

Задача: выделить фрагменты в тексте и протегировать. Метод классификации токенов (sequence labeling). $x_i \rightarrow y_i = \{0, 1\}^T$. Если $y_{i,j} = 1$, то i -й токен входит в фрагмент с тегом j . Обучаем модель минимизацией бинарной кросс-энтропии:

$$\mathcal{L}(y, \hat{y}) = \sum_{i,j} \text{logloss}(y_{i,j}, \hat{y}_{i,j}) \rightarrow \min$$

$$\text{logloss}(y, \hat{y}) = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$$

