

Lecture 14, Mar 8, 2024

Logistic Regression

- Try to estimate $P[Y = i | \mathbf{x}, \mathbf{w}]$ where $\mathbf{y} \in C = \{1, 2, \dots, c\}$ are classes, \mathbf{x} is a feature, and \mathbf{w} are linear model weights
- Example: binary hypothesis ($c = 2$) with Bernoulli probabilities
 - Then $p(y = 1 | \mathbf{x}) = \frac{p(x|y=1)p(y=1)}{p(x|y=1)p(y=1) + p(x|y=0)p(y=0)} = \frac{1}{1 + \frac{p(x|y=0)p(y=0)}{p(x|y=1)p(y=1)}}$
 - We can write this as $\frac{1}{1 + e^{-\alpha}}$ where $\alpha = \log \frac{p(x|y=0)p(y=0)}{p(x|y=1)p(y=1)}$
- $\sigma(\alpha) = \frac{1}{1 + e^{-\alpha}}$ is the *sigmoid function*, which maps $\mathbb{R} \rightarrow (0, 1)$ which is useful for probabilities
 - Has an S shape with value of $\frac{1}{2}$ at 0
 - Note $\sigma(-\alpha) = 1 - \sigma(\alpha)$ and $\alpha = \log \frac{\sigma(\alpha)}{1 - \sigma(\alpha)}$
 - $\frac{d\sigma}{d\alpha} = \sigma(\alpha)(1 - \sigma(\alpha))$
 - We can classify $\hat{y} = 1$ if $\sigma(\alpha) > \frac{1}{2}$ or $\hat{y} = 0$ otherwise
- Our model is then $\hat{p}(y = 1 | \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = \sigma(\mathbf{w}^T \mathbf{x})$, where we try to find the best weights \mathbf{w}
- Compared to Gaussian discriminant analysis, which has $2D$ for means and $D(D + 1)/2$ for covariances and priors, we only have D parameters and a lot less computation overall
- Consider a Bernoulli trial with $\theta = P[y = 1]$, so $P[y] = \theta^y(1 - \theta)^{1-y}$
 - Let $\theta = P[y = 1 | \mathbf{x}, \mathbf{w}] = \sigma(\mathbf{w}^T \mathbf{x})$
 - For n trials, the NLL is $-\log \prod_{i=1}^n P[y_i | \mathbf{x}_i, \mathbf{w}] = -\sum_{i=1}^n \log(\theta_i^{y_i}(1 - \theta_i)^{1-y_i}) = -\sum_{i=1}^n y_i \log \theta_i + (1 - y_i) \log(1 - \theta_i)$ where $\theta_i = \sigma(\mathbf{w}^T \mathbf{x}_i) = \sigma_i$
 - This is the *cross-entropy* loss function
 - $\frac{d}{d\mathbf{w}} \text{NLL} = -\sum_{i=1}^n \left(y_i \frac{1}{\theta_i} \theta'_i + (1 - y_i) \frac{1}{1 - \theta_i} (-\theta'_i) \right) = -\sum_{i=1}^n \left(y_i \frac{\theta'_i}{\theta_i} - (1 - y_i) \frac{\theta'_i}{1 - \theta_i} \right) = 0$
 - $\theta'_i = \sigma_i(1 - \sigma_i) \frac{d}{dw_j} \mathbf{w}^T \mathbf{x}_i = \sigma_i(1 - \sigma_i) x_{ij}$
 - $\frac{\theta'_i}{\theta_i} = (1 - \sigma_i) x_{ij} \implies \frac{\theta'_i}{1 - \theta'_i} = \sigma_i x_{ij}$
 - Therefore $\frac{d}{d\mathbf{w}} \text{NLL} = \sum_{i=1}^n (y_i - \sigma_i) x_{ij} = 0$
 - * This can be interpret as the error multiplied by the observation
 - No closed-form solution; we can use methods such as gradient descent
 - * The gradient vector is $\sum_{i=1}^n (y_i - \sigma_i) \mathbf{x}_i$
- Just like in linear regression, we're not restricted to just a single basis; we can change to e.g. a polynomial basis
 - Change of basis can make the space more linearly separable
 - Sometimes the problem is unsolvable as-is due to the data not being linearly separable
- For multiple classes, $p(c_k | \mathbf{x}, \mathbf{w}) = \frac{e^{\alpha_k}}{\sum_i e^{-\alpha_i}}$ where $\alpha_k = \mathbf{w}_k^T \mathbf{x}$
 - This is a softmax
 - This reduces to the same sigmoid function if we only have 2 classes
- We can also replace the sigmoid with tanh (and rescale to between 0 and 1)