

Problem Sheet 1

1 Nearest Neighbour Classification

In the lectures, we studied the perceptron, a linear classifier of the form $y = \text{sign}(\mathbf{w} \cdot \mathbf{x} + w_0)$, where $\text{sign}(z) = 1$ if $z \geq 0$ and $\text{sign}(z) = 0$ otherwise. The parameters to be learnt are \mathbf{w} and w_0 . The “Nearest neighbour classifier” (NN) is a different approach to learning from data. Suppose we are given N points $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ where $y_i \in \{0, 1\}$; for a parameter k and given a new point \mathbf{x}^* , the k -NN approach does the following: find $\mathbf{x}_{j_1}, \dots, \mathbf{x}_{j_k}$ the k -closest points to \mathbf{x}^* , then output \hat{y}^* as the majority label from the set $\{y_{j_1}, \dots, y_{j_k}\}$, *i.e.*, the most commonly occurring label among the k -nearest neighbours.

1. What advantage does the k -NN approach offer over a linear classifier like the perceptron?
2. How many parameters does the nearest neighbour model have? How much memory do you need to store the model? What is the computational cost of predicting the label \hat{y}^* ?
3. In this part, we’ll look at the setting where the vectors \mathbf{x} are points on the boolean hypercube, *i.e.*, $\mathbf{x} \in \{0, 1\}^D$. Fix $\mathbf{x}^* = (0, 0, \dots, 0)$ to be the origin and imagine that data consists of points drawn uniformly at random from the boolean hypercube. What is the distribution of the Hamming distance of data points from \mathbf{x}^* ? What happens as $D \rightarrow \infty$? (*Hint*: Use the central limit theorem.)
4. Let us now fix some numbers. Suppose the dimension of the data $D = 10,000$; let $\mathbf{x}^* = (0, 0, \dots, 0)$ and suppose we generated $N = 10,000$ data points. What do you expect the distance of \mathbf{x}^* from the nearest data-point to be? the furthest? How large does N need to be to get points that are reasonably close to \mathbf{x}^* , say within Hamming distance 50?

Remark: You do not have to write precise numbers or even mathematical expressions for the answers to part 4 above. Make sure you understand the behaviour qualitatively. The phenomenon explored in the last two parts of the question is referred to as the *curse of dimensionality*.

2 Normalization constant for a 1D Gaussian

The normalization constant for a zero-mean Gaussian is given by

$$Z = \int_{-\infty}^{\infty} \exp\left(-\frac{x^2}{2\sigma^2}\right) dx. \quad (2.1)$$

To compute this, consider its square

$$Z^2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp\left(-\frac{x^2 + y^2}{2\sigma^2}\right) dx dy. \quad (2.2)$$

Let us change variables from cartesian (x, y) to polar (r, θ) using $x = r \cos \theta$ and $y = r \sin \theta$. Since $dx dy = r dr d\theta$ (recall that r is the determinant of the Jacobian matrix in 2D) and $\cos^2 \theta + \sin^2 \theta = 1$, we have

$$Z^2 = \int_0^{2\pi} \int_0^\infty r \exp\left(-\frac{r^2}{2\sigma^2}\right) dr d\theta \quad (2.3)$$

Evaluate this integral and thus show that $Z = \sqrt{2\pi\sigma^2}$.

Hint 1: Separate the integral into a product of two integrands, the first of which (involving $d\theta$) is constant, so is easy.

Hint 2: If $u = \exp\left(-\frac{r^2}{2\sigma^2}\right)$ then $\frac{du}{dr} = -\frac{r}{\sigma^2} \cdot \exp\left(-\frac{r^2}{2\sigma^2}\right)$, so the second integral is also easy (since $\int u'(r) dr = u(r)$).

3 Reducing the cost of linear regression for large D , small N

The ridge method is a regularized version of least squares with objective function:

$$\min_{\mathbf{w} \in \mathbb{R}^D} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \quad (3.1)$$

Here λ is a scalar, the input matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$ and the output vector $\mathbf{y} \in \mathbb{R}^N$. The parameter vector $\mathbf{w} \in \mathbb{R}^D$ is obtained by differentiating the cost function, yielding the *normal equations*

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_D) \mathbf{w} = \mathbf{X}^T \mathbf{y}, \quad (3.2)$$

where \mathbf{I}_D is the $D \times D$ identity matrix. The predictions $\hat{\mathbf{y}} = \hat{\mathbf{y}}(\mathbf{X}_*)$ for new test points $\mathbf{X}_* \in \mathbb{R}^{N^* \times D}$ are obtained by evaluating the hyperplane

$$\hat{\mathbf{y}} = \mathbf{X}_* \mathbf{w} = \mathbf{X}_* (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H} \mathbf{y}. \quad (3.3)$$

The matrix \mathbf{H} is known as the *hat matrix* because it puts a “hat” on \mathbf{y} .

1. Show that the solution can be written as $\mathbf{w} = \mathbf{X}^T \tilde{\mathbf{w}}$, where $\tilde{\mathbf{w}} = \lambda^{-1}(\mathbf{y} - \mathbf{X}\mathbf{w})$.
2. Show that $\tilde{\mathbf{w}}$ can also be written as follows: $\tilde{\mathbf{w}} = (\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$ and, hence the predictions can be written as follows:

$$\hat{\mathbf{y}} = \mathbf{X}_* \mathbf{w} = \mathbf{X}_* \mathbf{X}^T \tilde{\mathbf{w}} = [\mathbf{X}_* \mathbf{X}^T] (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I}_N)^{-1} \mathbf{y}. \quad (3.4)$$

(This an *awesome trick* because if $N = 20$ patients with $D = 10,000$ gene measurements, the computation of $\tilde{\mathbf{w}}$ only requires inverting the $N \times N$ matrix, while the direct computation of \mathbf{w} would have required the inversion of a $D \times D$ matrix.)

4 Logical Gates Using Perceptrons

Recall that a perceptron with input features x_1, \dots, x_D , weights w_1, \dots, w_D and bias w_0 outputs the value:

$$y = \begin{cases} 1 & \text{if } w_0 + \sum_{i=1}^D w_i x_i \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (4.1)$$

1. Suppose there are at most two inputs and the inputs always take binary values, *i.e.*, $x_i \in \{0, 1\}$. Show how to construct AND, OR and NOT gates by suitably adjusting weights.
2. The constructions for AND and OR gates required only the bias term w_0 to be negative, all other weights were positive. Can you achieve a similar construction for the NOT gate? Why?
3. Can you construct an XOR (exclusive or) gate? If not, give reasons.
4. Often, instead of using a hard threshold we would like to use a continuous approximation. Recall the hyperbolic tangent function $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$. We consider another type of *artificial neuron* whose output is defined as

$$y = \tanh \left(w_0 + \sum_{i=1}^D w_i x_i \right). \quad (4.2)$$

Suppose you treat outputs above 0.99 as true and those below -0.99 as false. Show that similar constructions to the ones you had earlier can still be used to construct logic gates.