

Lecture 35, Apr 10, 2023

Support Vector Machines

- In normal regression we had $x \in \mathbb{R}^n$ and $y \in \mathbb{R}^n$; for classification problems, they generally have input $x \in \mathbb{R}^n$ with output $y \in \{-1, 1\}$, i.e. the output is a binary yes or no
- We can express a hyperplane as $\mathbf{w}^T \mathbf{x} - b = 0$, where \mathbf{w} is the normal vector defining the orientation and b is an offset from the origin
- The hyperplane divides all of the input space into 2 regions, $\mathbf{w}^T \mathbf{x} > b$ and $\mathbf{w}^T \mathbf{x} < b$; each region corresponds to a different value of y
- Given some data, we're looking for a hyperplane that separates the 2 types of data
 - We also want a hyperplane that's the most "in the middle" and divides the empty space between 2 types evenly
- We want to find $\mathbf{w}^T \mathbf{x} - b = 0$ that maximizes d , the distance on each side of the hyperplane, while separating the data
- Unlike in linear regression, this problem is not analytically solvable
 - In linear regression, a change in any data point is going to affect the total error and therefore change the solution; however in this problem moving a data point may not affect the solution at all
 - Depending on the orientation of the hyperplane \mathbf{w} , different data points will become relevant
- What is the expression for d ?
 - Consider 2 parallel hyperplanes, $\mathbf{w}^T \mathbf{x} = 1$ and $\mathbf{w}^T \mathbf{x} = 0$ and some point on the first hyperplane so $\mathbf{w}^T \mathbf{x}^* = 1$, then $d = \|\mathbf{x}^*\|$
 - We know $\mathbf{x}^* = \alpha \mathbf{w}$ so $\|\mathbf{x}^*\| = \frac{1}{\|\mathbf{w}\|}$, $\alpha = \frac{1}{\|\mathbf{w}\|^2} \implies \mathbf{x}^* = \frac{\mathbf{w}}{\|\mathbf{w}\|^2}$
 - $d = \frac{1}{\|\mathbf{w}\|}$
- To maximize d we want to minimize $\|\mathbf{w}\|$ or $\|\mathbf{w}\|^2$, subject to the constraint that if $y_i = 1$, then $\mathbf{w}^T \mathbf{x}_i - b > 0$, or if $y_i = -1$, then $\mathbf{w}^T \mathbf{x}_i - b < 0$
- The support vector machine is $\min_{\mathbf{w}, b} \|\mathbf{w}\|^2$ such that $y_i(\mathbf{w}^T \mathbf{x}_i - b) \geq 0$ for all training data
 - This is a quadratic program
 - If this is solvable, i.e. the data is separable, then we have an optimal classifier