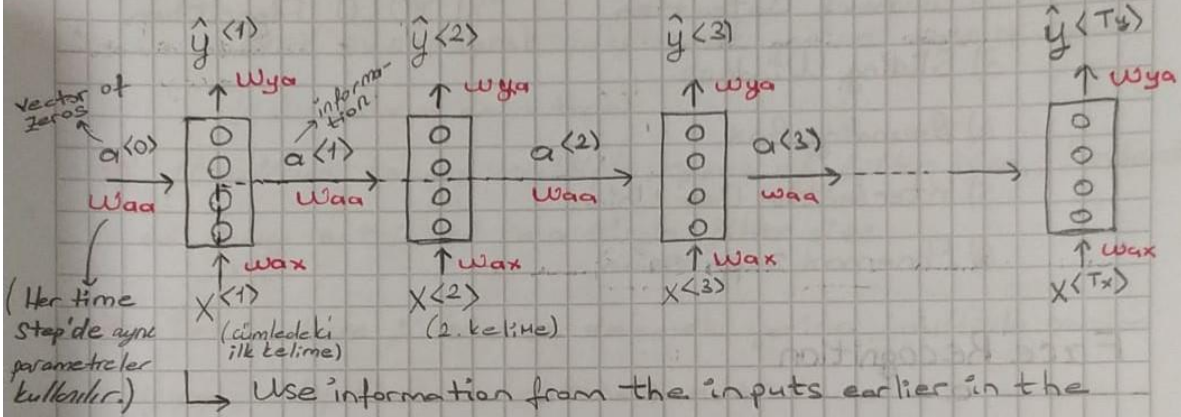


12) Recurrent Neural Networks

Recurrent Neural Network Model:

Simple Unidirectional

Neural Network Architecture :



Use information from the inputs earlier in the sequence; but not information later in the sequence.

$\hat{y}^{(3)}$ prediction'u yaparken sadece $x^{(3)}$ 'den değil; $x^{(1)}$ ve $x^{(2)}$ 'deki information'lardan da yararlanır.

Unidirectional modelde, $y^{(3)}$ 'ü tahmin ederken önceki information'ları kullanır; ancak x_4, x_5, x_6 info'larını kullanamaz.

Forward Propagation :

$$(a^{(1)} = h^{(1)})$$

$$a^{(0)} = \vec{0}$$

$$a^{(t)} = g_1(W_{aa} a^{(t-1)} + W_{ax} x^{(t)} + b_a)$$

$$\hat{y}^{(t)} = g_2(W_{ya} a^{(t)} + b_y)$$

$$a^{(t)} = g(W_{aa} a^{(t-1)} + W_{ax} x^{(t)} + b_a)$$

$$\hat{y}^{(t)} = g(W_{ya} a^{(t)} + b_y)$$

Backpropagation :

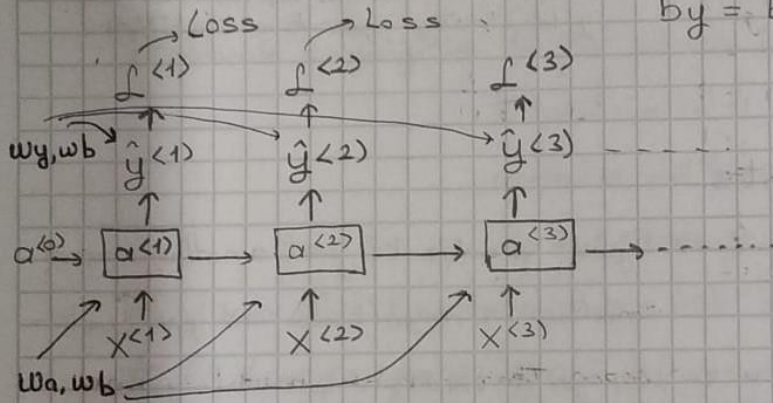
$$W_a = W_a - \frac{\alpha * \partial \text{Loss}(\hat{y}, y)}{\partial W_a}$$

$$W_x = W_x - \frac{\alpha * \partial \text{Loss}(\hat{y}, y)}{\partial W_x}$$

$$W_y = W_y - \frac{\alpha * \partial \text{Loss}(\hat{y}, y)}{\partial W_y}$$

$$b_a = b_a - \frac{\alpha * \partial \text{Loss}(\hat{y}, y)}{\partial b_a}$$

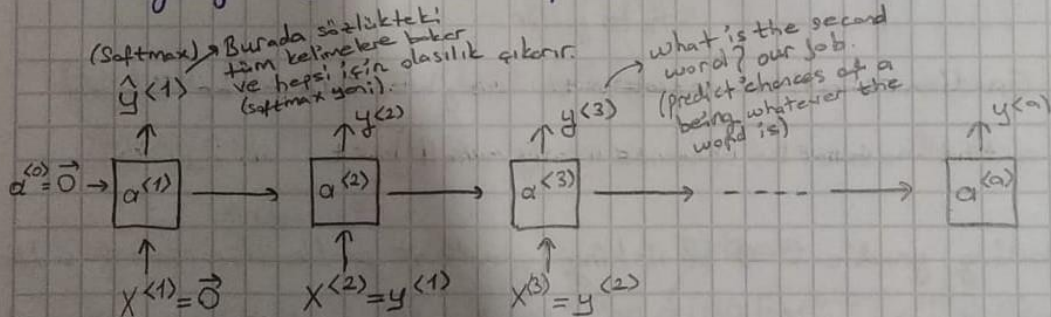
$$b_y = b_y - \frac{\alpha * \partial \text{Loss}(\hat{y}, y)}{\partial b_y}$$



$$L^{(t)}(\hat{y}^{(t)}, y^{(t)}) = -y^{(t)} \log \hat{y}^{(t)} - (1 - y^{(t)}) \log (1 - \hat{y}^{(t)})$$

$$L(\hat{y}, y) = \sum_{t=1}^{T_y} L^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

Language Model and Sequence Generation



we will give correct words. (Her adımda önceki adımdaki doğru sözcükleri de almış oluruz)



$$L(\hat{y}^{(t)}, y^{(t)}) = -\sum_i y_i^{(t)} \log \hat{y}_i^{(t)}$$

$$\text{Overall } L = \sum_t L^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

ör: $P(y^{(1)}, y^{(2)}, y^{(3)}) \rightarrow$ 3 sözcüklü cümle.

$$= P(y^{(1)}) \cdot P(y^{(2)} | y^{(1)}) \cdot P(y^{(3)} | y^{(1)}, y^{(2)})$$

\rightarrow Probability of that 3-word sentence

Sequence Generation

↳ Character-level language model veya word-level language model ile yapılabilir.

News

↳ Modeli haber-ler ile train edersek; "President says..." gibi cümle generate etmesi çok muhtemel.

Shakespeare

↳ Shakespeare ile train edersek "The mortal man..." gibi bir şey generate edilebilir.

Vanishing Gradients

Vanishing gradient \Rightarrow Your gradients get smaller and smaller in magnitude as you back-propagate through lower layers.

Her katmandan geçerken biraz daha düşecek olan bu gradyan değerleri, baştaki katmanlara doğru sıfıra yaklaşıncaya başlar, ve yapay sinir ağı öğrenemez hale gelir.

(Gradients with respect to weights in earlier layers of the network becomes really small)

Exploding gradient \Rightarrow Vanishing gradient probleminin tam tersidir.

(If we multiply a bunch of terms that are all greater than one, we gonna get something greater)



Backprop ile birlikte ağırlıkların yok olmasından dolayı RNN'de vanishing gradient problemi ile karşılaşmak çok olasıdır. Her bir katmandaki ağırlıklar zincir kuralı üzerinden ayarlandığından, gradient değerleri geriye doğru ilerledikçe katlanarak küçülecek ve sonunda

yok olacaktır. Örneğin; MLP uygulamasını ele alırsak, Fıktya göre "What" ve "time" 'dan gelen bilgilerin tamamen ortadan kalktığını düşündüğümüzde "is" ve "it" tahminini nasıl yapabiliriz?

↳ Problemin kaynağı \Rightarrow 10.000 zaman adımı büyüklüğünde bir dizi verisini işleyen bir RNN, optimize edilmesi zor olan 10.000 derin katmana sahiptir.

* Uzun dizi boyutuna sahip RNN'lerde ortaya çıkar.

!!! * Bu durumda, daha önceki katmanlar herhangi bir öğrenme yapamazlar. Bu da ilk Layer'lardaki bilgileri kaybetmelerine neden olur. (Basic RNN, long-range dependency'leri yakalayamaz)

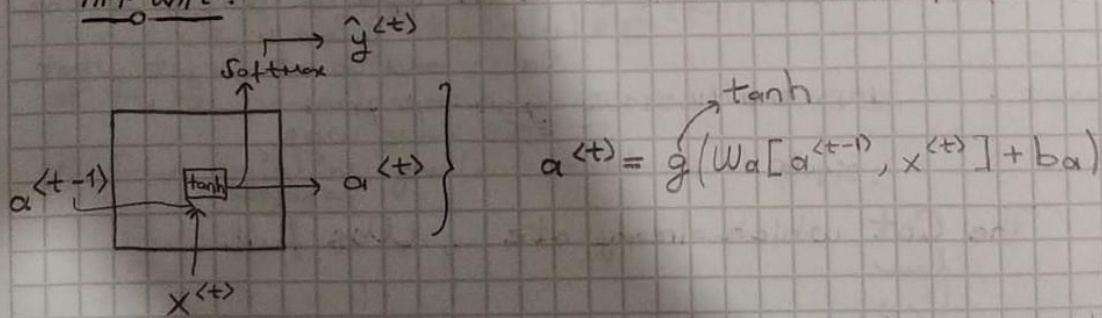
Vanishing (Kaybolan) Gradient problemi

• GRU

• LSTM yöntemleri ile çözebiliriz.

Gated Recurrent Unit (GRU)

RNN Unit:



GRU unit'i diğer sayfada gösterelim.

GRU (simplified):

* c : memory cell (provide a bit of memory to remember)

$$c^{(t)} = a^{(t)}$$

$$\tilde{c}^{(t)} = \tanh(W_c[c^{(t-1)}, x^{(t)}] + b_c)$$

Γ_u = Gate \Rightarrow u stands for update gate.
(gamma u)

$$\Gamma_u = \sigma(W_u[c^{(t-1)}, x^{(t)}] + b_u)$$

\rightarrow always between zero and one.

$\tilde{c}^{(t)} \rightarrow$ Candidate for replacing c of it.

* Gate decides whether or not we actually update the c value using \tilde{c} .

* Gate also decides when do you update these c values.

Specific equation used in GRU:

$$1) c^{(t)} = \Gamma_u * \tilde{c}^{(t)} + (1 - \Gamma_u) * c^{(t-1)}$$

actual c value

a) * when; $\Gamma_u = 0$ Do not update. Hang on to the value and do not forget what this value was.

The Cat which already ate ..., Was full.

Let's say: $c^{(t)} = 1$

GRU unit memorize

the value of c all the way until Was.

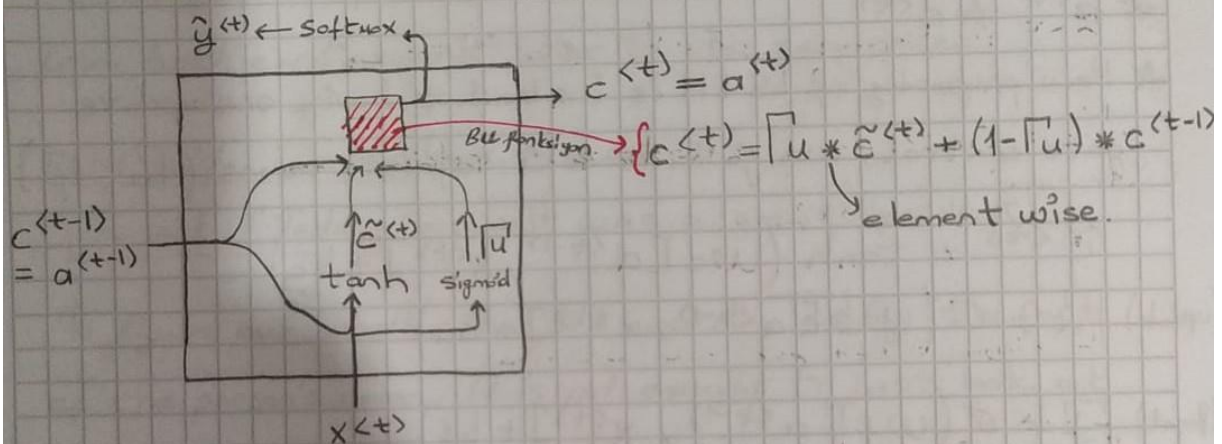
$\Gamma_u = 1$

$\Gamma_u = 0$

$\Gamma_u = 0$

decide when you update c value.

b) * when; $\Gamma_u = 1$; update the value of c with new candidate value.



* Γ_u allow neural network to go on even very long range dependencies. (Cat-----Was,) even they are separated by a lot of words in the middle.

Full GRU } önceki fonksiyonları çok ufak birlektirme.

$$\tilde{c}^{(t)} = \tanh(W_c[\Gamma_r * c^{(t-1)}, x^{(t)}] + b_c)$$

$$\Gamma_r = \sigma(W_r[c^{(t-1)}, x^{(t)}] + b_r)$$

$$\Gamma_u = \sigma(W_u[c^{(t-1)}, x^{(t)}] + b_u)$$

$$c^{(t)} = \Gamma_u * \tilde{c}^{(t)} + (1 - \Gamma_u) * c^{(t-1)}$$

* Try to have longer range connections and

* solve vanishing gradient problems.

Long Short Term Memory (LSTM)

Slightly more powerful and more general version of the GRU.

* $a^{(t)} \neq c^{(t)}$

* New property of LSTM is instead of having one update gate control, we gonna have 2 separate terms. (Γ_u, Γ_f)

$$\tilde{c}^{(t)} = \tanh(W_c[a^{(t-1)}, x^{(t)}] + b_c)$$

(update) $\Gamma_u = \sigma(W_u[a^{(t-1)}, x^{(t)}] + b_u)$

(forget) $\Gamma_f = \sigma(W_f[a^{(t-1)}, x^{(t)}] + b_f)$

(output) $\Gamma_o = \sigma(W_o[a^{(t-1)}, x^{(t)}] + b_o)$

$$c^{(t)} = \Gamma_u * \tilde{c}^{(t)} + \Gamma_f * c^{(t-1)}$$

↳ forget gate.

$$a^{(t)} = \Gamma_o * \tanh(c^{(t)})$$

Fig 1:

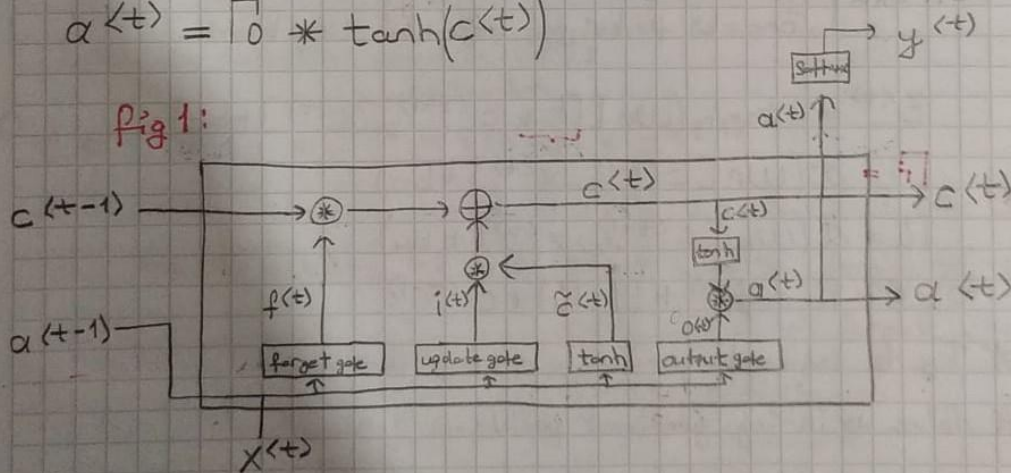
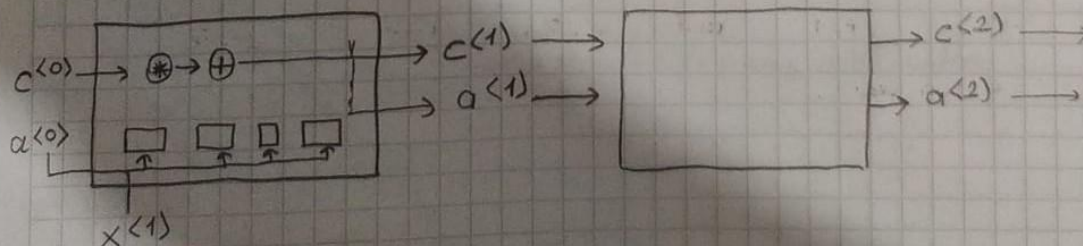


Fig 2:



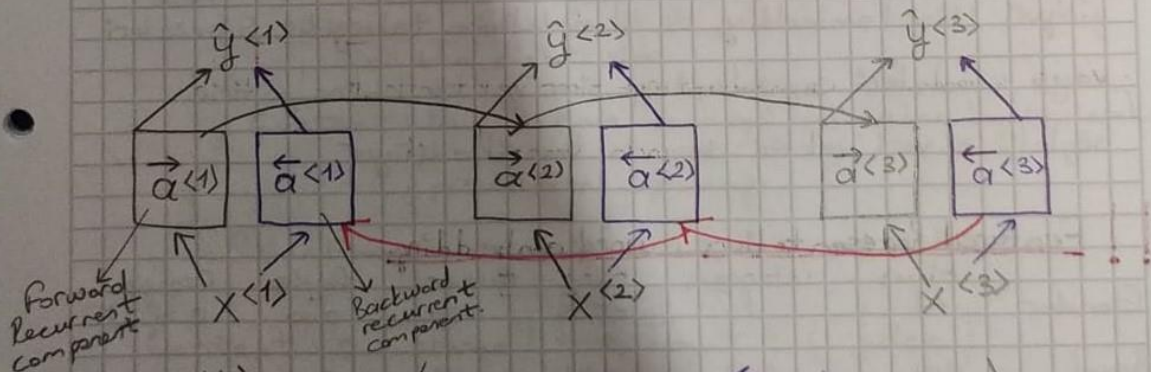
Bidirectional RNN (BRNN)

Lets you at a point in time to take information from both earlier and later in the sequence.

He said, "Teddy bears are on sale!"

He said, "Teddy Roosevelt was a great President!"

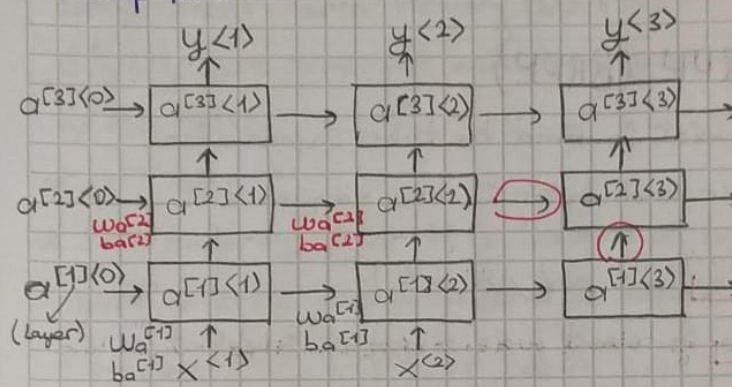
To figure out whether the 3rd word Teddy is a part of the person's name, it is not enough to just look at the first part of the sentence shown as red colour. So to tell, if $\hat{y}^{(3)} = \text{Teddy}$ should be zero or one, we need more information than just the 3 words, because the first 3 words does not tell you whether they are talking about Teddy bears or former US president, Teddy Roosevelt.



$$\hat{y}^{(t)} = g(w_y [\vec{a}^{(t)}, \overleftarrow{a}^{(t)}] + b_y)$$

Information from past, present and future is taken.

Deep RNN:



* Her layer kendi parametrelerine sahiptir ($w_a^{[l]}$, $b_a^{[l]}$)

$$a^{[2][3]} = g(w_a^{[2]}[a^{[2][2]}, a^{[1][3]}] + b_a^{[2]})$$

Natural Language Processing and Word Embeddings

Word Embedding: Metin verisini makinenin anlayacağı hale getirmemiz gerekiyor. Burada word embedding devreye giriyor. Word embedding bir dil modelleme tekniğidir. Sözcükleri veya cümleleri sayısalastırıp birer vektör haline getirir böylelikle veri vektör uzayında temsil ediliyor.

!! Featured Representation: word embedding

	Man	Women	King	Queen	Apple	Orange
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
Size						

Birçok feature all the same, apple ve orange için

"Man" kelimesini represent etmek için

300 dimensional vector kullanılır.

$g(\text{Feature size})$

Apple ve orange'in feature'larının benzer olması, öğrenme algoritmasının işini kolaylaştırır. Örneğin;

* I want a glass of orange juice
* I want a glass of apple juice] Generalize.

! Transfer Learning Word embedding'de sıkça kullanılan bir metot'dur.

! 2D dimension'da visualize etmek için (word'leri) ve gruplamak için t-sne algoritması kullanılır.

↳ it takes 300-D data and it maps it
in a very non-linear way to a 2D space.

Learning Word Embeddings:

1) Word2Vec: Taking as input one word like "orange" etc.
Algorithm and then trying to predict some words
(Skip-Gram) skipping a few words from the left or the right side.

↳ Context: Kaç kelime ve hangi kelimeler kullanılarak prediction yapılacağı.

↳ Skip-Gram modelde, just one word kullanılır.

Ör Context c ("orange") → Target t ("juice")
Target t ("my")

Take ^{context word} input one word

↳ predict what comes little bit before or after the context word.

Not ⇒ 2 farklı Word2Vec yaklaşımı vardır, biz Skip-Gram'ı işledik, (diğeri de CBOW'dur)

Skip-Grams Akış:

- 1) Pick a word to be the context word.
- 2) Pick another word within some window (örneğin ± 5 context word) as a target word.

↳ Bu algoritma Softmax Classifier kullanıyor
(Bu yüzden computationally expensive)

↳ Bu problemi çözmek için "Negative Sampling" metodu kullanılır ve binary classification yolu ile problem çözülmür.

2) GloVe Algorithm:

↳ Skip-Grams kadar yoğun kullanılmaz.
(Word2Vec'e eklenti olarak geliştirildi)

!!! Converting text data into numerical data, which is called **Vectorization** or in the NLP word, it is called **word embeddings**.