

STATISTICAL MACHINE LEARNING (WS2018)
SEMINAR 4

Assignment 1. a) Recall the negative log likelihood for the linear regression:

$$\mathcal{L}(\mathbf{w}) = \frac{m}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \quad (1)$$

where:

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1n} \\ 1 & \vdots & \ddots & \vdots \\ 1 & x_{m1} & \dots & x_{mn} \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_m^T \end{pmatrix} \quad (2)$$

is the input matrix (input dimension n , the number of data samples m),

$$\mathbf{y} = (y_1, \dots, y_m)^T, \quad y_i \in \mathbb{R}^m \quad (3)$$

the target vector and

$$\mathbf{w} = (w_0, \dots, w_n)^T, \quad w_i \in \mathbb{R}^m \quad (4)$$

is the weight vector. Prove that the MLE solution is:

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (5)$$

b) When L2 regularization is involved we get the following loss function:

$$\mathcal{L}(\mathbf{w}) = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w}, \quad \lambda > 0 \quad (6)$$

Prove that the MLE solution is now:

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (7)$$

Hint: use matrix calculus identities for differentiation to simplify your computations.

Assignment 2. a) Give the backward message for the multinomial cross entropy layer for a single dataset sample:

$$l(\mathbf{p}, \mathbf{t}) = \sum_{j=1}^K t_j \log(p_j) \quad (8)$$

where $\mathbf{p} = (p_1, \dots, p_K)^T$ is an input vector and \mathbf{t} a one-hot encoded target vector (for the target class $k \in \{1, \dots, K\}$ we have $t_k = 1$ and $t_l = 0$ for $l \neq k$).

b) Give the backward message for the softmax layer:

$$p_j(\mathbf{s}) = \frac{e^{s_j}}{\sum_{i=1}^K e^{s_i}}, \quad j \in \{1, \dots, K\} \quad (9)$$

c) Give the backward message for the linear layer:

$$s_j(\mathbf{x}, \mathbf{W}) = \sum_{i=1}^N w_{ij} x_i, \quad j \in \{1, \dots, K\} \quad (10)$$

where $\mathbf{x} = (x_1, \dots, x_N)^T$ is an input vector and w_{ij} a weight of the connection between input i and output j . Give also the parameter message $\frac{\partial s_j}{\partial w_{ij}}$.

d) Combine all layers to form the logistic regression model:

$$l(\mathbf{p}(\mathbf{s}(\mathbf{x}, \mathbf{W})), \mathbf{t}) \quad (11)$$

and use the previous results to get its gradient with respect to the parameters:

$$\frac{\partial l}{\partial w_{ij}}, \quad i \in \{1, \dots, N\}, j \in \{1, \dots, K\} \quad (12)$$

Assignment 3. Consider the following neural network:

$$y_k(\mathbf{x}, \mathbf{W}^{(1)}, \mathbf{W}^{(2)}) = \sum_{j=1}^M w_{jk}^{(2)} \tanh\left(\sum_{i=1}^N w_{ij}^{(1)} x_i\right), \quad (13)$$

having outputs denoted $k \in \{1, \dots, K\}$. Draw the network. Swapping positions of any two neurons in the first linear layer will not change the output of the network. Similarly changing sign of all input and output connection weights of a neuron in the same layer will not change the network's output (because $\tanh(-x) = -\tanh(x)$). Such configurations are called the *weight-space symmetries*. Compute the maximum number of networks preserving output for the same input based on the defined weight-space symmetries.

Assignment 4. Prove that for two independent random variables X and Y :

$$\text{Var}(XY) = [\mathbb{E}(X)]^2 \text{Var}(Y) + [\mathbb{E}(Y)]^2 \text{Var}(X) + \text{Var}(X)\text{Var}(Y). \quad (14)$$

We used this identity in the derivation of the Xavier initialization for a linear neuron.

Hint: use $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$ for X and Y independent.

Assignment 5. Define backward messages for the following layers. Note that all outputs of these layers are computed independently so we can omit subscripts, i.e., $y(x) = y_j(x)$ for each layer output j .

a) Sigmoid:

$$y(x) = \frac{1}{1 + e^{-x}}. \quad (15)$$

Express the backward message $y'(x)$ using the forward message $y(x)$.

b) Hyperbolic tangent:

$$y(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}. \quad (16)$$

Express the backward message $y'(x)$ using the forward message $y(x)$.

c) Rectified Linear Unit (ReLU):

$$y(x) = \max(0, x). \quad (17)$$

d) Parametric ReLU (PReLU):

$$y(x, a) = \max(0, x) + a \min(0, x). \quad (18)$$

Derive also the parameter message $\frac{\partial y(x, a)}{\partial a}$.