

ML - Exercise Sheet 5 - kernels

$$\text{Exercise 1 - logSumExp Trick} \quad P_{\times} = \prod_{i=1}^n p_i \quad P_{+} = \sum_{i=1}^n p_i$$

(a) Compute P_{\times} and P_{+} using $p_i \in [0, 1]$ - why numerical instabilities?
 When we work with probabilities $\in [0, 1]$, if we multiply many probabilities to compute P_{\times} , the result can become extremely tiny and may underflow to zero in the computer however the true result is not zero.

Example: let $n = 1000$ and $p_i = 0.01$ for $i = 1, \dots, n$

$$P_{\times} = \prod_{i=1}^{1000} p_i = (0.01)^{1000} = (10^{-2})^{1000} = 10^{-2000} \rightarrow \text{which is very small to be considered as a number not equal to zero in computer.}$$

(b) when we add numbers from $[0, 1]$, if we have very small numbers for sum of p_i 's these numbers maybe loose their effect when comparing but relatively high numbers in this interval.

$$\text{Example: } p_1 = 0.9 \quad p_2 = 0.05 \quad p_3 = (0.01)^{-10}$$

when we add these numbers in computer result occurs because:

$$P_{\times} = p_1 + p_2 + p_3 = 0.9 + 0.05 + 0.01^{-10} = 0.95 \rightarrow \text{the effect of } p_3 \text{ is not considered.}$$

(c) $\ln P_{\times}$ given $\ln p_i$'s $\Rightarrow p_i \in [k_1, k_2, \dots, k_n]$ s.t. $\ln p_i = \ln k_i + \ln p_i'$

$$P_{\times} = \prod_{i=1}^n p_i \Rightarrow \ln P_{\times} = \ln \left(\prod_{i=1}^n p_i \right) = \ln \left(\prod_{i=1}^n e^{k_i} \right) = \sum_{i=1}^n \ln k_i = \sum_{i=1}^n K_i$$

if we are given the values $k_i = \ln p_i$ we could use linear sum of the given values. Since we don't compute the product of tiny numbers, this method avoids underflow.

(d) Compute $\ln P_{+}$ given $\ln p_i$'s: $\ln P_{+} = \ln \sum_{i=1}^n \exp(\ln p_i)$ $e^{\ln p_i - c} \cdot e^c = e^{\ln p_i}$

(e) LogSumExp Trick?

$$\ln P_{+} = \ln \sum_{i=1}^n \exp(\ln p_i) \xrightarrow{c = \max \ln p_i} \ln P_{+} = \ln \left(\sum_{i=1}^n \exp(\ln p_i - c) \cdot \exp(c) \right)$$

$$= \ln(\exp(c)) + \ln \sum_{i=1}^n \exp(\ln p_i - c) = c + \ln \sum_{i=1}^n \exp(\ln p_i - c)$$

The result is the same with original formula, we just pull out the maximum value out of exponential.

(f) why this method reduce numerical stability issues?

when we subtract the max $\ln p_i$ from all p_i 's the all terms $(\ln p_i - \max \ln p_i)$ will be less or equal to zero and the exponentials will be less or equal to one always. Hence we will not have overflow for any term $(\ln p_i - c)$ and tiny probabilities ~~below zero~~ will become numbers in $[0, 1]$ preventing underflow.

Exercise Sheet 5 - Kernels

Exercise 2 : prove that Gaussian kernel is Mercer kernel

$$K(x,y) = \exp\left(-\frac{\|x-y\|^2}{2\sigma^2}\right)$$

$$\|x-y\|^2 = \|x\|^2 + \|y\|^2 - 2x^T y$$

$$\Rightarrow K(x,y) = \exp\left(-\frac{\|x\|^2 + \|y\|^2 - 2x^T y}{2\sigma^2}\right) = \exp\left(-\frac{\|x\|^2}{2\sigma^2}\right) \cdot \exp\left(-\frac{\|y\|^2}{2\sigma^2}\right) \cdot \exp\left(\frac{x^T y}{\sigma^2}\right)$$

Taylor Series for $\exp(x^T y / \sigma^2)$:

$$\exp\left(\frac{x^T y}{\sigma^2}\right) = \sum_{n=0}^{\infty} \frac{(x^T y)^n}{n! \sigma^{2n}}$$

multinomial theorem for $(x^T y)^n$:

$$(x^T y)^n = \left(\sum_{i=1}^d n_i y_i\right)^n = \sum_{k_1+k_2+\dots+k_d=n} \frac{n!}{k_1! k_2! \dots k_d!} \prod_{i=1}^d (n_i^{k_i} y_i^{k_i})$$

function of y

$$\Rightarrow (x^T y)^n = \sum_K f_K(x) \cdot f_K(y)$$

$\underbrace{\quad}_{\text{constant}} = \prod_{i=1}^d \frac{x_i^{k_i}}{k_i! k_i!}$

Now if we consider $\exp(-\frac{\|x\|^2}{2\sigma^2})$ together with $f_K(x)$ and also $\exp(-\frac{\|y\|^2}{2\sigma^2})$ with $f_K(y)$:

$$K(x,y) = \underbrace{\exp\left(-\frac{\|x\|^2}{2\sigma^2}\right)}_{\text{function of } x} \cdot \underbrace{\left(f_K(x)\right) * \exp\left(-\frac{\|y\|^2}{2\sigma^2}\right)}_{\text{function of } y} \cdot \left(f_K(y)\right)$$

$$K(x,y) = \langle F(x), F(y) \rangle$$

\hookrightarrow inner product in some other feature space