

Lecture 4 → Classification

- Supervised learning → learn to predict y from x given data x in feature space and labels y
- In classification, these labels are discrete-valued
- Classification → given an input vector x , assign it to one of K distinct classes C_k where $k=1, 2, \dots, K$

If $K=2$,

- $y=1$ means x is in C_1

- $y=0$ means x is in C_2

If $K>2$,

- use 1-of- K coding

- $y=(0,1,0,0)^T$ means x is in C_2

0-1 loss → Classification error = $\frac{1}{N} \sum_{j=1}^N \mathbb{1}[h(x_{\text{test}}^{(j)}) \neq y_{\text{test}}^{(j)}]$

Logistic Regression

- Decision boundary is modeled as a function of the input x
- Learn $P(C_k | x)$ over the data
- Directly predict class labels from inputs
- Logistic regression models the class probability (posterior) using a sigmoid function applied to a linear function of the feature vectors

$$p(C_1 | \theta) = h(\theta) = \sigma(w^T \phi(x))$$

Logistic sigmoid function → $\sigma(a) = \frac{1}{1 + \exp(-a)}$

Inverse of logistic sigmoid is logit function → $a = \ln\left(\frac{\sigma}{1-\sigma}\right)$ aka. odds ratio

Generalization is softmax → $p_i = \frac{\exp(\theta_i)}{\sum_j \exp(\theta_j)}$

We can define the likelihood of x as,

$$P(y=1 | x, w) = \sigma(w^T \phi(x))$$

$$P(y=0 | x, w) = 1 - \sigma(w^T \phi(x))$$

$$P(y | x | w) = \sigma(w^T \phi(x))^y (1 - \sigma(w^T \phi(x)))^{1-y}$$

The likelihood function is $p(y | w) = \prod_{n=1}^N (h^{(n)})^{y^{(n)}} (1 - h^{(n)})^{1-y^{(n)}}$

where $h^{(n)} = p(C_1 | \phi(x^{(n)})) = \sigma(w^T \phi(x^{(n)}))$

Define our loss function as $E[w] = -\log p(y | w)$, minimizing $E[w]$ will maximize the likelihood

MLE Derivation

$$\nabla_{\omega} \log P(y | x^{(1)}, \dots, x^{(N)}, \omega)$$

$$= \sum_{n=1}^N \nabla_{\omega} (y^{(n)} \log h(x^{(n)}, \omega) + (1 - y^{(n)}) \log (1 - h(x^{(n)}, \omega)))$$

$$= \sum_{n=1}^N \nabla_{\omega} y^{(n)} \left(\frac{\sigma^{(n)}(1 - \sigma^{(n)})}{\sigma^{(n)}} - (1 - y^{(n)}) \frac{\sigma^{(n)}(1 - \sigma^{(n)})}{1 - \sigma^{(n)}} \right) \nabla_{\omega} (\omega^T \phi(x^{(n)}))$$

$$= \sum_{n=1}^N (y^{(n)}(1 - \sigma^{(n)}) - (1 - y^{(n)})\sigma^{(n)}) \nabla_{\omega} (\omega^T \phi(x^{(n)}))$$

$$= \sum_{n=1}^N (y^{(n)} - \sigma^{(n)}) \phi(x^{(n)})$$

So we have $\nabla E[\omega] = \sum_{n=1}^N (h^{(n)} - y^{(n)}) \phi(x^{(n)})$, where $h^{(n)} = p(L=1 | \phi(x^{(n)}) = \sigma(\omega^T \phi(x^{(n)})))$

This is very similar to the gradient descent expression for linear regression, with the addition of the sigmoid function $\sigma(a)$

Newton's Method

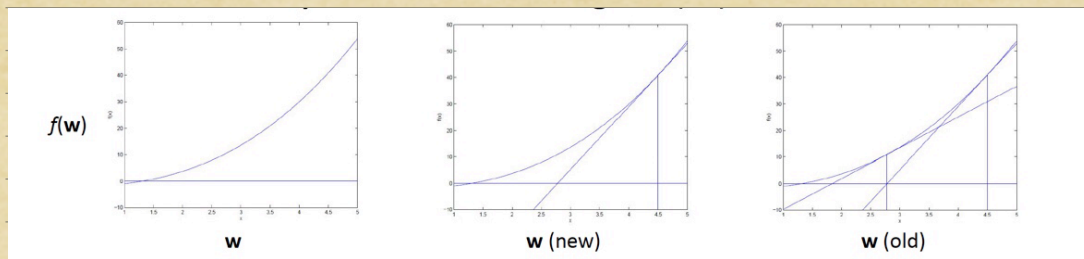
Goal is to minimize a general function $E[\omega]$

Solve for $f(\omega) = \frac{\partial E[\omega]}{\partial \omega} = 0$

We do so by repeating the following until convergence,

$$\omega := \omega - \frac{f(\omega)}{f'(\omega)}$$

$\xrightarrow{\text{Current value}}$ $\xrightarrow{\text{Slope}}$



In the multivariate case, we have $\omega := \omega - H^{-1} \nabla_{\omega} E$, where H is the Hessian matrix evaluated as $H_{ij}(\omega) = \frac{\partial^2 E[\omega]}{\partial \omega_i \partial \omega_j}$

LR closed form solution $\rightarrow \omega_{\text{ML}} = (\Phi^T \Phi)^{-1} \Phi^T y$

WLR closed form solution $\rightarrow \omega_{\text{WLS}} = (\Phi^T R \Phi)^{-1} \Phi^T R y$ R is an $N \times N$ diagonal weight matrix

In logistic regression, $h(x, \omega)$ is non linear and there is no closed form solution

- we can only estimate by repeatedly applying Newton steps

Iteratively applying NR involves least-squares with weights matrix R , $R_{nn} = h^{(n)}(1 - h^{(n)})$

R depends on ω , ω depends on R , so we have iterative reweighted least squares (IRLS)

$$\omega^{(\text{new})} = (\Phi^T R \Phi)^{-1} \Phi^T R z, \quad z = \Phi \omega^{(\text{old})} - R^{-1} (h - y)$$

K-Nearest Neighbor Classification

- Given a test example x , find the k training examples that are closest to x
- Predict as a label for the testing sample the most frequent class among all y 's from $KNN(x)$

$$KNN(x) = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(k)}, y^{(k)})\}$$

$$h(x) = \arg \max_y \sum_{(x', y') \in KNN(x)} 1[y' = y] \quad \text{majority vote}$$

- Larger k leads to smoother decision boundary
- Classification generally improves as N increases
- As $N \rightarrow \infty$, error rate of 1-NN is never more than twice the optimal error, which is obtained from true conditional class dist's

- Hyperparameters \rightarrow ① Distance metric $\rightarrow D(x, x')$
② Value of k

- Can turn kNN into a regression problem by assigning the label of the testing example as the average of the the nearest target labels,

$$h(x) = \frac{1}{k} \sum_{(x', y') \in KNN(x)} y'$$

- kNN advantages \rightarrow simple + flexible
effective for low dimensional outputs
- kNN disadvantages \rightarrow expensive as we have to compute + store distances for the whole training data for every prediction
all points are far away in high dimensions
not robust to irrelevant features/outliers