

Advanced Machine Learning Assignment

Anjila Budathoki

January 2024

1 Refreshers on Optimization and probability fundamentals.

1. (Continuous optimization) Deriving the optimal solution θ^* that minimizes $f(\theta)$.

Given, $f(\theta) = \sum_{i=1}^n w_i (x_i - \theta)^2$

Assumption, $w_i > 0$

To derive optimal solution θ^* ,

$$f(\theta) = \sum_{i=1}^n w_i (x_i - \theta)^2$$

$$\operatorname{argmin}_{\theta} f(\theta) =$$

Derivating $f(\theta)$ w.r.t θ ,

$$\frac{df(\theta)}{d\theta} = -2 \sum_{i=1}^n w_i (x_i - \theta)$$

Setting derivative to 0, to find optimal solution:

$$0 = -2 \sum_{i=1}^n w_i (x_i - \theta)$$

$$0 = \sum_{i=1}^n w_i x_i - \sum_{i=1}^n w_i \theta$$

$$\sum_{i=1}^n w_i \theta = \sum_{i=1}^n w_i x_i$$

$$n\theta^* = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

$$\theta^* = \frac{\sum_{i=1}^n x_i}{n}$$

What happens if some w_i are negative?

Answer: When the w_i are negative, these negative values might cancel out positive terms or reduce the actual function output. Hence, it wouldn't give the accurate measure.

2. (Counting and combinatorics)

The total ways of grouping $2n$ boys = $2nCn$ and total ways of grouping $2n$ boys in 2 groups is given by

$$\begin{aligned}
&= \frac{2nCn}{2!} \\
&= \frac{(2n)!}{2! * n! * n!} \\
&= \frac{(2n)!}{2 * n! * n!}
\end{aligned}$$

Now,

Removing tall boys from the group = (2n-2) and total boys in same group will be n-2, So, Total ways in which 2 tallest kids in the same group

$$\begin{aligned}
&= (2n-2)C(n-2) \\
&= \frac{(2n-2)!}{(n-2)! * (2n-2-n+2)!} \\
&= \frac{(2n-2)!}{(n-2)! * (n)!}
\end{aligned}$$

i. The probability that the two tallest kids will be in the same subgroup

$$\begin{aligned}
&= \frac{(2n-2)!}{(n-2)! * n!} * \frac{2 * n! * n!}{(2n)!} \\
&= \frac{(2n-2)!}{(n-2)!} * \frac{2n * (n-1) * (n-2)!}{2n * (2n-1) * (2n-2)!} \\
&= \frac{n-1}{2n-1}
\end{aligned}$$

ii. The probability that the two tallest kids will be in the different subgroups

$$\begin{aligned}
&= 1 - \frac{n-1}{2n-1} \\
&= \frac{2n-1-n+1}{2n-1} \\
&= \frac{n}{2n-1}
\end{aligned}$$

3. (Bayes rule)

Given:

Probability of a candidate knows the answer $P(KA) = p$

Probability of a candidate knows the answer $P(\overline{KA}) = 1 - p$

Probability of correct answer given that candidate knows the answer $P(CA | KA) = 0.99$

Probability of correct answer given that candidate does not know the answer $P(CA | \overline{KA}) = \frac{1}{k}$

To find the conditional probability that the candidate knew the answer to a question, given that she has made the correct answer: $P(KA | CA)$

$$\text{Using Bayes rule: } P(A|B) = \frac{P(B | A)P(A)}{P(A)}$$

For marginal probability $P(KA)$:

$$= P(CA | KA)P(KA) + P(CA | \overline{KA})P(\overline{KA})$$

$$= 0.99p + \frac{1}{k}(1-p)$$

$$= \frac{0.99pk + 1 - p}{k}$$

$$P(KA | CA) = \frac{P(CA | KA)P(KA)}{P(KA)}$$

$$= \frac{0.99pk}{0.99pk + 1 - p}$$

4. (Likelihood and maximum likelihood)

Given a biased coin with the probability of head = p .

Sequence of outcomes: T,H,H,T,T,H,H,H,T,H.

Probability (likelihood) of observing this sequence

$$L(p) = (1-p) \cdot p \cdot p \cdot (1-p) \cdot (1-p) \cdot p \cdot p \cdot p \cdot (1-p) \cdot p = p^6(1-p)^4$$

Calculate the value of p that maximizes the likelihood $L(p)$.

Using log on both sides:

$$\begin{aligned} \log L(p) &= \log(p^6(1-p)^4) \\ &= \log(p^6(1-p)^4) \\ &= \log(p^6) + \log(1-p)^4 \\ &= 6\log p + 4\log(1-p) \end{aligned}$$

Finding derivative of $\log L(p)$

$$\begin{aligned} \frac{d \log L(p)}{dp} &= 6\log p + 4\log(1-p) \\ &= \frac{6}{p} - \frac{4}{1-p} \end{aligned}$$

Setting derivative to 0:

$$\begin{aligned} 0 &= \frac{6}{p} - \frac{4}{1-p} \\ 0 &= \frac{6 - 6p - 4p}{p(1-p)} \end{aligned}$$

$$0 = \frac{6 - 10p}{p(1 - p)}$$

$$0 = 6 - 10p$$

$$10p = 6$$

$$p^* = \frac{6}{10}$$

The value of p that maximizes the likelihood $L(p) = \frac{6}{10}$.

Maximizing $\log L(p)$ will maximize $L(p)$ because log function is monotonic in nature, Hence, maximizing log likelihood is same as maximizing likelihood itself.

5. (Calculus, gradients)

$$\nabla_w f = \left[\frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2}, \frac{\partial f}{\partial w_3}, \dots, \frac{\partial f}{\partial w_n} \right]^T$$

Using chain-rule:

$$f(w) = \sum_{i=1}^n (x_i^T w - y_i)^2 + \lambda \sum_{i=1}^d w_i^2$$

$$\frac{\partial f}{\partial w_1} = \frac{\partial \sum_{i=1}^n (x_i^T w - y_i)^2}{\partial (x_i^T w - y_i)} * \frac{\partial (x_i^T w - y_i)}{\partial w_1} + \lambda \frac{\partial \sum_{i=1}^d (w_i)^2}{\partial w_i} * \frac{\partial (w_i)}{\partial w_1}$$

$$\frac{\partial f}{\partial w_1} = 2 \sum_{i=1}^n (x_i^T w - y_i) * x_{i1} + 2\lambda \sum_{i=1}^d w_i$$

Similarly,

$$\frac{\partial f}{\partial w_2} = 2 \sum_{i=1}^n (x_i^T w - y_i) * x_{i2} + 2\lambda \sum_{i=1}^d w_i$$

So,

$$\nabla_w f = \left[\frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2}, \frac{\partial f}{\partial w_3}, \dots, \frac{\partial f}{\partial w_n} \right]^T$$

$$\nabla_w f = 2 \sum_{i=1}^n (x_i^T w - y_i) \sum_{i=1}^d w_i [x_{i1} + \lambda, x_{i2} + \lambda, \dots, x_{in} + \lambda]^T$$

6. (Chain-rule and softmax function) Calculate the gradient of f w.r.t. vector $x = (x_1, \dots, x_n)^T$. Here, $f(x_1, x_2, \dots, x_n) = \log \sum_{i=1}^n e^{x_i}$

We know that,

$$\nabla_x f = \left[\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \frac{\partial f}{\partial x_3}, \dots, \frac{\partial f}{\partial x_n} \right]^T$$

Using chain-rule:

$$\frac{\partial f}{\partial x_1} = \frac{\partial \log(e^{x_1} + e^{x_2} + \dots + e^{x_n})}{\partial e^{x_1} + e^{x_2} + \dots + e^{x_n}} * \frac{\partial e^{x_1} + e^{x_2} + \dots + e^{x_n}}{\partial x_1}$$

$$\frac{\partial f}{\partial x_1} = \frac{1}{e^{x_1+x_2+x_3+\dots+x_n}} * e^{x_1}$$

Similarly,

$$\frac{\partial f}{\partial x_2} = \frac{\partial \log(e^{x_1} + e^{x_2} + \dots + e^{x_n})}{\partial e^{x_2} + e^{x_2} + \dots + e^{x_n}} * \frac{\partial e^{x_2} + e^{x_2} + \dots + e^{x_n}}{\partial x_2}$$

$$\frac{\partial f}{\partial x_2} = \frac{1}{e^{x_1+x_2+x_3+\dots+x_n}} * e^{x_2}$$

So,

$$\nabla_x f = [\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \frac{\partial f}{\partial x_3}, \dots, \frac{\partial f}{\partial x_n}]^T$$

$$\nabla_x f = \frac{1}{\sum_{i=1}^n e^{x_i}} [e^{x_1}, e^{x_2}, e^{x_3}, \dots, e^{x_n}]^T$$

7. (Simple mathematical proof) For the soft-max function in part (f). Prove that $\max_i x_i \leq f(x_1, \dots, x_n) \leq \max_i(x_i) + \log n$. Here, $f(x_1, x_2, \dots, x_n) = \log \sum_{i=1}^n e^{x_i}$

Taking: Lower bound i.e $\max_i x_i \leq f(x_1, \dots, x_n)$

We know that,

$$= \log a \leq \log b \text{ if } a \leq b$$

$$= e^{x_i} \geq 0$$

$$= e^{\max_i x_i} \leq \sum_{i=1}^n e^{x_i} \text{ [As, max component of } x \text{ is there along with other components of } x, \text{ so only max will be smaller than the whole component sum.]}$$

Taking log on both sides:

$$= \log(e^{\max_i x_i}) \leq \log(\sum_{i=1}^n e^{x_i})$$

$$= \max_i x_i \leq \log(\sum_{i=1}^n e^{x_i})$$

$$= \max_i x_i \leq f(x_1, \dots, x_n)$$

Taking upper bound: $f(x_1, \dots, x_n) \leq \max_i(x_i) + \log n$

We know that,

$$\log(\sum_{i=1}^n e^{x_i}) = \log(e^{\max_i(x_i)}) + \log \sum_{i=1}^n e^{x_i} - \log(e^{\max_i(x_i)})$$

$$\log(\sum_{i=1}^n e^{x_i}) = \log(e^{\max_i(x_i)}) + \log\left(\frac{\sum_{i=1}^n e^{x_i}}{e^{\max_i(x_i)}}\right)$$

$$\log(\sum_{i=1}^n e^{x_i}) \leq \max_i(x_i) + \log n$$

$$f(x_1, \dots, x_n) \leq \max_i(x_i) + \log n$$

So, combining both we get:

$$\max_i x_i \leq f(x_1, \dots, x_n) \leq \max_i(x_i) + \log n$$