

Homework #2: Machine Learning

Kyler Little

February 3, 2018

Problem #1

(a) More generally, if we are learning from ± 1 data to predict a noisy target $P(y|\mathbf{x})$ with candidate hypothesis h , show that the maximum likelihood method reduces to the task of finding h that minimizes:

$$E_{\text{in}}(\mathbf{w}) = \sum_{n=1}^N \mathbb{I}[y_n = +1] \ln\left(\frac{1}{h(\mathbf{x}_n)}\right) + \mathbb{I}[y_n = -1] \ln \frac{1}{1 - h(\mathbf{x}_n)}$$

Formally, we are trying to learn the target function:

$$f(\mathbf{x}) = P[y = +1|\mathbf{x}]$$

In non-mathematical terms, this target function represents the probability of getting $y_n = +1$ from the data (\mathbf{x}_n) . The function f is generated by a noisy target function, so the data doesn't give us the value of f explicitly. Instead, the data is generated by a noisy target function $P(y|\mathbf{x})$. Our goal is to minimize the error between the classifications of actual training data and the predicted classifications from our hypothesis model h . We can do this by talking about the notion of likelihood. The likelihood that the target distribution $P(y|\mathbf{x})$ is captured by our hypothesis $h(\mathbf{x})$ is:

$$P(y|\mathbf{x}) = \begin{cases} h(\mathbf{x}) & \text{for } y = +1; \\ 1 - h(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

If we assume that the data points were independently generated (a very fair assumption to make), then we can express the probability of getting all of

the y_n 's in the data set from their corresponding x_n 's as:

$$\prod_{n=1}^N P(y_n | \mathbf{x}_n)$$

The method of maximum likelihood method would select the hypothesis h which maximizes the probability above. We can actually express this maximization problem as a minimization problem by using the properties of 'ln'. The natural logarithm is a monotonically increasing function, and so taking the natural logarithm of the probability above will not affect the maximization problem. By the same token, ' $-\ln$ ' is a monotonically decreasing function, so we can take the natural logarithm of the probability above and convert the maximization problem to a minimization problem without affecting global minima of the probability distribution. Using a simple logarithm rule ($-\ln x = \ln x^{-1}$), we find:

$$\prod_{n=1}^N P(y_n | \mathbf{x}_n) = \sum_{n=1}^N \ln \frac{1}{P(y_n | \mathbf{x}_n)}$$

Substituting in what $P(y|\mathbf{x})$ actually equals, we arrive at what we wanted. We now find h that minimizes:

$$E_{\text{in}}(\mathbf{w}) = \sum_{n=1}^N [\![y_n = +1]\!] \ln \frac{1}{h(\mathbf{x}_n)} + [\![y_n = -1]\!] \ln \frac{1}{1 - h(\mathbf{x}_n)}$$

(b) For the case $h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$, argue that minimizing the in-sample error in part (a) is equivalent to minimizing the one in (3.9).

Substitute in $h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x}) = \frac{e^{w^T \mathbf{x}}}{1 + e^{w^T \mathbf{x}}}$ to the in-sample error from part (a). This yields:

$$E_{\text{in}}(\mathbf{w}) = \sum_{n=1}^N [\![y_n = +1]\!] \ln \frac{1 + e^{w^T \mathbf{x}}}{e^{w^T \mathbf{x}}} + [\![y_n = -1]\!] \ln(1 + e^{w^T \mathbf{x}})$$

Next, use the rule $b \log(x) = \log(x^b)$ to bring the y_n 's inside the logarithms and $\log a + \log b = \log(ab)$ to reduce the two expression to a single expression:

$$E_{\text{in}}(\mathbf{w}) = \sum_{n=1}^N \ln(1 + e^{-y_n w^T \mathbf{x}})$$

Multiplying by a constant factor of $1/N$ will not change the minimization problem. Thus, the argument is complete.

Side Note (from LFD):

For two probabilistic distributions $p, 1 - p$ and $q, 1 - q$ with binary outcomes, the cross-entropy (from information theory) is:

$$p \log \frac{1}{q} + (1 - p) \log \frac{1}{1 - q}$$

The in-sample error in part (a) corresponds to a cross-entropy error measure on the data point (\mathbf{x}_n, y_n) , with $p = \llbracket y_n = +1 \rrbracket$ and $q = h(\mathbf{x}_n)$.

Problem #2

Recall the objective function for linear regression can be expressed as

$$E(w) = \frac{1}{N} \|Xw - y\|^2$$

as in Equation (3.3) of LFD. Minimizing this function with respect to w leads to the optimal w as $(XTX)^{-1}X^Ty$. This solution holds only when X^TX is nonsingular. To overcome this problem, the following objective function is commonly minimized instead:

$$E_2(w) = \|Xw - y\|^2 + \lambda \|w\|^2$$

where $\lambda > 0$ is a user-specified parameter. Please do the following:

(a) Derive the optimal w that minimize $E_2(w)$.

This is fairly straightforward. We differentiate $E_2(w)$, set the result to zero, and solve for w . The reason we can do this is because $E_2(w)$ is a convex function. This means its derivative is a monotonically increasing function, signifying that there is a unique global minimum. Before differentiating, it's necessary to reduce $E_2(w)$ to an expression that is more easily differentiated. To start, we recall that $\|x\| = \sqrt{x^Tx}$. Applying this yields:

$$E_2(w) = \frac{1}{N} ((Xw - y)^T(Xw - y) + \lambda w^Tw)$$

Expanding this out yields:

$$E_2(w) = \frac{1}{N} (w^TX^TXw - w^TX^Ty - y^TXw + y^Ty + \lambda w^Tw)$$

Note that $w^T X^T y = (y^T X w)^T$, so we may combine the middle two terms, since they are identical real numbers after evaluation, resulting in $-w^T X^T y - y^T X w = -2w^T X^T y$. Now, we simply must take the derivative of $E_2(w)$, set it to 0, and solve for w . For the first term, we use the fact that $\frac{\partial(w^T A w)}{\partial w} = (A + A^T)w$. Since A is symmetric, we have that $(A + A^T)w = 2Aw$. In our situation, A is replaced by $X^T X$, so $\frac{\partial(w^T X^T X w)}{\partial w} = 2X^T X w$. For the second term, we know that $\frac{\partial(2w^T X^T y)}{\partial w} = 2X^T y$, as discussed in class. The third term disappears since it is constant with respect to w . The last term seems tricky at first, but an easy solution is illuminated by writing $w^T w$ as $w_1^2 + \dots + w_n^2$. Thus, $\partial \lambda w^T w / \partial w_i = 2\lambda w_i$. This holds $\forall i \in 1, \dots, N$, so $\partial \lambda w^T w / \partial w = 2\lambda w$. Thus, we have

$$\nabla E_2(w) = \frac{1}{N}(2X^T X w - 2X^T y + 2\lambda w)$$

Now, we set $\nabla E_2(w)$ to zero and isolate w .

$$\begin{aligned} 0 &= \frac{1}{N}(2X^T X w - 2X^T y + 2\lambda w) \\ 0 &= X^T X w - X^T y + \lambda w \\ X^T y &= X^T X w + \lambda w \\ X^T y &= (X^T X + \lambda I)w \\ (X^T X + \lambda I)^{-1} X^T y &= (X^T X + \lambda I)^{-1} (X^T X + \lambda I)w \\ w &= (X^T X + \lambda I)^{-1} X^T y \end{aligned}$$

(b) Explain how this new objective function can overcome the singularity problem of $X^T X$.

Now, just prove it's nonsingular by showing it's positive definite.

Problem #3

In logistic regression, the objective function can be written as:

$$E(w) = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n w^T x_n})$$

- (a) (10 points) Compute the first-order derivative $\nabla E(w)$. You will need to provide the intermediate steps of derivation.
- (b) (10 points) Once the optimal w is obtained, it will be used to make predictions as follows:

$$\text{Predicted class of } x = \begin{cases} 1 & \text{if } \theta(w^T x) \geq 0.5 \\ -1 & \text{if } \theta(w^T x) < 0.5 \end{cases}$$

where the function $\theta(z) = \frac{1}{1+e^{-z}}$ looks like

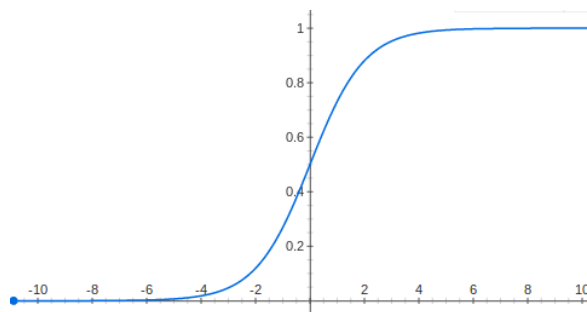


Figure 1: A basic sigmoid curve

Explain why the decision boundary of logistic regression is still linear, though the linear signal $w^T x$ is passed through a nonlinear function θ to compute the outcome of prediction. Is the decision boundary still linear if the prediction rule is changed to the following? Justify briefly.

$$\text{Predicted class of } x = \begin{cases} 1 & \text{if } \theta(w^T x) \geq 0.9 \\ -1 & \text{if } \theta(w^T x) < 0.9 \end{cases}$$

In light of your answers to the above two questions, what is the essential property of logistic regression that results in the linear decision boundary?

Problem #4

As your final deliverable to a customer, would you use the linear model with or without the 3rd order polynomial transform? Briefly explain your reasoning.