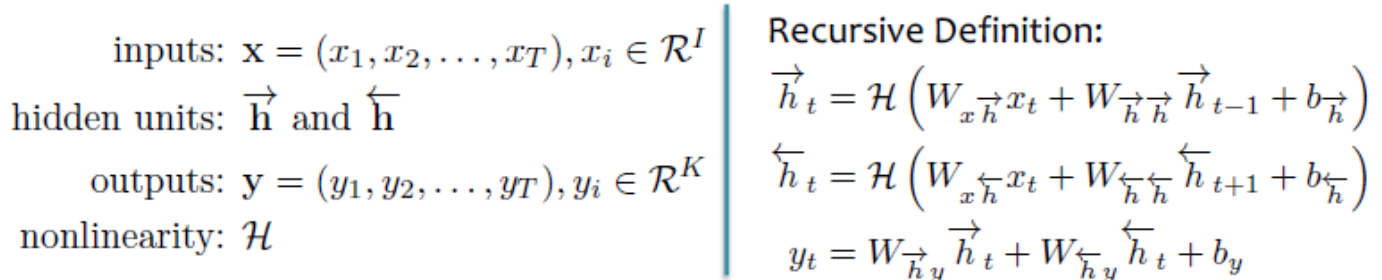


2. Run backprop through the resulting feed-forward network

(Robinson & Fallside, 1987)  
(Werbos, 1988)  
(Mozer, 1995)



## Bidirectional RNN



كدا وانت رايح لقدام انت h3 هو هيددين ريبيرزنتيشن لحد x3 .. وانت راجع بقا h3 شمال ده ريبيرزنتيشن لكل اللي علي اليمين فانت لما بت predict y3 بتأخذ الاتنين h3 left and right arrow وتجييب ال distribution over output class y3 .. طيب لو عوزنا نعمل deep RNNs ... عندنا layers كتيره فوق بعض ..

## Deep RNNs

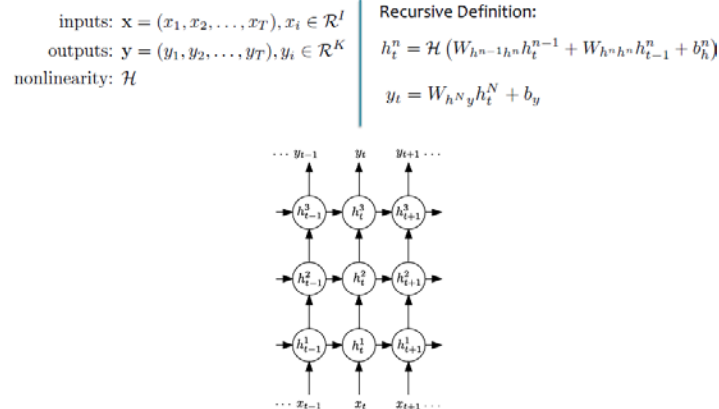


Figure from (Graves et al., 2013)

23

## Deep Bidirectional RNNs

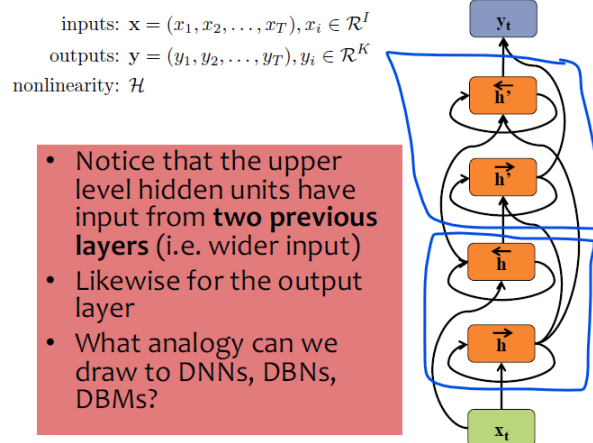


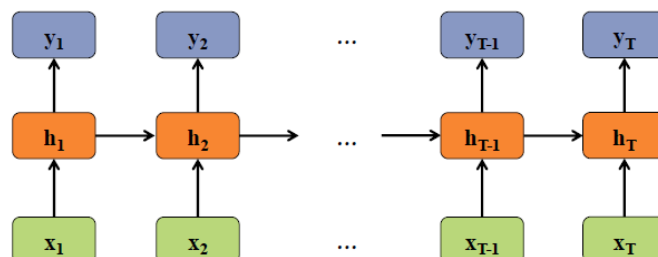
Figure from (Graves et al., 2013)

24

## Long Short-Term Memory (LSTM)

Motivation:

- Standard RNNs have trouble learning long distance dependencies *problem in Elman*
- LSTMs combat this issue



الـ LSTMs ببيعالج الـ long distance dependencies ... افترض ان عندك second-to-last tag وانت محتاج تبقا فاهم اول كلمتين فعلا بشكل كويس .. غير كذا انت مش هتفهم الإنترنتيشن ... فكر في "time flies like an arrow" ... لو انت فهمت غلط الـ verb flies فده هي break بقيت الفهم بتاع الجملة .. في الانجليزي تقدر تجيب dependencies that stretch the string .. يعني مثلا

"time flies, as you know it always has, like an arrow"

فالمشكلة ديه لو جيت تبص ع السلايد .. هتلاقي ان مشكلة الـ vanishing gradient problem عندك .. بص علي الـ influence بتاع اول كلمه .. وبص علي الكلمه اللي قبل الأخير .. هتلاقي ان الـ standard RNN بتوضح في الصوره ان الـ influence بيقل قدام في النتورك بتاعتك .. معناه ان لو انت فهمت الـ tag اللي قبل الأخير غلط .. بيقتا محتاج تعدل الـ weight بتاع أول كلمه .. وده هيبقا صعب انك تعملو back-prop ..

## Long Short-Term Memory (LSTM)

Motivation:

- Vanishing gradient problem for Standard RNNs
- Figure shows sensitivity (darker = more sensitive) to the input at time  $t=1$

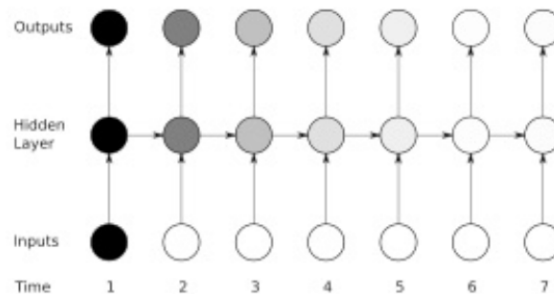


Figure from (Graves, 2012)

26

## Long Short-Term Memory (LSTM)

Motivation:

- LSTM units have a rich internal structure
- The various “gates” determine the propagation of information and can choose to “remember” or “forget” information

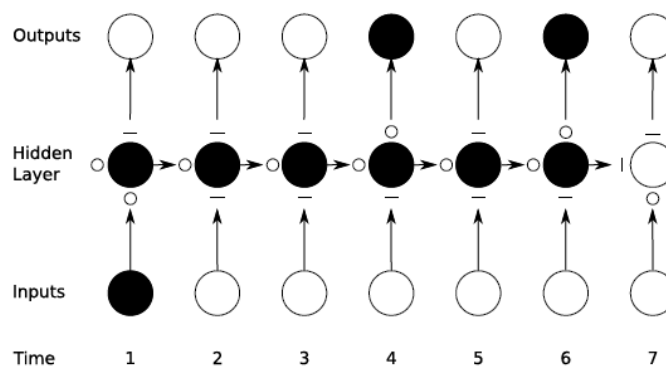


Figure from (Graves, 2012)

27



# Long Short-Term Memory (LSTM)

هنا بستعمل بوابات gates

Motivation:

- LSTM units have a rich internal structure
- The various “gates” determine the propagation of information and can choose to “remember” or “forget” information

open gate

closed gate

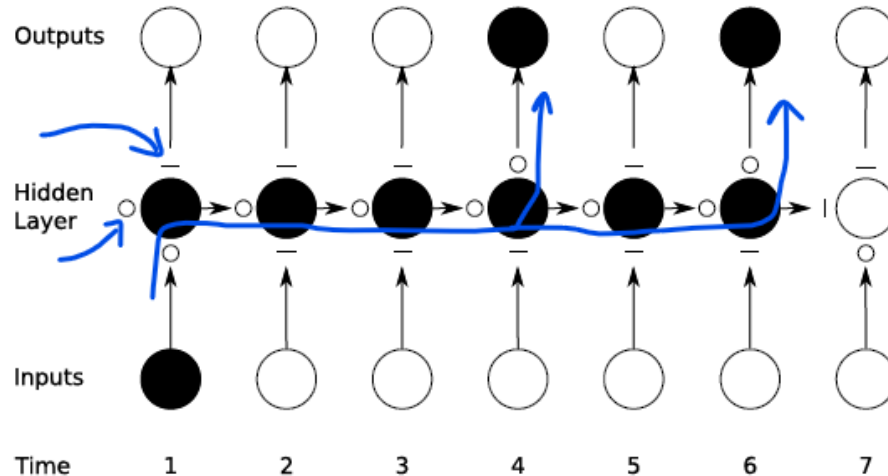
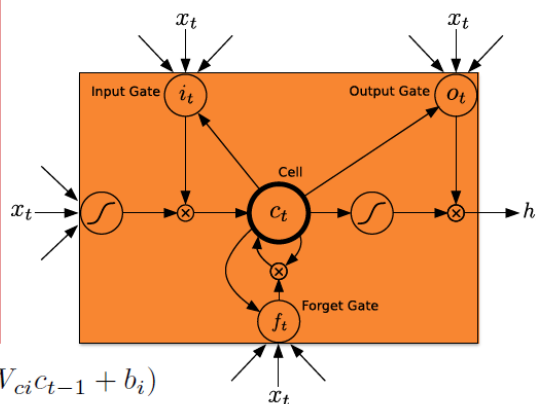


Figure from (Graves, 2012)

27

## Long Short-Term Memory (LSTM)

- Input gate:** masks out the standard RNN inputs
- Forget gate:** masks out the previous cell
- Cell:** stores the input/forget mixture
- Output gate:** masks out the values of the next hidden



$$\begin{aligned}
 i_t &= \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i) \\
 f_t &= \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f) \\
 c_t &= f_t c_{t-1} + i_t \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \\
 o_t &= \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o) \\
 h_t &= o_t \tanh(c_t)
 \end{aligned}$$

Figure from (Graves et al., 2013)

29

طيب دلوقت هنتكلم علي ال learning theory .. الهدف من ده نفهم هو ليه اللي قلناه قبل كذا في المحاضرات فعلا بيشتغل ... "انا بحب الدكتور ده جداً .. بيفكرني بديفيد ميلان بتاع CS50" ... طيب لو عندك كلاسيفير ليه صفر ترينج ابرور .. ايه هو اللي تقدر تقولو علي الجينيريللايزيشن ابرور او الترو



ايروور ... او سؤال ثاني .. لو عندك كلاسيفير ليه تريننج ايروور قليل ... ايه اللي يتقال في الترو ايروور .. منين ما بتتكلم علي ترو ايروور وتست ايروور .. انت دائماً مهتم بالاوفر فيتنج ... وواحد من الحاجات اللي انت بتعملها هو ال regularization المهم الدكتور بيقول هل احنا نقدر نجيب حاجه ك theoretical ... justification for our ways of combating overfitting

## Questions For Today

1. Given a classifier with **zero training error**, what can we say about **true error** (aka. generalization error)?  
(Sample Complexity, Realizable Case)
2. Given a classifier with **low training error**, what can we say about **true error** (aka. generalization error)?  
(Sample Complexity, Agnostic Case)
3. Is there a **theoretical justification for regularization** to avoid overfitting?  
(Structural Risk Minimization)

39

تعال نجواب علي اول سؤال .. عندك داتا سورس .. بيديك unlabeled examples وده بيخش ل expert ... وبعدين يقوم موديه لل learning algorithm .. وبعدين ده يتعلم decision boundary ... تعال نبص علي نوعين الإيروورز الي عندنا ..

## Two Types of Error

1. True Error (aka. **expected risk**)

$$R(h) = \underbrace{P_{\mathbf{x} \sim p^*}(\mathbf{x})}_{\text{Prob}} (c^*(\mathbf{x}) \neq h(\mathbf{x}))$$

This quantity is always **unknown**

2. Train Error (aka. **empirical risk**)

$$\begin{aligned} \hat{R}(h) &= P_{\mathbf{x} \sim S} (c^*(\mathbf{x}) \neq h(\mathbf{x})) \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{1}(c^*(\mathbf{x}^{(i)}) \neq h(\mathbf{x}^{(i)})) \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{1}(y^{(i)} \neq h(\mathbf{x}^{(i)})) \end{aligned}$$

We can **measure** this on the training data

where  $S = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}_{i=1}^N$  is the training data set, and  $\mathbf{x} \sim S$  denotes that  $\mathbf{x}$  is sampled from the empirical distribution.

42

الترو ايروور هنا هو الدكتور قال انو بنحسب ال probability of making a mistake ل example randomly sample from  $p^*$  و ال  $p^*$  هنا هي ال unknown distribution علي الإكسات اللي عندك ... فالتريننج ايروور والتست ايروور هم زي بعض .. الفرق بس ببينقا where the actual distribution over  $X$ s comes from ... هل هي ال unknown  $p^*$  و لا هي ال training sample  $S$  ... لو رجعنا لل PAC/SLT ... model

# PAC / SLT Model

We've also referred to this as the "Function Approximation View"

1. Generate instances from *unknown* distribution  $p^*$

$$\mathbf{x}^{(i)} \sim p^*(\mathbf{x}), \forall i \quad (1)$$

2. Oracle labels each instance with *unknown* function  $c^*$

$$y^{(i)} = c^*(\mathbf{x}^{(i)}), \forall i \quad (2)$$

3. Learning algorithm chooses hypothesis  $h \in \mathcal{H}$  with low(est) training error,  $\hat{R}(h)$

$$\hat{h} = \operatorname{argmin}_h \hat{R}(h) \quad (3)$$

4. Goal: Choose an  $h$  with low generalization error  $R(h)$

## Three Hypotheses of Interest

①

The true function  $c^*$  is the one we are trying to learn and that labeled the training data:

②

$$y^{(i)} = c^*(\mathbf{x}^{(i)}), \forall i \quad (1)$$

The expected risk minimizer has lowest true error:

③

$$h^* = \operatorname{argmin}_{h \in \mathcal{H}} R(h)$$

min True error

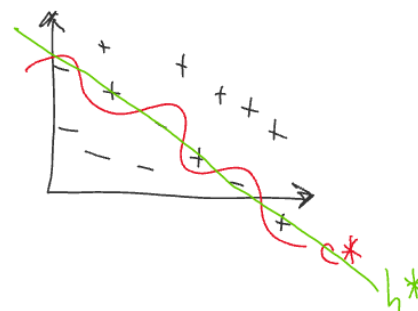
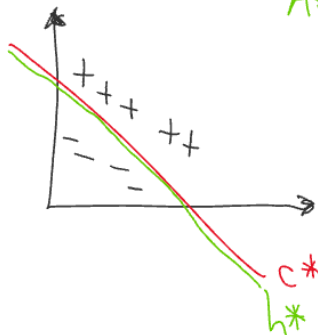
Question: True or False:  $h^*$  and  $c^*$  are always equal.

The empirical risk minimizer has lowest training error:

$$\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} \hat{R}(h) \quad (3)$$

Q: Roll  $h^*$  vs.  $c^*$

Assume  $\mathcal{H}$  = linear D.B.s



الدكتور رجع للسلاید بتاعت الایرور وقال انت متقدرش تحسب ال  $h$  star عشان انت اصلا متعرفش ايه هي ال  $R(h)$  .. فال Pac learning هو الهدف منو انو نجابو علي سؤال هيخلصنا نلف وندور في الكلام شويه .. Can we bound  $R(h)$  in terms of  $\hat{R}(h)$  ... الاجابه اه ...  $R(h)$  ديه انت مش عارفها

انما ال  $\hat{R}(h)$  ده عارفينو .. فاحنا هدفنا ن bound unknown quantity in terms of known quantity .. ال PAC هنا اختصار ل probably approximately correct .. وال PAC learner هو هيقا hypothesis  $h$  from set of Hypothesis  $H$  .. وديه بتبقا approximately correct باحتماليه عاليه ...

## PAC Learning

Q: Can we bound  $R(h)$  in terms of  $\hat{R}(h)$ ?

A: Yes!

PAC stands for Probably Approximately Correct

PAC Learner yields hypothesis  $h \in H$  which is

- approximately correct  $R(h) \approx 0$
- with high probability  $P(R(h) \approx 0) \approx 1$

Def: PAC

ايه المقصود من approximately correct .. اللي احنا مهتمين بيه هو ال true error فبالنسبه لل true error ايه هو معني ان نبقى approximately correct .. معناه ان  $R(h)$  تقريبا بصفر ... ايه المقصود ب high probability .. هنا  $P(R(h))$  close to 1 ... طيب هنتكلم علي تعريف ال PAC criterion وهو ان probability of all hypotheses having true error minus training error absolute value اقل من او بتساوي ايسيلون .. و ده ببيا اكبر من او بيساوي واحد ماينص دلتا .. الدلتا و ايسيلون ارقام صغيره مثلا 0.01 ..... اللي بنقولو ان عشان نبقا تقريبا صح .. لو التریننج ایرور كان صغير اللي هو هنا ال  $\hat{R}(h)$  .. احنا مش عاوزين ال true error يبقا بعيد عنو انم بيقا قريب منو بمقدار ايسيلون .. فال  $R(h)$  هي fixed unknown true error ... وعندك رينج من سالب ايسيلون ل موجب ايسيلون .... التریننج ایرور هيقا في الرينج ده ... احنا ليه دخلنا في سكة ال probability .. هي الاحتماليه هنا بتبقا ف ايفنت .. ايه الايفنت اللي هنا .... فين ال randomness اللي في الموضوع ...

Def: PAC Criterion

$$Pr(\forall h, |R(h) - \hat{R}(h)| \leq \epsilon) \geq 1 - \delta$$



Q: What is random here? A:  $\hat{R}(h)$  is based on a random sample of training data from  $p^*(x)$

estimation of  $R(h)$  randomness

Def: sample

Def: sample complexity is min number of training examples  $N$  st. the PAC criterion is satisfied for  $\epsilon$  and  $\delta$

Def: a hypothesis is consistent w/ the training data if  $\hat{R}(h) = 0$   
 $h \in H$   $D$

### Four Bounds

Two Cases for  $C^*$ :

A) Realizable case:  $C^* \in H$

B) Agnostic case:  $C^* \notin H$  or  $C^* \in H$

Two Cases for  $|H|$ :

A) Finite:  $|H| < +\infty$

B) Infinite:  $|H| = +\infty$

in hypothesis space

مش مشترك

← Agnostic

أول نظرية:

### Thm.1: Sample Complexity (Realizable, Finite $|H|$ )

$N \geq \frac{1}{\epsilon} [\ln(|H|) + \ln(1/\delta)]$  labeled examples  
 are sufficient to ensure that w/prob.  $(1-\delta)$   
 all  $h \in H$  with  $\hat{R}(h) = 0$  have  $R(h) \leq \epsilon$

هنا عندنا  $H$  finite .. فاحنا نفكر ده زي ايه بالظبط .. حاجة كذا زي الديسيجن ترييز ... فانت طلع عدد كثير من الترييز لحد ديبيث معين يعني .. أول نظرية بتقول لو عندنا  $N$  كان اكبر من اللوج بتاع حجم ال hypotheses space و  $N$  ده هو عدد ال labeled examples .. وكان الأمثلة ديه كانت كفيلا انها ت ensure ان باحتماليه عاليه حوالي واحد ماينص دلنا .. كل ال hypothesis in  $H$  ليهم صفر تريينج ايرور اللي هو  $\hat{R} = 0$  ليهم true error  $R(h)$  اقل من او بيساوي  $\epsilon$  .. الجملة ديه قويه .. بتقول ان لو احنا عدد كافي من الأمثلة اللي ليها labels .. فبالتالي باحتماليه  $1 - \delta$  كل ال hypotheses اللي ليهم صفر تريينج ايرور برضو هيبقا ليهم low true error ... ازاى نجيب الاثبات بتاع theorem 1 ...

تعال نفترض عندنا  $K$  bad hypotheses  $h_1, h_2, \dots, h_K$  .. معني انو وحش ان التروو ايرور اكبر من ال  $\epsilon$  .. افترض انك اخترت one specific bad  $h_i$  ... انا عاوز اعرف ايه احتمالية ان  $h_i$  هتبقا consistent مع أول training example .. ايه هو الرقم ده .. ده معناه يعني انها تبقا

صح في التريننج اكزامل اللي عندك دلوقت يعني .. ايه هي احتمالية انك تعمل ايرور في أول مثال عندك بالكلاسيفير السيء ده ... ؟ ايه هي ال  $R(h_i)$  ...  
 $\text{true error} = \text{probability of making a mistake on a randomly chosen example}$

فبالتالي ايه هي الاحتمالية ان  $a \text{ bad } h(i) \text{ will make a mistake on a randomly chosen example?}$  .. هتبقا أكبر من إبسيلون .. ديه احتماليه انك تعمل غلطة ... طيب ايه هي الاحتماليه انك تبقا consistent مع أول مثال عندك يعني انك فعلا عرفت انو صح .. هتبقا اقل من او ببساوي واحد ماينص ابسيلون

طيب ايه احتمالية ان ال  $h_i$  هتبقا consistent مع أول  $N$  training examples ... هتبقا اقل من او ببساوي واحد ماينص ابسيلون الكل مرفوع لل  $N$  ...

طيب ايه احتمالية ان علي الاقل واحده من ال  $h_i$  تبقا consistent مع اول  $N$  training examples .. الدكتور راح يتكلم علي ال union bound اللي هو  $P(A \cup B)$  هي اقل من  $P(A) + P(B)$  .. وقال ان  $A$  هو ان عندك  $B$  .. first of the bad hypotheses be consistent with the data .. فهنا انت بتقول ان احتمالية ان واحد علي الاقل يبقا وحش .. هي ال union .. فهنا انت هتجمع عدد  $k$  ..

### Proof of Thm 1

① Assume  $k$  bad hypotheses  $h_1, h_2, \dots, h_k$  with  $R(h_i) > \epsilon$

② Pick bad  $h_i$ : Prob  $h_i$  is consistent w/first training ex-ple  $\leq 1 - \epsilon$

Prob  $h_i$  is consistent w/first  $N$  training ex-ples  $\leq (1 - \epsilon)^N = \overbrace{(1 - \epsilon)(1 - \epsilon) \dots (1 - \epsilon)}^N$

③ Prob that at least one bad  $h_i$  is consistent w/first  $N$  training ex-ples  $\leq k(1 - \epsilon)^N$

$\Rightarrow \exists h$  s.t.  $\hat{R}(h) = 0$  and  $R(h) > \epsilon$

by Union Bound  
 $P(A \cup B) \leq P(A) + P(B)$

which is bigger?  
 $\leq |H| (1 - \epsilon)^N$

④ Fact:  $(1 - x) \leq \exp(-x) \Rightarrow |H| (1 - \epsilon)^N \leq |H| \exp(-\epsilon N)$

⑤ Calculate value of  $N$  and  $\delta$  s.t.  $|H| \exp(-\epsilon N) \leq \delta$

⑥ Solve for  $N$ :  $\rightarrow \frac{\ln(|H|/\delta)}{\epsilon}$

⑥ Solve for  $N$ :  $\Rightarrow |H| (1/\delta) \leq \frac{1}{\exp(-\epsilon N)}$

Statement 1  $\Rightarrow \ln(|H|) + \ln(1/\delta) \leq \ln(1) - \ln(\exp(-\epsilon N))$

$\Rightarrow \frac{1}{\epsilon} [\ln(|H|) + \ln(1/\delta)] \leq N$   $\epsilon N$

Assume Stat. 1, then:

with  $\swarrow$  w/prob  $\delta$   
probability  $\exists h$  s.t.  $R(h) > \epsilon$  and  $\hat{R}(h) = 0$

"bad"

w/prob  $(1-\delta)$

all  $h \in H$  with  $\boxed{R(h) > \epsilon}$  have  $\boxed{\hat{R}(h) > 0}$

$\Rightarrow$  all  $h \in H$  with  $\boxed{\hat{R}(h) = 0}$  have  $\boxed{R(h) \leq \epsilon}$

$\nearrow$  bad  
not bad

Contrapositive  
 $A \rightarrow B$   
 $\neg B \rightarrow \neg A$

if learner returns  $h \in H$  w/zero training error  
then w/prob  $(1-\delta)$   $h$  has  $\leq \epsilon$  true error