

O objetivo deste trabalho é de estudar uma implementação dos algoritmos do classificador perceptron com uma estrutura de rede de tipo *sparse* e comparar com a sua versão *full*. Um template de script acompanhado com funções de visualização é disponibilizado.

1 Definições e notações

O espaço dos atributos é $\mathcal{A} = \mathbb{R}^I$ onde notamos $x = (x_1, \dots, x_I)^T$ as componentes do vetor x . O espaço da classe para o Perceptron é $\mathcal{C} = \{-1, 1\}$. Consideramos o conjunto de dados, $D = (x^n, y^n)_{n=1, \dots, N}$ onde $x^n \in \mathbb{R}^I$ e $y^n \in \mathcal{C}$. Recordamos que a função $s \rightarrow \text{sign}(s)$ representa a função sinal.

1.1 Fully Connected Perceptron (FCP)

O classificador *Perceptron* é caracterizado pela arquitectura

$$x \in \mathbb{R}^I \rightarrow \hat{y}(x; \tilde{w}) = \text{sign}(\tilde{w}^T x)$$

onde $\tilde{w} \in \mathbb{R}^{I+1}$ são os parâmetros do classificador. Para avaliar a qualidade do classificador $\hat{y}(x; \tilde{w})$ com a base de dados D , introduzimos a função custo baseada no erro quadrado

$$E(\tilde{w}; D) = \frac{1}{N} \sum_{n=1}^N \frac{1}{2} (y^n - \hat{y}^n)^2, \quad \hat{y}^n = \hat{y}(x^n; \tilde{w}).$$

Recordamos que y^n é o valor real, dada pela base de dados enquanto \hat{y}^n representa a predição da classe. O objetivo é de determinar os argumentos \tilde{w} que minimizam a funcional:

$$\tilde{w}^* = \arg \min_{\tilde{w} \in \mathbb{R}^{I+1}} E(\tilde{w}; D).$$

1.2 Sparsy Connected Perceptron (SCP)

Alem do vetor \tilde{w} , introduzimos o conjunto de indices $S \subset \llbracket 1, I \rrbracket$ correspondente aos coeficientes ativos. Caracterizamos S com o vetor $\tilde{s} = (s_i)_{i=0}^I$ tal como $s_0 = 1$, $s_i = 1$ se $i \in S$, senão $s_i = 0$. Por este modo, temos

$$w_0 + \sum_{i \in S} w_i x_i = w_0 s_0 x_0 + \sum_{i=1}^I w_i s_i x_i = \sum_{i=0}^I w_i s_i x_i.$$

Definimos os *sparsy connected Perceptron* como

$$x \in \mathbb{R}^I \rightarrow \hat{y}(x; \tilde{w}, \tilde{s}) = \text{sign}\left(\sum_{i=0}^I s_i w_i x_i\right).$$

A função custo associada é dada por

$$E(\tilde{w}, \tilde{s}; D) = \frac{1}{N} \sum_{n=1}^N \frac{1}{2} (y^n - \hat{y}^n)^2, \quad \hat{y}^n = \hat{y}(x^n; \tilde{w}, \tilde{s}).$$

Nota que, se $S = \llbracket 1, I \rrbracket$, obtemos de novo o *Perceptron* usual onde todos os coeficientes de \tilde{w} são ativos.

O objetivo é de determinar os argumentos \tilde{w} que minimizam as funcionais:

$$\tilde{w}^*, \tilde{s}^* = \arg \min_{\substack{\tilde{w} \in \mathbb{R}^{I+1} \\ |\tilde{s}|_0 = K}} E(\tilde{w}, \tilde{s}; D),$$

onde K é um número inteiro dado, $K \leq I$ e $|\tilde{s}|_0$ representa a soma das componentes de \tilde{s} . Por outro lado, notamos por S^* the index subset associado to vector \tilde{s}^* .

2 Implementação do sparsy Perceptron

Recordamos o algoritmo (primal) do Perceptron. Construir uma sequência $\tilde{w}(t)$ tal como

$$\tilde{w}(t+1) = \tilde{w}(t) + \eta \frac{y^m - \hat{y}^m}{2} \tilde{x}^m$$

onde $m = m(t)$ é escolhido aleatoriamente para cada t e $\eta > 0$ a taxa de aprendizagem. A versão *sparsy* é baseada na construção do vetor \tilde{s} para desativar algumas células. A avaliação da qualidade é realizada pelo cálculo das matrizes de confusão usando o *in-samples* e *out-samples error*.

2.1 Algoritmo S -minimization

Seja S , um conjunto de índices ativos caracterizados pelo vetor \tilde{s} . Definimos o problema de minimização

$$\tilde{w}^S = \arg \min_{\tilde{w} \in \mathbb{R}^{S+1}} E(\tilde{w}, \tilde{s}; D)$$

onde o minimizante depende da escolha de S . Para este efeito, criamos uma sequência $\tilde{w}(t) \rightarrow \tilde{w}^S$ em \mathbb{R}^{S+1} com o algoritmo seguinte.

1. Escolher (x^m, y^m) , $m = m(t)$ da base de dados;
2. Calcular um predictor

$$\tilde{w}(t+1) = \tilde{w}(t) + \eta \frac{y^m - \hat{y}^m}{2} \tilde{x}^m, \quad \hat{y}^m = \hat{y}(x^m; \tilde{w}(t), \tilde{s}).$$

3. Projectar no espaço \mathbb{R}^{S+1} : $w(t+1) \in \mathbb{R}^{I+1} \rightarrow w(t+1) \in \mathbb{R}^{S+1}$

$$w_i(t+1) = w_i(t+1)s_i, \quad i = 0, \dots, I.$$

Quando convergir, notamos por $\tilde{w}^S \in \mathbb{R}^{S+1}$ a solução. Os trabalhos para realizar são:

- implementar o algoritmo;
- Usando a base de dados ..., experimentar diversos conjuntos S ;
- Produzir as tabelas de confusões usando um *training set* e um *test set*;
- Experimentar também com a base de dados ...

2.2 Algoritmo estocastico

Seja $K \leq I$. O método estocástico ou de Monte-Carlo consiste em experimentar vários conjuntos S escolhidos aleatoriamente tal como $|S| \leq K$ e criar uma sequência \tilde{s}^ℓ caracterizando o conjunto S^ℓ e $\tilde{w}^\ell \in \mathbb{R}^{S^\ell}$ tal como

$$E(\tilde{w}^{\ell+1}, \tilde{s}^{\ell+1}; D) \leq E(\tilde{w}^\ell, \tilde{s}^\ell; D).$$

O algoritmo é o seguinte.

1. Escolher um primeiro conjunto $|S^0| \leq K$ e determinar \tilde{w}^0 .
2. Seja S tal como $|S| \leq K$. Determinar \tilde{w}^S usando o algoritmo de S -minimização. Se $E(\tilde{w}^S, \tilde{s}; D) \leq E(\tilde{w}^\ell, \tilde{s}^\ell; D)$ então $\tilde{s}^{\ell+1} = \tilde{s}$ e $\tilde{w}^{\ell+1} = \tilde{w}^S$.
3. Até convergência.

O trabalho para realizar consiste em implementar o algoritmo e experimentá-lo usando a base de dados

- cria um conjunto S aleatoriamente tal como $|S| \leq K$;
- determinar \tilde{w}^S e compara-o para realizar a atualização como indicado no ponto (2).
- Produzir as tabelas de confusões usando um *training set* e um *test set*;
- Experimentar com a outra de base de dados ...

O algoritmo será utilizado para realizar a classificação da base de dados ...

2.3 Algoritmo Hard Thresholding

Seja $K \leq I$. Definimos a sequência $\tilde{w}(t), \tilde{s}(t)$ como

1. Escolher $(x^m, y^m) \in D$. Define

$$\tilde{w}(t+1) = \tilde{w}(t) + \eta \frac{y^m - \hat{y}^m}{2} \tilde{x}^m$$

com $\hat{y}^m = \hat{y}(x^m; \tilde{w}(t), s(t))$;

2. Definir $\tilde{s}(t+1)$ com $s_0(t+1) = 1$ e $s_i(t+1) = 1$ para as K maiores entradas $|w_i(t+1)|$, $i = 1, \dots, I$.

3. Calcular $\tilde{w}(t+1)$ como $w_i(t+1) = s_i(t+1)w_i(t+1)$.

O algoritmo acaba quando o erro $E(\tilde{w}(t), \tilde{s}(t); D)$ fica bloqueado por um valor limite por baixo.

- implementar o algoritmo de Hard Thresholding.
- Usando a base de dados ..., experimentar diversos valores de K . Produzir as tabelas de confusões usando um *training set* e um *test set*. Conclusões?
- Experimentar com a outra de base de dados ...

2.4 Algoritmo Iterative Hard Thresholding (IHT)

Este algoritmo é uma alternativa ao algoritmo anterior onde calculamos de novo o minimizante para cada iteração.

1. Escoler $(x^m, y^m) \in D$. Define

$$\tilde{w}(t+1) = \tilde{w}(t) + \eta \frac{y^m - \hat{y}^m}{2} \tilde{x}^m$$

com $\hat{y}^m = \hat{y}(x^m; \tilde{w}(t), \tilde{s}(t))$;

2. Definir $\tilde{s}(t+1)$ com $s_0(t+1) = 1$ e $s_i(t+1) = 1$ para as K maiores entradas $|w_i(t+1)|$, $i = 1, \dots, I$.
3. Determinar $\tilde{w}(t+1)$ como minimizante

$$\tilde{w}(t+1) = \arg \min_{\tilde{w} \in \mathbb{R}^{S(t+1)}} E(\tilde{w}, \tilde{s}(t+1); D)$$

calculado com o algoritmo de S -minimização.

Usando a base de dados ..., os trabalhos para realizar são de

- implementar o algoritmo de IHT;
- experimentar diversos valores de K e produzir as tabelas de confusões usando um *training set* e um *test set*;
- experimentar com a outra de base de dados ...