

# Decomposição em Valores Singulares (SVD)

Uma aplicação para *prunning* de redes neurais

André Pacheco

*Universidade Federal do Espírito Santo*

The diagram illustrates the SVD decomposition of a matrix  $A$ . Matrix  $A$  is shown as a vertical rectangle with height  $m$  and width  $n$ . It is equal to the product of three matrices:  $U$ ,  $\Sigma$ , and  $V^T$ . Matrix  $U$  is a square with side length  $m$ . Matrix  $\Sigma$  is a vertical rectangle with height  $m-r$  and width  $n-r$ ; its top-left  $r \times r$  sub-region is shaded with a diagonal line, representing the non-zero singular values. Matrix  $V^T$  is a square with side length  $n$ .

$$\begin{matrix} n \\ m \end{matrix} A = \begin{matrix} m \\ m \end{matrix} U \begin{matrix} r & n-r \\ m-r & \end{matrix} \Sigma \begin{matrix} n \\ n \end{matrix} V^T$$

Novembro de 2016

# SUMÁRIO

1. Introdução
2. Objetivo
3. Decomposição em Valores Singulares (SVD)
  - 3.1 Conceitos Básicos
  - 3.2 Algoritmo de aplicação
4. *Prunning* de rede neural
5. Resultados
6. Considerações finais

# INTRODUÇÃO

- ▶ A decomposição em valores singulares (SVD) é uma importante ferramenta na ciência da computação, análise de dados e estatística
- ▶ É utilizada em diversas aplicações, dentre elas:
  - ▶ Reconhecimento de padrões
  - ▶ Filtro de sinais
  - ▶ Problemas de otimização
  - ▶ Compressão de dados
  - ▶ Recuperação da informação
- ▶ Além das aplicações práticas, os resultados da SVD produz diversos corolários na álgebra linear

# INTRODUÇÃO

- ▶ A ideia principal da SDV é decompor uma matriz alvo no produto de três outras matrizes
  - ▶ O intuito é facilitar a resolução de um determinado problema
- ▶ O *pruning* de redes neurais é uma técnica que visa encontrar um número adequado de neurônios em uma camada oculta
  - ▶ Importante para o treinamento da rede
  - ▶ Evitar *overfitting* e convergência lenta

# OBJETIVO

- ▶ O objetivo deste seminário é introduzir os conceitos da SVD e apresentar uma aplicação para *pruning* de redes neurais

# CONCEITOS BÁSICOS

- ▶ Seja  $\vec{v}$  um vetor de tamanho  $n$ , sua **norma** é definida como:

$$|\vec{v}| = \sqrt{\sum_{i=1}^n v_i^2} \quad (1)$$

- ▶ Seja  $\vec{v}^1$  e  $\vec{v}^2$  dois vetores de tamanho  $n$ , o **produto interno** (produto escalar ou *dot*) entre eles é definido como:

$$\text{dot}(\vec{v}^1, \vec{v}^2) = \sum_{i=1}^n v_i^1 v_i^2 \quad (2)$$

# CONCEITOS BÁSICOS

- ▶ Dois vetores são **ortogonais** se o produto interno entre eles é igual a 0
  - ▶ Em um espaço bidimensional é equivalente a dizer que os vetores são perpendiculares entre si, ou seja, o ângulo formado entre eles é de  $90^\circ$
- ▶ Um **vetor normal** (ou unitário) é aquele com norma = 1. Qualquer vetor com norma  $> 0$  pode ser normalizado dividindo o mesmo pela norma
  - ▶ Seja  $\vec{v} = [2, 4, 1, 2]$ , então  $|\vec{v}| = \sqrt{2^2 + 4^2 + 1^2 + 2^2} = \sqrt{25} = 5$
  - ▶ Fazendo  $\vec{u} = [\frac{2}{5}, \frac{4}{5}, \frac{1}{5}, \frac{2}{5}]$ ,  $\vec{u}$  se torna normal, pois:
$$|\vec{u}| = \sqrt{\frac{2^2}{5} + \frac{4^2}{5} + \frac{1^2}{5} + \frac{2^2}{5}} = \sqrt{1} = 1$$

# CONCEITOS BÁSICOS

- ▶ Sejam dois vetores **normais**  $\vec{u}$  e  $\vec{v}$ . Se esses vetores **ortogonais** entre si, eles são chamados de **vetores ortonormais**. Por exemplo:
  - ▶ Seja  $\vec{u} = [\frac{2}{5}, \frac{1}{5}, \frac{-2}{5}, \frac{4}{5}]$  e  $\vec{v} = [\frac{3}{\sqrt{65}}, \frac{-6}{\sqrt{65}}, \frac{4}{\sqrt{65}}, \frac{2}{\sqrt{65}}]$
  - ▶ Ambos são normalizados, ou seja,  $|\vec{u}| = 1$  e  $|\vec{v}| = 1$
  - ▶ O produto interno  $dot(\vec{u}, \vec{v}) = 0$
  - ▶ Os vetores  $\vec{u}$  e  $\vec{v}$  são **ortonormais**



# CONCEITOS BÁSICOS

- ▶ Uma matriz  $A$  é **ortogonal** se os vetores que formam sua coluna são **ortogonais**. Por exemplo, seja

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 3/5 & -4/5 \\ 0 & 4/5 & 3/5 \end{bmatrix}$$

- ▶  $A^T A = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 3/5 & -4/5 \\ 0 & 4/5 & 3/5 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 3/5 & 4/5 \\ 0 & -4/5 & 3/5 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$
- ▶ A matriz ortogonal  $A$  satisfaz as seguintes propriedades:
  - ▶  $A^T A = I$
  - ▶  $A^{-1} = A^T$
  - ▶ O produto de duas matrizes ortogonais gera outra matriz ortogonal

# CONCEITOS BÁSICOS

- ▶ Um **autovetor** é um vetor não nulo que satisfaz a equação

$$A\vec{v} = \lambda\vec{v} \quad (3)$$

- ▶  $A$  é uma matriz quadrada,  $\lambda$  são os **autovalores** e  $\vec{v}$  os autovetores
  - ▶ Os autovalores podem ser calculados solucionando o sistema linear obtido a partir da equação  $\det(A - \lambda I) = 0$
  - ▶ Os autovetores associados aos autovalores são obtidos por meio da equação  $(\lambda I - R)\vec{v} = 0$
- ▶ Os **valores singulares** de uma matriz são definidos como:

$$\sigma_j = \sqrt{\lambda_j(A^T A)} \quad (4)$$

# CONCEITOS BÁSICOS

## ► Processo de ortogonalização de Gram-Schmidt

- Tem como objetivo transformar um conjunto de vetores em vetores ortogonais
- O algoritmo possui as seguintes regras:
  1. Organizar os vetores em uma matriz. Cada vetor será uma coluna;
  2. Normalizar o primeiro vetor;
  3. Iterativamente reescrever os demais vetores em termos deles mesmos menos os vetores já normalizados:

$$\vec{w}_k = \vec{v}_k - \sum_{i=1}^{k-1} \text{dot}(\vec{u}_i, \vec{v}_k) \times \vec{u}_i \quad (5)$$

sendo  $\vec{w}_k$  o vetor resultante,  $\vec{v}_k$  o vetor atual,  $\vec{u}_i$  o vetor ortonormal anterior e  $k$  o número de vetores

## CONCEITOS BÁSICOS

- ▶ Considere 3 vetores organizados na forma da matriz  $A$ , sendo cada coluna, um vetor:

$$A = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 2 & 0 \\ 2 & 3 & 1 \\ 1 & 1 & 0 \end{bmatrix} \quad (6)$$

- ▶ **Passo 1:** Normalizar o vetor  $\vec{v}_1 = [1, 0, 2, 1]$ . Como  $|\vec{v}_1| = \sqrt{6}$ , o vetor normalizado de  $\vec{v}_1$  será  $\vec{u}_1 = [\frac{1}{\sqrt{6}}, 0, \frac{2}{\sqrt{6}}, \frac{1}{\sqrt{6}}]$
- ▶ **Passo 2:** Calcular os vetores restantes de acordo com a equação 4. Neste exemplo:

$$\vec{w}_2 = \vec{v}_2 - \text{dot}(\vec{u}_1, \vec{v}_2) \times \vec{u}_1$$

$$\vec{w}_3 = \vec{v}_3 - \text{dot}(\vec{u}_2, \vec{v}_3) \times \vec{u}_2$$

## CONCEITOS BÁSICOS

- ▶ **Passo 2:** Calcular os vetores restantes de acordo com a equação 4.

$$\begin{aligned}\vec{w}_2 &= \vec{v}_2 - \text{dot}(\vec{u}_1, \vec{v}_2) \times \vec{u}_1 \\ &= [2, 2, 3, 1] - \text{dot}\left(\left[\frac{1}{\sqrt{6}}, \frac{1}{\sqrt{6}}, 0, \frac{2}{\sqrt{6}}\right], [2, 2, 3, 1]\right) \times \left[\frac{1}{\sqrt{6}}, \frac{1}{\sqrt{6}}, 0, \frac{2}{\sqrt{6}}\right] \\ \vec{w}_2 &= \left[\frac{1}{2}, 2, 0, \frac{-1}{2}\right]\end{aligned}$$

$$\text{Normalizando } \vec{w}_2: \vec{u}_2 = \left[\frac{\sqrt{2}}{6}, \frac{2\sqrt{2}}{3}, 0, \frac{-\sqrt{2}}{6}\right]$$

- ▶ Fazendo esse mesmo processo para

$$\vec{w}_3 = \vec{v}_3 - \text{dot}(\vec{u}_2, \vec{v}_3) \times \vec{u}_2, \text{ é obtido: } \vec{u}_3 = \left[\frac{2}{3}, \frac{-1}{3}, 0, \frac{-2}{3}\right]$$

- ▶ A matriz ortogonal será:  $A = \begin{bmatrix} \frac{1}{\sqrt{6}} & \frac{\sqrt{2}}{6} & \frac{2}{3} \\ 0 & \frac{2\sqrt{2}}{3} & \frac{-1}{3} \\ \frac{2}{\sqrt{6}} & 0 & 0 \\ \frac{1}{\sqrt{6}} & \frac{-\sqrt{2}}{6} & \frac{-2}{3} \end{bmatrix}$

# DECOMPOSIÇÃO EM VALORES SINGULARES (SVD)

- ▶ Dado uma matriz  $A \in \mathbb{R}$ , a SVD de  $A$  é a fatorização de  $A$  representada pelo produto de três matrizes:

$$A_{m \times n} = U_{m \times m} D_{m \times n} V_{n \times n}^T \quad (7)$$

- ▶  $U$  é uma matriz ortogonal, sendo que suas colunas são autovetores ortonormais de  $AA^T$
- ▶  $V$  é uma matriz ortogonal, sendo que suas colunas são autovetores ortonormais de  $A^T A$
- ▶  $D$  é uma matriz diagonal contendo a raiz quadrada dos autovalores de  $U$  ou  $V$  em ordem decrescente
- ▶ A ortogonalização das matrizes é obtida por meio do **Processo de Gram-Schmidt**

# DECOMPOSIÇÃO EM VALORES SINGULARES (SVD)

- ▶ Demonstrando o algoritmo por meio de um exemplo, considera a matriz:

$$A = \begin{bmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{bmatrix}_{2 \times 3} \quad (8)$$

- ▶ Primeiramente, vamos calcular  $U_{2 \times 2}$

$$AA^T = \begin{bmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{bmatrix} \times \begin{bmatrix} 3 & -1 \\ 1 & 3 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 11 & 1 \\ 1 & 11 \end{bmatrix}_{2 \times 2}$$

- ▶ Na sequência, precisamos calcular os autovalores e autovetores de  $AA^T$
- ▶ Realizando cálculos mais rápido do que a velocidade da luz, temos os autovalores 12 e 10, com autovetores associados sendo  $[1, 1]$  e  $[-1, 1]$
- ▶ Esses autovetores se tornam colunas da matriz  $\check{U}$ , ordenados pelos autovalores correspondentes (decrecente)

## DECOMPOSIÇÃO EM VALORES SINGULARES (SVD)

- ▶ Portanto:  $\check{U} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}$
- ▶ Para obter a matriz ortogonal  $U$ , é realizado o processo de Gram-Schmidt (calculando novamente mais rápido que a velocidade da luz):

$$U = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix}_{2 \times 2} \quad (9)$$

- ▶ De maneira similar, calculamos a matriz  $V^T_{3 \times 3}$  :

$$A^T A = \begin{bmatrix} 3 & -1 \\ 1 & 3 \\ 1 & 1 \end{bmatrix} \times \begin{bmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{bmatrix} = \begin{bmatrix} 10 & 0 & 2 \\ 0 & 10 & 4 \\ 2 & 4 & 2 \end{bmatrix}_{3 \times 3}$$



## DECOMPOSIÇÃO EM VALORES SINGULARES (SVD)

- ▶ Os autovalores de  $A^T A$  são: 0, 10 e 12. Os autovetores associados são:  $[1, 2, -5]$ ,  $[2, -1, 0]$  e  $[1, 2, 1]$ , respectivamente

- ▶ Portanto:  $\check{V} = \begin{bmatrix} 1 & 2 & 1 \\ 2 & -1 & 2 \\ 1 & 0 & -5 \end{bmatrix}$

- ▶ Aplicando o processo de Gram-Schmidt e em seguida transpondo a matriz ortogonalizada, obtemos:

$$V^T = \begin{bmatrix} \frac{1}{\sqrt{6}} & \frac{2}{\sqrt{6}} & \frac{1}{\sqrt{30}} \\ \frac{2}{\sqrt{5}} & \frac{-1}{\sqrt{5}} & 0 \\ \frac{1}{\sqrt{30}} & \frac{2}{\sqrt{30}} & \frac{-5}{\sqrt{30}} \end{bmatrix}_{3 \times 3} \quad (10)$$

# DECOMPOSIÇÃO EM VALORES SINGULARES (SVD)

- ▶ Por fim, vamos obter a matriz  $D_{2 \times 3}$ 
  - ▶ Utilizamos a raiz dos autovalores de  $U$  ou  $V$
  - ▶ Os autovalores são colocados na diagonal principal em ordem decrescente
  - ▶ Os autovalores diferentes de zeros de  $U$  e  $V$  serão sempre iguais
  - ▶ A última coluna de zeros mantém as propriedades da multiplicação de matrizes

$$D = \begin{bmatrix} \sqrt{12} & 0 & 0 \\ 0 & \sqrt{10} & 0 \end{bmatrix}_{2 \times 3} \quad (11)$$

- ▶ A diagonal de  $D$  contém os valores singulares de  $A$

# DECOMPOSIÇÃO EM VALORES SINGULARES (SVD)

- Portanto,  $A$  foi decomposta nas 3 matrizes  $U$ ,  $D$  e  $V^T$ :

$$A = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix}_{2 \times 2} \times \begin{bmatrix} \sqrt{12} & 0 & 0 \\ 0 & \sqrt{10} & 0 \end{bmatrix}_{2 \times 3} \times \begin{bmatrix} \frac{1}{\sqrt{6}} & \frac{2}{\sqrt{6}} & \frac{1}{\sqrt{30}} \\ \frac{2}{\sqrt{5}} & \frac{-1}{\sqrt{5}} & 0 \\ \frac{1}{\sqrt{30}} & \frac{2}{\sqrt{30}} & \frac{-5}{\sqrt{30}} \end{bmatrix}_{3 \times 3} \quad (12)$$

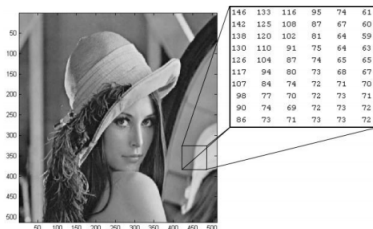
- Podemos visualizar  $A$ , da seguinte forma:

$$A = \sum_{j=1}^r \sigma_j u_j v_j^T \quad (13)$$

Sendo  $r$  o posto da matriz

# COMPRESSÃO DE DADOS COM SVD

- ▶ Assim como PCA, SVD também pode ser utilizado como uma maneira de comprimir dados
- ▶ Considere a imagem  $512 \times 512$  *pixels* em níveis de cinza:



- ▶ Ao decompor a matriz  $X$ , que representa a imagem, os 512 valores singulares se distribuem entre  $[5 \times 10^4, 2 \times 10^{-2}]$
- ▶ Portanto:  $X = \sum_{j=1}^{512} \sigma_j u_j v_j^T$
- ▶ Cada componente vai contribuir de forma bem diferente devido aos valores singulares

# COMPRESSÃO DE DADOS COM SVD

- Considerando os 20 primeiros (e maiores) valores singulares e na sequência os 50 primeiros, obtemos:



a) 20 valores singulares



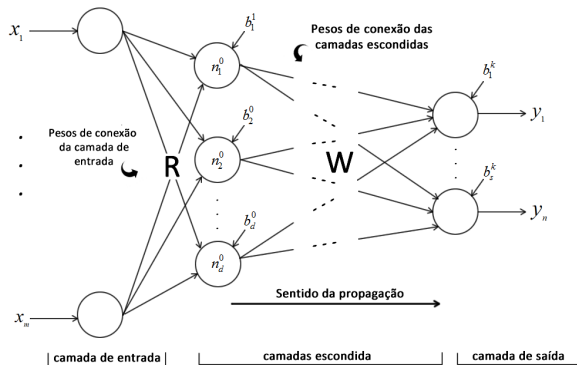
b) 50 valores singulares

## *Pruning* DE REDES NEURAIIS UTILIZANDO SVD

- ▶ Determinar o **número de neurônios de uma camada** de uma rede neural é estritamente ligado ao desempenho da mesma
  - ▶ Uma má escolha pode ocasionar *overfitting* ou não convergência do treinamento
- ▶ Uma maneira atacar este problema é **alargar** o número de neurônios no início do algoritmo e de acordo com o resultado **diminuir** este número (*pruning*)
- ▶ Em 1994, Psychogios e Ungar [2], propuseram o SVD-NET
  - ▶ Algoritmo que aplicava *pruning* de maneira automática para remodelar camadas escondidas

# Pruning DE REDES NEURAIS UTILIZANDO SVD

- Considere uma rede neural de uma camada oculta



- A priori é escolhido um número grande de neurônios escondidos

## *Pruning* DE REDES NEURAIS UTILIZANDO SVD

- ▶ Inicializado os pesos  $R$  e  $W$  de maneira aleatória, a SVD-NET atua da seguinte forma:
  1. Os pesos de  $R$  são ajustados por meio da minimização do erro quadrático. Os valores de  $W$  se mantêm fixos;
  2. Conhecidos os valores de  $R$ , cada neurônio oculto  $n$  pode ser calculado a parti das entradas  $x$ , gerando uma matriz  $N_{m \times d}$ ;
  3. Os pesos de  $W$  são ajustados minimizando  $E = |R \times W - \hat{y}|$
  4. Aplicar SVD na matriz  $A$  em busca de neurônios ocultos redundantes
    - 4.1 Um número não nulo de valores singulares pequenos indicam redundância
    - 4.2 Aplicar *pruning* nos neurônios ocultos
- ▶ Os autores afirmam que o mesmo procedimento pode ser aplicado para redes de mais camadas



# RESULTADOS

- ▶ No trabalho Abid e Najim [3], os autores aplicam essa metodologia para dois problemas simples: aproximar a função seno e reconhecer padrões de uma entrada
- ▶ Para ambos os casos, foram utilizadas uma rede neural tradicional e a SVD-NET. Ambas com apenas uma camada.

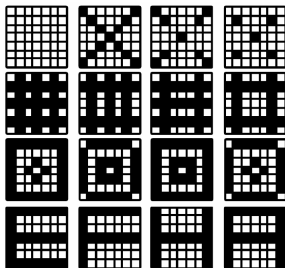
**Problema I:** aproximar  $f(x) = \sin(x)$  no intervalo  $0 < x < 4\pi$

	Neurônios	Iteracoes para converg.	Tempo (s)
NN	15	2238	12.13
SVD-NET	11	1776	10.88

Table: Resultado para aproximação da função seno

# RESULTADOS

**Problema II:** reconhecer os padrões apresentados na figura abaixo



	Neurônios	Iteracoes para converg.	Tempo (s)
NN	38	632	34.25
SVD-NET	29	479	29.36

Table: Resultado para o reconhecimento dos padrões

## CONSIDERAÇÕES FINAIS

- ▶ O SVD é uma metodologia clássica da álgebra linear com aplicações em diversos problemas reais
- ▶ Utilizar o SVD para *pruning* de redes neurais foi proposto há mais de 2 décadas
  - ▶ É um problema relevante se tratando de redes neurais
  - ▶ Todavia, a metodologia não teve tanta aceitação
- ▶ Um dos motivos é nova técnicas mais recentes, como *dropout* que trata do mesmo problema de maneira mais fácil e rápida
- ▶ Por outro lado, SDV continua sendo bastante utilizada para compressão de dados, juntamente com PCA

## REFERÊNCIAS

1. Baker, Kirk. "Singular value decomposition tutorial." The Ohio State University 24 (2005).
2. Psychogios, Dimitris C., and Lyle H. Ungar. "SVD-NET: An algorithm that automatically selects network structure." IEEE Transactions on Neural Networks 5.3 (1994): 513-515.
3. Abid, S., F. Fnaiech, and M. Najim. "A new Neural Network pruning method based on the singular value decomposition and the weight initialisation." Signal Processing Conference, 2002 11th European. IEEE, 2002.