

Homework #2: Machine Learning

Kyler Little

February 4, 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 $\partial(\lambda w^T w)/\partial w = 2\lambda w$. As a result, 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 $X^T X$.

We are not necessarily guaranteed that $X^T X$ is invertible. However, we will show that $(X^T X + \lambda I)$ is always invertible. We can do this by demonstrating that it's nonsingular. In this proof, we will show that $(X^T X + \lambda I)$ is a positive definite matrix, meaning that it's nonsingular.

First, write $(X^T X + \lambda I)$ as $A^T A$. $A^T A$ is still symmetric because $X^T X$ is symmetric, and subtracting λI won't affect A 's symmetry because it only affects the diagonal of $X^T X$. Now, in order for some matrix A to be positive

definite, we must have that $x^T Ax > 0, \forall x \neq 0$. Let's now prove this.

$$\begin{aligned} x^T A^T Ax &> 0 \\ (Ax)^T Ax &> 0 \\ \|Ax\|^2 &> 0 \end{aligned}$$

Why is $\|Ax\|^2$ strictly greater than zero? It's obvious that $\|Ax\|^2$ is greater than or equal to zero because the two-norm is nonnegative. However, we also know that Ax cannot be equal to zero because subtracting the constant λ from the diagonal of $X^T X$ gives us the liberty to always choose λ such that Ax is nonzero.

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.

Using the chain rule, we first note that $\frac{d(\ln x)}{dx} = \frac{1}{x}$. Thus, we automatically have part of the derivative as

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

Next, we differentiate the argument of the natural logarithm.

$$\frac{d(1 + e^{-y_n w^T x_n})}{dw} = -y_n x_n e^{-y_n w^T x_n}$$

This is easy to derive. If we write $w^T x$ as $w_1 x_1 + \dots + w_n x_n$, then it becomes obvious why $\partial(w^T x)/\partial w = x$. In the end, we are left with:

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

(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

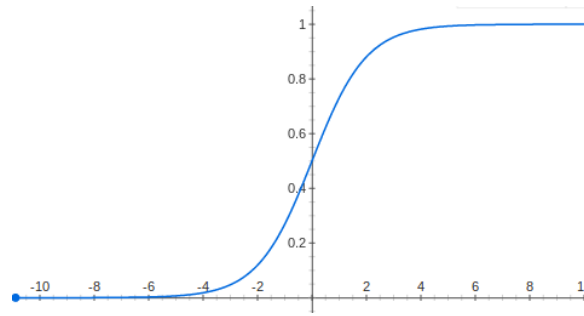


Figure 1: A basic sigmoid curve

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.

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

The reason why the decision boundary of logistic regression is still linear is because the nonlinear function that processes the linear signal $w^T x$ is a monotonically increasing function. In other words, the output can be linearly separated because monotonically increasing/decreasing functions will never return to a value in their codomains after they have left it. Thus, even though the linear signal is passed through a nonlinear function θ , the output of this nonlinear function is still linearly separable.

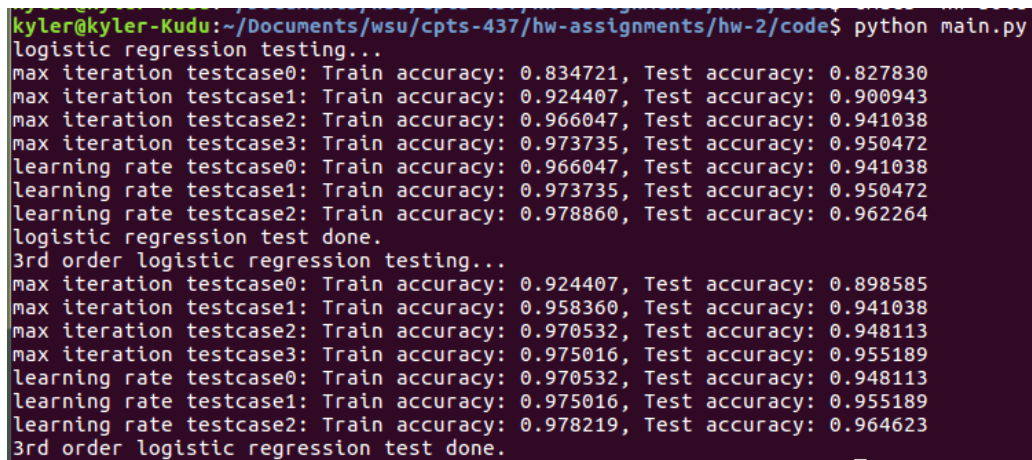
The new decision boundary listed above is still linear because we are merely translating the boundary (which is already linear).

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? The functional transformation must be monotonically increasing or decreasing. That is what allows logistic regression to retain a 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.

Based on my results, I would use the linear model without the 3rd order polynomial transform. The results I obtained with the 3rd order polynomial transform were not significantly better than those without the transform, as you can see in the figure below. Although the 3rd order polynomial trans-



```
kyler@kyler-Kudu:~/Documents/wsu/cpts-437/hw-assignments/hw-2/code$ python main.py
logistic regression testing...
max iteration testcase0: Train accuracy: 0.834721, Test accuracy: 0.827830
max iteration testcase1: Train accuracy: 0.924407, Test accuracy: 0.900943
max iteration testcase2: Train accuracy: 0.966047, Test accuracy: 0.941038
max iteration testcase3: Train accuracy: 0.973735, Test accuracy: 0.950472
learning rate testcase0: Train accuracy: 0.966047, Test accuracy: 0.941038
learning rate testcase1: Train accuracy: 0.973735, Test accuracy: 0.950472
learning rate testcase2: Train accuracy: 0.978860, Test accuracy: 0.962264
logistic regression test done.
3rd order logistic regression testing...
max iteration testcase0: Train accuracy: 0.924407, Test accuracy: 0.898585
max iteration testcase1: Train accuracy: 0.958360, Test accuracy: 0.941038
max iteration testcase2: Train accuracy: 0.970532, Test accuracy: 0.948113
max iteration testcase3: Train accuracy: 0.975016, Test accuracy: 0.955189
learning rate testcase0: Train accuracy: 0.970532, Test accuracy: 0.948113
learning rate testcase1: Train accuracy: 0.975016, Test accuracy: 0.955189
learning rate testcase2: Train accuracy: 0.978219, Test accuracy: 0.964623
3rd order logistic regression test done.
```

Figure 2: Logistic Regression Console Output

form seems to do much better than the unmodified linear model when the maximum iterations are low, they achieve roughly equivalent results when the maximum number of iterations are higher. Furthermore, they achieve essentially the same accuracy ratings for the same learning rates. In my opinion, the very marginal improvement in the model is not worth it due to

the sacrifice in performance. The 3rd order polynomial transform adds much unnecessary overhead, and it also causes poor generalization.