

ML1 HOME ASSIGNMENT 3

id: 12179078

VOLODYMYR MEDENTSIR

11

1.1. Data likelihood without independence assumption:

$$p(T, X | \Theta) = [\text{product rule}] = p(X | T, \Theta) \cdot p(T | \Theta) = \prod_{n=1}^N p(x_n | t_n, \Theta)$$

$$\cdot p(x_n | t_n, \Theta) = \prod_{h=1}^N \prod_{k=1}^K \left(\pi_k p(x_n | t_n = c_k, \Theta) \right)^{I(t_n = c_k)}$$

Data likelihood with naive assumption:

$$p(X, T | \Theta) = \prod_{h=1}^N \prod_{k=1}^K \left(\pi_k p(x_n | t_n = c_k, \Theta_k) \right)^{I(t_n = c_k)} = [\text{naive Bayes assumption}] =$$

$$= \prod_{h=1}^N \prod_{k=1}^K \prod_{d=1}^D \left(\pi_k \cdot p(x_{nd} | t_n = c_k, \Theta_{dk}) \right)^{I(t_n = c_k)}$$

1.2. Number of parameters before NB assumption: $K \times (2^D - 1)$

(For each class there is 2^D possible combinations of the features, but the probability of the last combination could be derived from all others $\Rightarrow 2^D - 1$ parameters.)

Number of parameters after NB assumption: $K \times D$

(For each class there are D features, and each of them has its own probability to occur $\Rightarrow D$ parameters.)

This assumption is called naive because it usually does not hold in real life, but conditional independence makes training the model much easier.

Example: in text classification problem we assume that words are independent of a context which obviously does not hold in natural language.

1.3.

$$\begin{aligned} \ln p(T, X | \Theta) &= \sum_{k=1}^K \sum_{n=1}^N \sum_{d=1}^D I(t_n = c_k) \cdot (\ln \pi_k + \ln p(x_{nd} | t_n = c_k, \Theta_{dk})) = \\ &= \sum_{k=1}^K \sum_{n=1}^N \sum_{d=1}^D I(t_n = c_k) \cdot (\ln \pi_k + x_{nd} \ln \Theta_{dk} + (1-x_{nd}) \ln (1-\Theta_{dk})). \end{aligned}$$

1.4.

 $\hat{\theta}_{dk}^{MLE}$:

$$\frac{d \ln p(X, T | \Theta)}{d \theta_{dk}} = \sum_{n \in C_k} x_{nd} \frac{1}{\theta_{dk}} - \frac{(1-x_{nd})}{1-\theta_{dk}} = 0 \Rightarrow$$

$$\Rightarrow (1-\theta_{dk}) \sum_{n \in C_k} x_{nd} - \theta_{dk} \sum_{n \in C_k} (1-x_{nd}) = 0 \Rightarrow$$

$$\Rightarrow \theta_{dk} \left(- \sum_{n \in C_k} x_{nd} - \sum_{n \in C_k} (1-x_{nd}) \right) = - \sum_{n \in C_k} x_{nd} \Rightarrow$$

$$\Rightarrow \hat{\theta}_{dk} = \frac{\sum_{n \in C_k} x_{nd}}{\sum_{n \in C_k} (x_{nd} + 1-x_{nd})} = \frac{\sum_{n \in C_k} x_{nd}}{|C_k|}, \text{ with } |C_k| - \text{number of documents in class } k.$$

$$\frac{d^2 \ln p(X, T | \Theta)}{d \theta_{dk}^2} = - \frac{1}{\theta_{dk}^2} \sum_{n \in C_k} x_{nd} - \frac{1}{(1-\theta_{dk})^2} \sum_{n \in C_k} (1-x_{nd}) \Big|_{\theta_{dk} = \hat{\theta}_{dk}} =$$

$$= - \frac{(|C_k|)^2}{\sum_{n \in C_k} x_{nd}} - \frac{\sum_{n \in C_k} (1-x_{nd})}{\frac{((|C_k|) - \sum_{n \in C_k} x_{nd})^2}{|C_k|^2}} = - \frac{|C_k|^2}{\sum_{n \in C_k} x_{nd}} - \frac{|C_k|^2}{\sum_{n \in C_k} 1-x_{nd}} < 0.$$

$$\Rightarrow \hat{\theta}_{dk}^{MLE} = \frac{\sum_{n \in C_k} x_{nd}}{|C_k|}.$$

$\hat{\theta}_{dk}^{MLE}$ could be interpreted as the average number of word d per document in class k.

1.5. $p(C_1|x)$ for general κ classes naive Bayes classifier, $x \in \mathbb{R}^D$

$$p(C_1|x) = \frac{p(x|C_1) \cdot p(C_1)}{p(x)} = \frac{\pi_1 \prod_{d=1}^D p(x_d|C_1)}{\sum_{k=1}^{\kappa} \pi_k \prod_{d=1}^D p(x_d|\Theta_{dk})}$$

1.6. $p(C_1|x)$ for Bernoulli case.

$$p(C_1|x) = \frac{\pi_1 \prod_{d=1}^D \theta_{d1}^{x_d} \cdot (1-\theta_{d1})^{1-x_d}}{\sum_{k=1}^{\kappa} \pi_k \prod_{d=1}^D \theta_{dk}^{x_d} \cdot (1-\theta_{dk})^{1-x_d}}$$

1.7.

$$x \in C_1 \Leftrightarrow p(C_1|x) > p(C_K|x) \quad \forall K \neq 1.$$

$$\pi_1 \prod_{d=1}^D \theta_{d1}^{x_d} (1-\theta_{d1})^{1-x_d} > \pi_K \prod_{d=1}^D \theta_{dk}^{x_d} (1-\theta_{dk})^{1-x_d}$$

$$\prod_{d=1}^D \left(\frac{\theta_{d1}}{\theta_{dk}} \right)^{x_d} \left(\frac{1-\theta_{d1}}{1-\theta_{dk}} \right)^{1-x_d} > \frac{\pi_k}{\pi_1} \Rightarrow$$

$$\Rightarrow \sum_{d=1}^D x_d \ln \frac{\theta_{d1}}{\theta_{dk}} + (1-x_d) \ln \left(\frac{1-\theta_{d1}}{1-\theta_{dk}} \right) > \ln \frac{\pi_k}{\pi_1} \Rightarrow$$

$$\Rightarrow \sum_{d=1}^D x_d \left(\ln \frac{\theta_{d1}}{\theta_{dk}} - \ln \left(\frac{1-\theta_{d1}}{1-\theta_{dk}} \right) \right) > \ln \frac{\pi_k}{\pi_1} - \sum_{d=1}^D \ln \left(\frac{1-\theta_{d1}}{1-\theta_{dk}} \right) \Rightarrow$$

$$\Rightarrow x^T \alpha_K > c_K \quad \forall K \neq 1 \quad \text{with} \quad x^T = (x_1, \dots, x_D),$$

$$a_k = \left(\ln \left(\frac{\theta_{11}}{\theta_{01}} \cdot \frac{1-\theta_{1k}}{1-\theta_{01}} \right), \dots, \ln \left(\frac{\theta_{01}}{\theta_{1k}} \cdot \frac{1-\theta_{01}}{1-\theta_{11}} \right) \right)^T$$

$$c_k = \ln \frac{\pi_k}{\pi_1} - \sum_{d=1}^D \ln \left(\frac{1-\theta_{d1}}{1-\theta_{dk}} \right)$$

1.8. Yes, it is possible to consider feature vectors with continuous values. Then those features should be described by continuous random variables. One of the possible example is a Gaussian Bayes classifier, where $p(x|C_k) = N(x|\mu_k, \sigma^2_k)$.

In the document classification task continuous feature could represent the frequency of occurrence of the words in a class and then it could be modelled with a triangular distribution.

d2

2.0.

$$1) p(T | \varphi, w_1, \dots, w_K) = \prod_{n=1}^N \prod_{k=1}^K p(C_k | \varphi_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_k(\varphi_n)^{t_{nk}}$$

$$2) \ln p(T | \varphi, w_1, \dots, w_K) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_k(\varphi_n).$$

$$3) \frac{\partial y_k(\varphi_n)}{\partial w_j}$$

$$j \neq k : \frac{\partial y_k(\varphi_n)}{\partial w_j} = \exp(w_k^\top \varphi_n) \left(- \sum_i \exp(w_i^\top \varphi_n) \right)^{-2} \exp(w_j^\top \varphi_n) \cdot \varphi_n^\top$$

$$= - \frac{\exp(w_k^\top \varphi_n)}{\sum_i \exp(w_i^\top \varphi_n)} \cdot \frac{\exp(w_j^\top \varphi_n)}{\sum_i \exp(w_i^\top \varphi_n)} \cdot \varphi_n^\top = - y_k(\varphi_n) \cdot y_j(\varphi_n) \cdot \varphi_n^\top.$$

$$j = k \quad \frac{\partial y_k(\varphi_n)}{\partial w_j} = \frac{\exp(w_k^\top \varphi_n)}{\sum_i \exp(w_i^\top \varphi_n)} \cdot \varphi_n^\top - \frac{\exp(w_k^\top \varphi_n) \cdot \exp(w_k^\top \varphi_n)}{\left(\sum_i \exp(w_i^\top \varphi_n) \right)^2} \cdot \varphi_n^\top$$

$$= y_K(\varphi_n) (1 - y_K(\varphi_n)) \varphi_n^T$$

So in general case: $\frac{\partial y_k(\varphi_n)}{\partial w_j} = y_K(\varphi_n) (t_{kj} - y_j(\varphi_n)) \cdot \varphi_n^T$,

with $I_{kj} = \begin{cases} 1 & , k=j \\ 0 & , k \neq j \end{cases}$

$$4) \quad \nabla_{w_j} \text{ Lee } p(T | \varphi, w_1, \dots, w_K) = \left[y_{nk} := y_K(\varphi_n) \right] = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \frac{\partial \ln y_{nk}}{\partial w_j} =$$

$$= \sum_{n=1}^N \sum_{k=1}^K \frac{t_{nk}}{y_{nk}} \frac{\partial y_{nk}}{\partial w_j} = \sum_{n=1}^N \sum_{k=1}^K \frac{t_{nk}}{y_{nk}} \cdot y_{nk} (I_{kj} - y_j) \varphi_n^T =$$

$$= \sum_{n=1}^N \varphi_n^T \cdot \left(\sum_{k=1}^K t_{nk} \cdot I_{kj} - y_j \underbrace{\sum_{k=1}^K t_{nk}}_1 \right) =$$

$$= \sum_{n=1}^N \varphi_n^T \cdot (t_{nj} - y_j)$$

$$5) \quad \nabla_{w_j} \text{ Lee } p(T | \varphi, w_1, \dots, w_K) = d_j^T \varphi$$

with $d_j = (t_{1j} - y_{1j}, \dots, t_{nj} - y_{nj})^T$

(if $W \in \mathbb{R}^{M \times K}$ -weight matrix, then we can also consider:

$$\nabla_w \text{ Lee } p(T | \varphi, w_1, \dots, w_K) = (T - Y)^T \varphi \in \mathbb{R}^{K \times M}$$

It is useful because matrix multiplications are more computationally efficient.

2.2.

$$E = -\ln p(T|q, w_1, \dots, w_K) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_k(q_n).$$

E - is a cross-entropy.

The relationship between E and log-likelihood: Instead of maximizing log-likelihood we need to minimize E .

What changes in terms of weight update?

Because we minimize E and maximize log-likelihood:

$$w^{z+1} = w^z + \eta \cdot \triangledown_w (\text{log-likelihood})^\top, \text{ but}$$

$$w^{z+1} = w^z - \eta \triangledown_w E^\top$$

2.3.

B - the batch size, n_B - number of batches.

1. Initialize: η - learning rate, w_1, \dots, w_K - weights

2. For $z = 1$ to n_B

i. Randomly sample a minibatch M of size B .

$$\text{ii. } w_j = w_j - \eta \left(\triangledown_{w_j} \sum_{i \in M} E_i \right)^\top \text{ if } j = 1, \dots, K$$

$$\text{with } E_i = -\sum_{k=1}^K t_{ik} \ln y_k(q_i)$$

iii. (optional) decrease η .

3. Return w_1, \dots, w_K .

Advantage in comparison with SGD with single data point: the estimate of the gradient is less noisy.

Comparison with full batch gradient descent:

Though full batch provide a more accurate estimate of the gradient, it will be more computationally efficient to use mini-batch GD. Some argue that small batches can also provide a regularizing effect.

2.4. Multiclass logistic regression:

$$y_j = \frac{\exp(\tilde{w}_j^T q(x))}{\sum_{i=1}^3 \exp(\tilde{w}_i^T q(x))}$$

$$\tilde{w}_1 = (w_5 \ w_6)^T, \tilde{w}_2 = (w_7 \ w_8)^T, \tilde{w}_3 = (w_9 \ w_{10})^T$$

The activation function on the last layer should be softmax.

$$q(x) = (f_1(w_1x_1 + w_2x_2) \ f_2(w_3x_1 + w_4x_2))^T$$

Because $q(x)$ - fixed basis function $\Rightarrow w_1, w_2, w_3, w_4$ - should be fixed and non-zero, but activation functions f_1, f_2 - could be of any choice.

(Note :

However if we need a bias component $q_0(x_1, x_2) = 1$, then as an activation function f_1 we can use $f_1(x_1, x_2) \equiv 1$.

$$2.5. \bar{y} = \text{Softmax}(W_2 \cdot \text{ReLU}(W_1 \bar{x}))$$

$$1. \text{ Forward step. } \bar{x} = (0,3 \ 0,7)^T \quad y = (0 \ 0 \ 1)^T$$

$$W_1 \bar{x} = \begin{pmatrix} 0,4 & 0,87 \\ 0,58 & 0,34 \end{pmatrix} \begin{pmatrix} 0,3 \\ 0,7 \end{pmatrix} = \begin{pmatrix} 729/1000 \\ 103/250 \end{pmatrix} = \begin{pmatrix} 0,729 \\ 0,412 \end{pmatrix}$$

$$\text{ReLU} \begin{pmatrix} 729/1000 \\ 103/250 \end{pmatrix} = \begin{pmatrix} 0,729 \\ 0,412 \end{pmatrix}$$

$$W_2 \text{ReLU}(W_1 \bar{x}) = \begin{pmatrix} 0,12 & 0,87 \\ 0,82 & 0,31 \\ 0,34 & 0,9 \end{pmatrix} \begin{pmatrix} 0,729 \\ 0,412 \end{pmatrix} = \begin{pmatrix} 2787/6250 \\ 1451/2000 \\ 93213/100000 \end{pmatrix} \approx$$

$$\approx \begin{pmatrix} 0,446 \\ 0,726 \\ 0,932 \end{pmatrix}$$

$$\text{Softmax}(W_2 \text{ReLU}(W_1 x)) = \begin{bmatrix} e^{0,446} + e^{0,726} + e^{0,932} \approx \\ \approx 6,168 \end{bmatrix} =$$

$$= \begin{pmatrix} e^{0,446} / 6,168 \\ e^{0,726} / 6,168 \\ e^{0,932} / 6,168 \end{pmatrix} \approx \begin{pmatrix} 0,253 \\ 0,335 \\ 0,412 \end{pmatrix}$$

$$\text{The loss: } E = -\ln 0,412 \approx 0,887$$

2) Compute the derivative:

$$\frac{dE}{dw_5} = \frac{d}{dw_5} (- (t_1 \ln y_1(p) + t_2 \ln y_2(p) + t_3 \ln y_3(p))) =$$

$$= \frac{d}{dw_5} \left(- \left(t_1 \cdot (w_5 p_1 + w_6 p_2) + t_2 (w_7 p_1 + w_8 p_2) + t_3 (w_9 p_1 + w_{10} p_2) - \right. \right.$$

$$\left. \left. - \sum_{i=1}^3 t_i \underbrace{\ln (\exp(w_5 p_1 + w_6 p_2) + \exp(w_7 p_1 + w_8 p_2) + \exp(w_9 p_1 + w_{10} p_2))}_{n_1} \right) \right) =$$

$$= -t_1 p_1 + \frac{\exp(w_5 p_1 + w_6 p_2) \cdot p_1}{\exp(w_5 p_1 + w_6 p_2) + \exp(w_7 p_1 + w_8 p_2) + \exp(w_9 p_1 + w_{10} p_2)},$$

$$\text{with } p_1 = \text{ReLU}(w_1 x_1 + w_2 x_2) = 0,729, p_2 = \text{ReLU}(w_3 x_1 + w_4 x_2) = 0,412.$$

$$\frac{dE}{dw_5} = 0,729 \cdot 0,253 \approx 0,184.$$

$$w_i = w_i - 0,05 \cdot \frac{\partial E}{\partial w_i}$$

$$w_1 = 0,4 + 0,05 \cdot 0,044 \approx 0,402$$

$$w_2 = 0,87 + 0,05 \cdot 0,103 \approx 0,875$$

$$w_3 = 0,58 + 0,05 \cdot 0,062 \approx 0,583$$

$$w_4 = 0,34 + 0,05 \cdot 0,144 \approx 0,347$$

$$w_5 = 0,12 - 0,05 \cdot 0,184 \approx 0,111$$

$$w_6 = 0,87 - 0,05 \cdot 0,104 \approx 0,865$$

$$w_7 = 0,82 - 0,05 \cdot 0,244 \approx 0,808$$

$$w_8 = 0,31 - 0,05 \cdot 0,138 \approx 0,303$$

$$w_9 = 0,77 + 0,05 \cdot 0,429 \approx 0,791$$

$$w_{10} = 0,9 + 0,05 \cdot 0,242 \approx 0,912$$

9) $W_1 x = \begin{pmatrix} 0,402 & 0,875 \\ 0,583 & 0,347 \end{pmatrix} \begin{pmatrix} 0,3 \\ 0,7 \end{pmatrix} = \begin{pmatrix} 733/10000 \\ 2089/5000 \end{pmatrix} \approx \begin{pmatrix} 0,733 \\ 0,418 \end{pmatrix}$

$$\text{ReLU}(W_1 x) = \begin{pmatrix} 0,733 \\ 0,418 \end{pmatrix}$$

$$W_2 \text{ReLU}(W_1 x) = \begin{pmatrix} 0,111 & 0,865 \\ 0,808 & 0,303 \\ 0,791 & 0,912 \end{pmatrix} \begin{pmatrix} 0,733 \\ 0,418 \end{pmatrix} \approx \begin{pmatrix} 0,443 \\ 0,719 \\ 0,961 \end{pmatrix}$$

$$\text{Softmax}(W_2 \text{ReLU}(W_1 x)) = \left[e^{0,443} + e^{0,719} + e^{0,961} \approx 6,224 \right] =$$

$$= \begin{pmatrix} e^{0,443}/6,224 \\ e^{0,719}/6,224 \\ e^{0,961}/6,224 \end{pmatrix} \approx \begin{pmatrix} 0,25 \\ 0,33 \\ 0,42 \end{pmatrix}$$

Loss: $E = -\ln(0,42) \approx 0,868 \Rightarrow \text{the loss decreased.}$