

Deep Learning

Seminario Minería de Datos

Álvaro Riascos
Mónica Ribero

2 de mayo de 2017

Contenido

- 1 Introducción
- 2 Modelo Logit
- 3 Redes Neuronales
- 4 Aplicaciones
- 5 Redes en el contexto de NLP
- 6 Otras arquitecturas
 - Convolutional Networks
 - Recurrent Networks

Introducción

Problema

- Sea $(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n)$ una muestra marcada con variables independientes $\vec{x}_i = (x_{i1}, x_{i2}, \dots, x_{im})$ y dependiente $y_i, i = 1, \dots, n$.
- Queremos predecir el valor de y_i dado x_i .
- Asumimos que existe una función tal que $f(x_i) = y_i$ y queremos encontrar su mejor aproximación

Problema – Regresión

$\vec{x}_i = (\text{salario, cantidad de personas con las que vive})$

$y_i = \text{precio de su vivienda}$

Problema – Clasificación

\vec{x}_i = información pixeles en imagen de resonancia

y_i = tiene tumor

Problema – Clasificación

$\vec{x}_i =$ (cuentas de palabras en un documento)

$y_i =$ sentimiento (positivo, neutro, negativo)

Problema – Clasificación

$\vec{x}_i =$ (calificaciones de un usuario a varios restaurantes)

$y_i =$ perfil de cliente

Modelos

- Regresión lineal
- Regresión logística
- Árboles de decisión
- Random Forests y Boosting de árboles

Modelo Logit

Regresión Logística – Clasificación

- En vez de representar el valor de y , representar la Probabilidad de que y pertenezca a cierta categoría dado x

$$\begin{aligned}Pr(y = 1|x) &= g(\theta^t x) \\ &= \frac{1}{1 + e^{-\theta^t x}}\end{aligned}$$

- $\hat{y}_i = h_{\theta}(x_i) = \begin{cases} 1 & g(\theta^t x_i) \geq 0,5 \\ 0 & g(\theta^t x_i) < 0,5 \end{cases}$

Regresión Logística

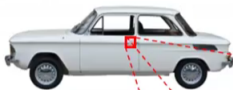
Función de costo

$$J(\theta) = -\frac{1}{n} \sum_{i=1}^n (y_i \log(h_{\theta}(x_i))) + (1 - y_i)(1 - \log(h_{\theta}(x_i))) \\ + \frac{\lambda}{2m} \sum_{i=1}^m \theta_i^2$$

$$\hat{\theta} = \operatorname{argmin}(J(\theta))$$

Desventajas del modelo Logit

- Problemas de reconocimiento visual pequeños tienen 50×50 píxeles.
- Hipótesis no lineales requieren muchas *features*
- 2500 *features* + términos cuadráticos ($\frac{n^2}{2}$) \approx 3,1 millones



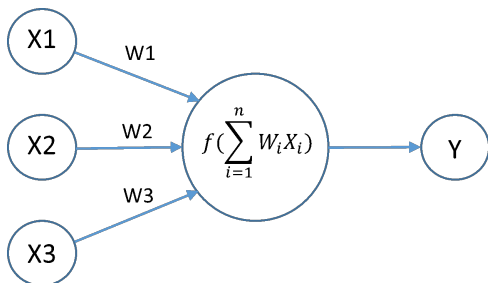
But the camera sees this:

194	210	201	212	199	213	215	195	178	158	182	209
180	189	190	221	209	205	191	167	147	115	129	163
114	126	140	188	176	165	152	140	170	106	78	88
87	103	115	154	143	142	149	153	173	101	57	57
102	112	106	131	122	138	152	147	128	84	58	66
94	95	79	104	105	124	129	113	107	87	69	67
68	71	69	98	89	92	98	95	89	88	76	67
41	56	68	99	63	45	60	82	58	76	75	65
20	43	69	75	56	41	51	73	55	70	63	44
50	50	57	69	75	75	73	74	53	68	59	37
72	59	53	66	84	92	84	74	57	72	63	42
67	61	58	65	75	76	76	73	59	75	69	50

Redes Neuronales

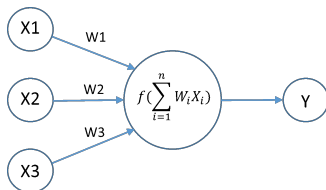
Neurona

- Unidad que recibe inputs y devuelve un output



- Por ejemplo una regresión logística

Feed Forward Networks



Cada neurona

- 1 Recibe un input x
- 2 Realiza una transformación lineal $(w_i x + b_i)$
- 3 Aplica una transformación no lineal g , para obtener $out = g(w_i x + b_i)$

Redes Neuronales

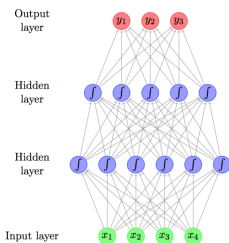


Figure 2: Feed-forward neural network with two hidden layers.

- Modelos no paramétricos de Machine Learning compuestos de unidades computacionales (**neuronas**)
- Aproximar la función $f : \mathcal{F} \rightarrow \mathcal{O}$ que relaciona un conjunto de features $\mathcal{F} \subset \mathbb{R}^n$ con un conjunto de outputs $\mathcal{O} \subset \mathbb{R}$

Redes Neuronales

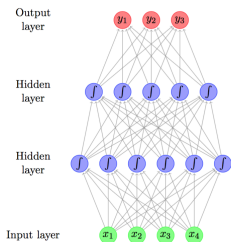


Figure 2: Feed-forward neural network with two hidden layers.

- Concatenación de capas de neuronas
 - Capa *input*
 - Capas ocultas
 - Capa *output*

Capa externa

- Regresión: Una neurona
- Clasificación Binaria: Probabilidad de éxito
- Clasificación k -categorías: Probabilidad de pertenecer a cada categoría. Para el entrenamiento, las etiquetas se convierten en vectores:

$$y_i = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

Feed Forward Networks

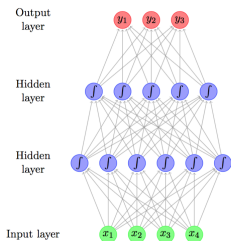


Figure 2: Feed-forward neural network with two hidden layers.

$$NN_{MLP2}(x) = y$$

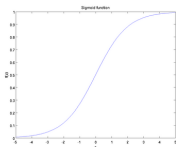
$$\mathbf{h}^1 = g^1(\mathbf{W}^1\mathbf{x} + \mathbf{b}^1)$$

$$\mathbf{h}^2 = g^2(\mathbf{W}^2\mathbf{h}^1 + \mathbf{b}^2)$$

$$\mathbf{y} = \mathbf{h}^2\mathbf{W}^3$$

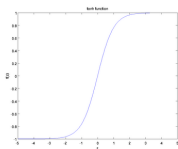
- Si las capa j y $j + 1$ tienen s_j y s_{j+1} neuronas respectivamente, $\mathbf{W}_{j+1} \in \mathbb{R}^{s_{j+1} \times s_j}$
- Si usa un término de sesgo, $\mathbf{W}_{j+1} \in \mathbb{R}^{s_{j+1} \times s_j + 1}$

Funciones de Activación



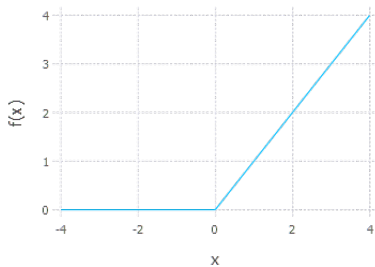
$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}}$$

Funciones de Activación



$$\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$$

Funciones de Activación



$$\text{ReLU}(x) = \max(0, x)$$

Función de Costo

$$J(\theta) = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K (y_i^k \log(h_{\theta}(x_i)))^k + (1 - y_i^k)(1 - \log(h_{\theta}(x_i)))_k \\ + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_j} \sum_{j=1}^{s_{j+1}} \Theta^{(l)}_{i,j}^2$$

$$\hat{\theta} = \operatorname{argmin}(J(\theta))$$

Entrenamiento

- Mejor más neuronas que menos, capturar más *no linearidades*
- Evitar overfitting con regularización
- Escalar las variables
- Mínimos locales: Hacer varias inicializaciones.

Minimización de la función de costo

Algoritmo *Backpropagation*

Aplicaciones

Carros autónomos

<https://www.coursera.org/learn/machine-learning/lecture/zYS8T/autonomous-driving>

Medicina

- Diagnóstico a partir de imágenes
- Diagnóstico a partir de lenguaje
- Genómica

Redes en el contexto de NLP

Input x

- En las aplicaciones de NLP, el input x codifica palabras, POS, información lingüística.
- Para cada documento se cuenta, en general, con un conjunto de vectores.
- Se utiliza una función $c : \mathcal{F} \rightarrow \mathbb{R}^d$
- c 's comunes incluyen concatenación y suma.

Embeddings

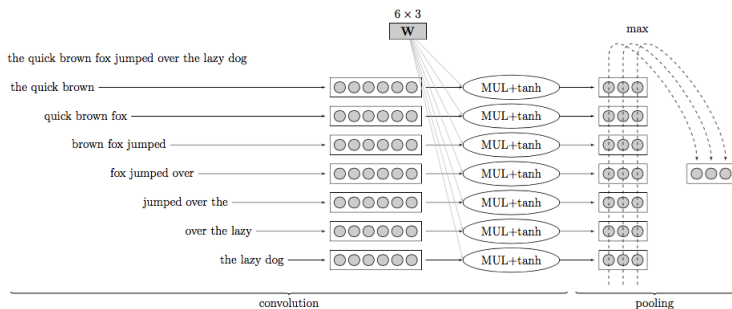
- Supervisados: Entrenar redes para contexto particular ('one hot representations')
- No supervisados: GloVe, LDA y modelos de tópicos.

Otras arquitecturas

Convolutional Networks

- Arquitectura que permite detectar features importantes sin importar la ubicación
- Ha permitido avances importantes en NLP y procesamiento de imágenes
- Aplica un *filtro* (función no lineal aprendida) a cada k -ventana de palabras para crear un vector
- Se hace *pooling* para combinar estos vectores en uno solo que concentre las características más importantes sin importar la ubicación

Convolutional Networks



Convolutional Networks

- Sea x_1, x_2, \dots, x_n una secuencia de palabras con vectores $v(x_1), \dots, v(x_n)$
- Para cada ventana i sea $w_i = [v_{i+1}, \dots, v_{i+k}]$
- Convolución: Calcule $p_i = g(w_i W + b)$
- Pooling: Calcule

$$\mathbf{c} = [c_j]$$

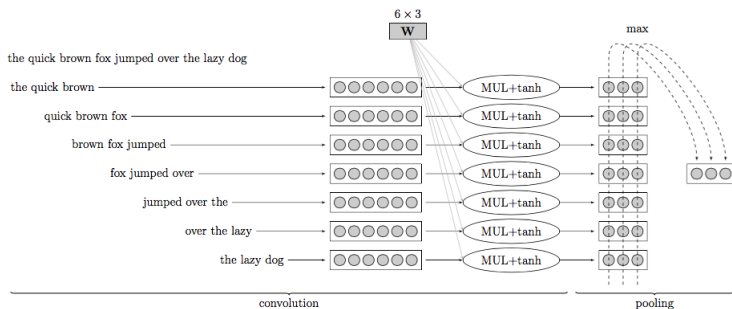
donde

$$c_j = \max_{1 < i \leq m} p_{ij}$$

Convolutional Networks

- Cada dimensión se especializará en una clase de predictores y la operación *max* escogerá el más importante de cada tipo

Convolutional Networks



Recurrent Networks

- Si bien las redes convolucionales tienen en cierta forma el orden, se restringen a patrones locales.
- Las Redes Neuronales Recurrentes permiten representar inputs de tamaños arbitrarios en un vector de tamaño fijo teniendo en cuenta propiedades estructurales.

Recurrent Networks

- Dados vectores x_1, \dots, x_n y un vector s_0 que será el estado inicial

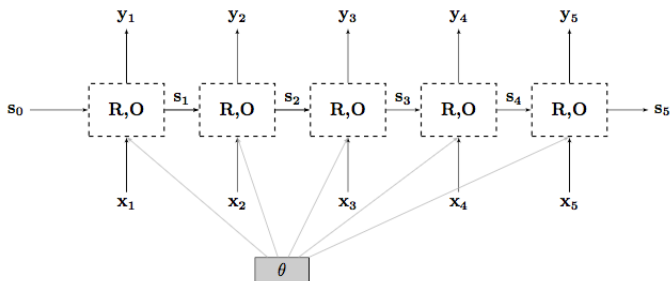
$$RNN(s_0, x_{1:n}) = s_{1:n}, y_{1:n}$$

$$s_i = R(s_{i-1}, x_{1:i})$$

$$y_i = O(s_i)$$

Recurrent Networks

Figura: Red Neuronal recurrente



Recurrent Networks

El output puede ser únicamente el último estado o utilizar también los intermedios

- **Acceptor:** Solo se utiliza el último estado para determinar y calcular la función de pérdida a partir de y_n (Análisis de sentimiento)
- **Encoder:** Creación de features para otra tarea (por ejemplo resumir un documento). Solo se utiliza y_n
- **Transducer:** Se utilizan todos los estados y outputs. Por ejemplo, encontrar la distribución de la palabra ii a partir de las palabras $1 : (i - 1)$
- **Encoder - Decoder:** Se utiliza un encoder para producir y_n que se entrega como input a otra RNN que utiliza las palabras ya traducidas y el último output para decodificar. Muy buenos resultados para traducir (o poner tags a las palabras)

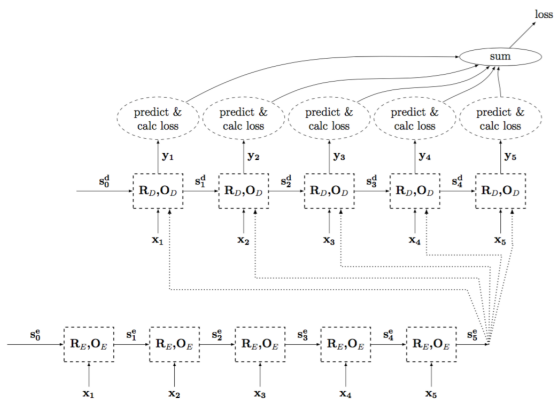


Figure 9: Encoder-Decoder RNN Training Graph.

Recurrent Networks – S-RNN

- **Red Recurrente Simple:** Muy buenos resultados en *sequence tagging* y modelaje de lenguaje.

$$s_i = g(x_i, W^x + s_{i-1}W^s + b)$$

$$y_i = O(s_i) = s_i$$

Recurrent Networks –S-RNN

- Simple RNN es difícil de entrenar por problema de *vanishing gradients*
- Difícil capturar dependencias muy lejanas.

Recurrent Networks – LSTM

- Introducir en el estado s celdas de memoria que preserven los gradientes.
- El acceso a la memoria está controlado por *puertas*. Estas son funciones con rango en $[0, 1]^n$

Recurrent Networks – LSTM

- **Long Short Term Memory:**

$$\mathbf{s}_j = R_{LSTM}(\mathbf{s}_{j-1}, \mathbf{x}_j) = [\mathbf{c}_j; \mathbf{h}_j]$$

$$\mathbf{c}_j = \mathbf{c}_{j-1} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i}$$

$$\mathbf{h}_j = \tanh(\mathbf{c}_j) \odot \mathbf{o}$$

$$\mathbf{i} = \sigma(\mathbf{x}_j \mathbf{W}^{xi} + \mathbf{h}_{j-1} \mathbf{W}^{hf})$$

$$\mathbf{f} = \sigma(\mathbf{x}_j \mathbf{W}^{xf} + \mathbf{h}_{j-1} \mathbf{W}^{hf})$$

$$\mathbf{o} = \sigma(\mathbf{x}_j \mathbf{W}^{xo} + \mathbf{h}_{j-1} \mathbf{W}^{ho})$$

$$\mathbf{g} = \tanh(\mathbf{x}_j \mathbf{W}^{xg} + \mathbf{h}_{j-1} \mathbf{W}^{hg})$$

$$y_j = O(s_j) = h_j$$

Recurrent Networks – GRU

- LSTM es efectiva pero complicada y computacionalmente costosa.
- **GRU** (Gated Recurrent Unit)

Recurrent Networks – GRU

$$\mathbf{s}_j = R_{GRU}(\mathbf{s}_{j-1}, \mathbf{x}_j) = (\mathbf{1} - \mathbf{z}) \odot \mathbf{s}_{j-1} + \mathbf{z} \odot \mathbf{h}$$

$$\mathbf{z} = \sigma(\mathbf{x}_j \mathbf{W}^{xz} + \mathbf{h}_{j-1} \mathbf{W}^{hz})$$

$$\mathbf{r} = \sigma(\mathbf{x}_j \mathbf{W}^{xr} + \mathbf{h}_{j-1} \mathbf{W}^{hr})$$

$$\mathbf{h} = \tanh(\mathbf{x}_j \mathbf{W}^{xh} + (\mathbf{h}_{j-1} \odot \mathbf{r}) \mathbf{W}^{hg})$$

$$y_j = O(\mathbf{s}_j) = h_j$$

GRACIAS