Quiz 4 - Supervised Learning Algorithms

Name *

Junhao Cheng

Email *

junhao.cheng24@imperial.ac.uk

Tree based Models

What is the primary criterion used to find the optimal split in a Classification Decision 1 point Tree?

Algorithm Decision Tree Learning Algorithm

Require: Training data $\{(F_i, y_i)\}_{i=1}^n$, stopping criteria **Ensure:** Decision tree T

- 1: Initialize tree with single root node containing all data
- 2: while nodes can be split and stopping criteria not met do
- 3: for each leaf node with region \mathcal{R} do
- 4: Find (j^*, τ^*) that maximizes:
- 5: $IG(j,\tau) = I(\mathcal{R}) \frac{|\mathcal{R}_L|}{|\mathcal{R}|}I(\mathcal{R}_L) \frac{|\mathcal{R}_R|}{|\mathcal{R}|}I(\mathcal{R}_R)$
- 6: Where $\mathcal{R}_L = \{F \in \mathcal{R} : F_j \leq \tau\}$ and $\mathcal{R}_R = \{F \in \mathcal{R} : F_j > \tau\}$
- 7: Split node using rule $F_{j^*} > \tau^*$
- 8: end for
- 9: end while
- 10: Assign prediction to each leaf node (majority class)
- 11: return T

Minimize Mean Squared Error

- Maximize Information Gain
- Minimize Sum of Squared Errors
- Gradient Descent

In a Regression Tree, what is the optimization objective when choosing a feature and 1 point threshold for splitting?

Algorithm Regression Tree Learning Algorithm

Require: Training data $\{(F_i, y_i)\}_{i=1}^n$, stopping criteria **Ensure:** Regression tree T

- 1: Initialize tree with single root node containing all data
- 2: while nodes can be split and stopping criteria not met do
- 3: for each leaf node with region \mathcal{R} do
- 4: Find (j^*, τ^*) that minimizes:
- 5: $SSE(j,\tau) = \sum_{i:F_i \in \mathcal{R}_L} (y_i \bar{y}_{\mathcal{R}_L})^2 + \sum_{i:F_i \in \mathcal{R}_R} (y_i \bar{y}_{\mathcal{R}_R})^2$

6: Where
$$\mathcal{R}_L = \{F \in \mathcal{R} : F_j \leq \tau\}$$
 and $\mathcal{R}_R = \{F \in \mathcal{R} : F_j \leq \tau\}$

- $F_j > \tau$ } 7: Split node using rule $F_{j^*} > \tau^*$
- 8: end for
- 9: end while
- 10: Assign prediction $\bar{y}_{\mathcal{R}_m}$ to each leaf node (average of y_i in the region)
- 11: return T

Maximize Information Gain

Minimize Sum of Squared Errors in resulting regions

O Maximize the number of samples in each region

O Minimize the entropy in each region





For the binary classification with labels {0, 1}, the weight update formula in AdaBoost 1 point can be written as:

$w_{t+1}(i) \propto w_t(i) \cdot e^{lpha_t y_i f_t(x_i)}$	$w_{t+1}(i) \propto w_t(i) \cdot e^{-lpha_t y_i f_t(x_i)}$
○ A	ОВ
$w_{t+1}(i) \propto w_t(i) \cdot e^{lpha_t (1-2\cdot \mathbb{I}\{y_i=f_t(x_i)\})}$	$w_{t+1}(i) \propto w_t(i) \cdot e^{-lpha_t(1-2\cdot \mathbb{I}\{y_i=f_t(x_i)\})}$

Neural Networks

In a shallow neural network with one hidden layer, what are the components of the parameter set θ ?



- Only the weights connecting input to hidden layer
- Only the weights connecting hidden to output layer
- Weights connecting input to hidden layer and hidden to output layer
- Weights and biases for both hidden and output layers

Which activation function has the range [0,1] and is commonly used as the output 1 point activation for binary classification problems?	nt	
O ReLU		
O Tanh		
Sigmoid		
O ELU		
For a multi-class classification problem with K classes, which loss function and final 1 point layer activation would typically be used?	nt	
O Mean Squared Error with linear activation		
O Binary Cross-Entropy with sigmoid activation		
Categorical Cross-Entropy with softmax activation		
O Mean Absolute Error with tanh activation		
In a regression task with potential outliers in the data, which loss function would be most 1 poin appropriate?	nt	
O Mean Squared Error (MSE)		
Mean Absolute Error (MAE) or Huber Loss		
O Binary Cross-Entropy		
O Categorical Cross-Entropy		

```
Algorithm Gradient Descent Algorithm
```

```
Require: Training data \{(\mathbf{x}_i, y_i)\}_{i=1}^n, loss function \mathcal{L}, learning rate \alpha, iterations \mathcal{T}
```

Ensure: Optimized parameters θ

- 1: Initialize parameters $\theta^{(0)}$ randomly
- 2: for t = 1 to T do $\theta^{(t)} = \theta^{(t-1)} - \alpha \cdot \nabla_{\theta} \mathcal{L}(\theta^{(t-1)})$
- 3: **if** Convergence criteria met **then**
- 4: break
- 5: end if
- 6: end for
- 7: return Final parameters $\theta^{(T)}$

Convergence will be very slow

• The algorithm may overshoot the minimum and potentially diverge

) The algorithm will always converge to the global minimum

Questions

Any comment ?

This content is neither created nor endorsed by Google.

Google Forms