



DECSAI

Departamento de Ciencias de la Computación e I.A.

Universidad de Granada



Deep Stacking Networks [DSNs]

Fernando Berzal, berzal@acm.org

Deep Stacking Networks (DSNs)

- Las técnicas convencionales para entrenar DNNs, en su etapa final de ajuste [fine tuning] requieren la utilización del gradiente descendente estocástico, que resulta difícil de paralelizar.
- La arquitectura de las DSNs (originalmente llamadas "Deep Convex Networks" o DCNs) se diseñó con el problema de la escalabilidad del aprendizaje en mente.

Li Deng & Dong Yu: "Deep convex network: A scalable architecture for speech pattern classification". Interspeech'2011.



Deep Stacking Networks (DSNs)

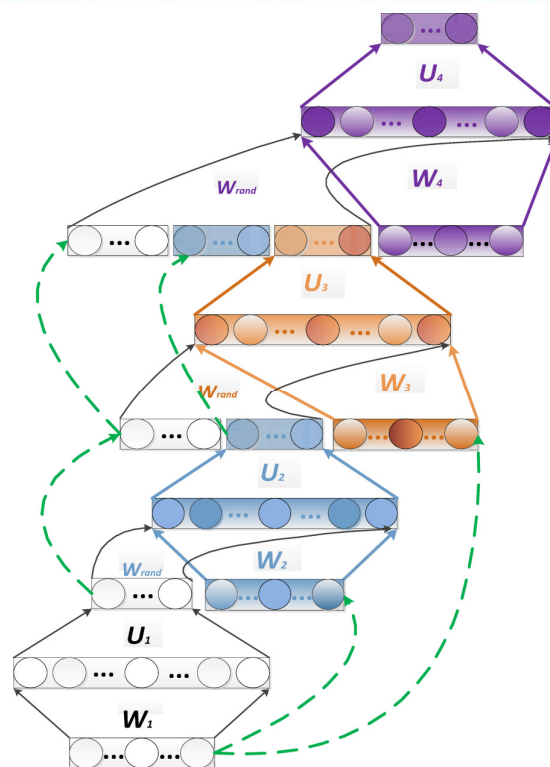
- Las DSNs apilan redes multicapa como módulos básicos en los que las unidades de salida son lineales y las unidades ocultas son sigmoidales (no lineales).
- La linealidad de las unidades de salida permite una estimación eficiente y paralelizable de los pesos de la capa de salida dada la actividad de la capa oculta (al tratarse de un problema de optimización convexa).
- Las restricciones impuestas por la arquitectura hacen mucho más fácil el aprendizaje de los demás parámetros de la red (los pesos de la capa oculta).



Deep Stacking Networks (DSNs)

La matriz de pesos W conecta la capa lineal de entrada con la capa oculta no lineal.

La matriz de pesos U conecta la capa oculta (no lineal) con la capa lineal de salida.



Deep Stacking Networks (DSNs)

Algoritmo de aprendizaje

Dados los vectores de salida para el conjunto de entrenamiento T , los pesos U y W se ajustan para minimizar el error:

$$E = \frac{1}{2} \sum_i \|y_i - t_i\|^2 = \frac{1}{2} \text{Tr}[(Y - T)(Y - T)^T]$$

donde Y es la salida de la red:

$$y_i = U^T h_i = U^T \sigma(W^T x_i) = G_i(UW)$$



Deep Stacking Networks (DSNs)

Algoritmo de aprendizaje

U puede aprenderse una vez que se tiene la matriz de actividad H para la capa oculta (o lo que es equivalente, cuando se conoce W).

Fijando la derivada del error con respecto a U a cero, obtenemos

$$U = (HH^T)^{-1}HT^T = F(W), \quad \text{where } h_i = \sigma(W^T x_i).$$

Esto nos proporciona una restricción explícita entre U y W que podemos aprovechar en la minimización del error



Deep Stacking Networks (DSNs)

Algoritmo de aprendizaje

Dada la restricción $U=F(W)$,
aplicamos el método de los multiplicadores de Lagrange:

$$E = \frac{1}{2} \sum_i \|G_i(U, W) - t_i\|^2 + \lambda \|U - F(W)\|$$

Ahora, podemos utilizar "full-batch learning":

$$\frac{\partial E}{\partial W} = 2X[H^T \circ (1 - H)^T \circ [H^\dagger(HT^T)(TH^\dagger) - T^T(TH^\dagger)]]$$

$$H^\dagger = H^T(HH^T)^{-1}$$



Deep Stacking Networks (DSNs)

Método de Lagrange

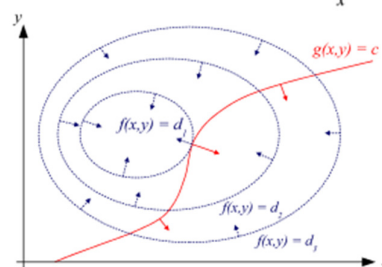
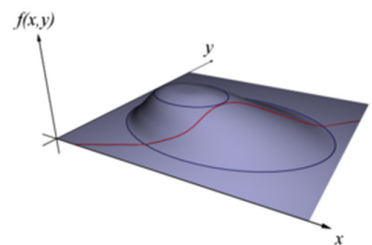
Método de optimización de funciones
de múltiples variables sujetas a restricciones

Función objetivo $f(\theta)$

Restricción $g(\theta) = 0$

Los gradientes $\nabla f(\theta)$ y $\nabla g(\theta)$
son paralelos, por lo que

$$\nabla f(\theta) = \lambda \nabla g(\theta)$$



Deep Stacking Networks (DSNs)

Método de Lagrange

Podemos incorporar la restricción $g(\theta)=0$ a nuestra función objetivo, añadiendo una nueva variable λ :

$$\mathbf{F}(\theta,\lambda) = \mathbf{f}(\theta) - \lambda\mathbf{g}(\theta)$$

F se denomina Lagrangiano.

Optimizamos, como siempre, igualando a cero su gradiente:

$$\nabla\mathbf{F}(\theta,\lambda) = \mathbf{0}$$



Deep Stacking Networks (DSNs)

Método de Lagrange

EJEMPLO

Cómo minimizar la superficie de una caja de volumen V.

Función objetivo: $f(x,y,z) = 2(xy+xz+yz)$

Restricción: $g(x,y,z) = xyz - V = 0$

Gradiente $\nabla f(x,y,z) = \langle 2(y+z), 2(x+z), 2(x+y) \rangle$

$$\nabla g(x,y,z) = \langle yz, xz, xy \rangle$$

Al hacer $\nabla f(x,y,z) = \lambda \nabla g(x,y,z)$

obtenemos 3 ecuaciones que nos dan la solución:

$$x = y = z = 4/\lambda \quad (\text{un cubo :-})$$

