

Solutions to the Final Exam for Uncertainty in AI class, Fall 2025

1. Suppose that we know the measurement results $\tilde{x}_1, \dots, \tilde{x}_n$, the data processing algorithm $y = f(x_1, \dots, x_n)$, and the standard deviations $\sigma_1, \dots, \sigma_n$.
 - a) Describe how to best estimate the standard deviation σ of the data processing result y when n is small (provide the formula) and how to estimate it when n is large (just explain the idea).
 - b) What will be the standard deviation when $y = x_1 - x_2$, $\sigma_1 = 0.3$, and $\sigma_2 = 0.4$?

Answer:

- a) When n is small, we can use an analytical formula

$$\sigma = \sqrt{c_1^2 \cdot \sigma_1^2 + \dots + c_n^2 \cdot \sigma_n^2},$$

where

$$c_i \stackrel{\text{def}}{=} \frac{\partial f}{\partial x_i}.$$

When n is large, it is more effective to use Monte-Carlo method, when we simulate the measurement errors and then use the simulated values Δy to estimate σ .

- b) For the function $f(x_1, x_2) = x_1 - x_2$, we have $c_1 = 1$, $c_2 = -1$, so

$$\sigma = \sqrt{1^2 \cdot 0.3^2 + (-1)^2 \cdot 0.4^2} = \sqrt{0.09 + 0.16} = \sqrt{0.25} = 0.5.$$

2. Suppose that we know the measurement results $\tilde{x}_1, \dots, \tilde{x}_n$, the data processing algorithm $y = f(x_1, \dots, x_n)$, and the upper bounds $\Delta_1, \dots, \Delta_n$ on the absolute values of the measurement errors.

- a) Describe how to best estimate the upper bound Δ of the absolute value of the approximation error of the measurement result y when n is small (provide the formula) and how to estimate it when n is large (just explain the idea).
- b) Describe two situations in which we only know such upper bounds and briefly explain why in these situations, we do not know the probabilities.
- c) What will be the value Δ when $y = x_1 - x_2$, $\Delta_1 = 0.3$, and $\Delta_2 = 0.4$?

Answer:

a) When n is small, we can use an analytical formula

$$\Delta = |c_1| \cdot \Delta_1 + \dots + |c_n| \cdot \Delta_n,$$

where

$$c_i \stackrel{\text{def}}{=} \frac{\partial f}{\partial x_i}.$$

When n is large, it is more effective to use Monte-Carlo method that uses Cauchy distribution, when we simulate the measurement errors and then use the simulated values Δy to estimate Δ .

b) The first situation when we cannot have probabilities, when we only know the upper the bound, is state-of-the-art measurements. The usual way to find the probabilities of different measurement errors is to compare the tested instrument with a more accurate one, but when the instrument that we are testing is the most accurate there is no more accurate one.

Another situation is measurements in industry. We can, in principle, calibrate every sensor, but calibration is expensive, so it is only done when absolutely necessary.

c) For the function $f(x_1, x_2) = x_1 - x_2$, we have $c_1 = 1$, $c_2 = -1$, so

$$\Delta = |1| \cdot 0.3 + |-1| \cdot 0.4 = 0.3 + 0.4 = 0.7.$$

3. Suppose that we use bisection to compute $\sqrt{3}$, i.e., the solution to the equation $f(x) = 0$ when $f(x) = x^2 - 3$. We know that $f(1) = 1^2 - 3 = 1 - 3 = -2 < 0$ and that $f(3) = 3^2 - 3 = 9 - 3 = 6 > 0$, so we know that the solution is somewhere on the interval $[1, 3]$. When we follow bisection method, what would we do next, and what will be the resulting new narrower interval?

Answer. We take the midpoint $m = (1 + 3)/2 = 2$ of the given interval, and compute $f(m) = f(2) = 2^2 - 3 = 4 - 3 = 1 > 0$. Since $f(1) < 0$ and $f(2) > 0$, we know the function $f(x)$ changes sign on the interval $[1, 2]$, so it must attain the value 0 somewhere on this interval.

4. Suppose that $y = f(x_1, x_2) = x_1 - x_2$. Suppose that with confidence 0.5, experts believe that the actual value of x_1 is in the interval $[1, 2]$, and that actual value of x_2 is in the interval $[2, 3]$. Describe the corresponding alpha-cut for y .

Answer. The above function $y = f(x_1, x_2)$ is strictly increasing with respect to x_1 and strictly decreasing with respect to x_2 – when both x_i are positive. So, when the inputs x_i are located in intervals, this function:

- attains its smallest value when x_1 is the smallest possible and x_2 is the largest possible, and
- attains its largest value when x_1 is the largest possible and x_2 is the smallest possible.

The alpha-cut is $[1 - 3, 2 - 2] = [-2, 0]$.

5.

- a) If 7 experts out of 10 believe that the statement A is true, what is the resulting degree of confidence?
- b) Suppose that our degree of confidence in a statement A is 0.8, in a statement B is 0.7. Suppose that we use min as “and” and max as “or”. What is our estimate for the degree of confidence in a composite statement $A \vee \neg B$?

Answer.

a) In general, if m out of n experts believe in some statement, then we assign, to this statement, the degree of certainty m/n . In our case, $m = 7$ and $n = 10$, so the degree of certainty is $7/10 = 0.7$.

b) In general, the desired degree is equal to $d = f_{\vee}(a, f_{\neg}(b))$, where a and b are our degrees of confidence in statements A and B . For $f_{\neg}(x) = 1 - x$ and for our choice of “and”- and “or”-operations, we have $d = \max(a, 1 - b)$. For given degrees of confidence a and b , we get

$$d = \max(0.8, 1 - 0.7) = \max(0.8, 0.3) = \max(0.8, 0.3) = 0.8.$$

6.

- a) How many binary questions do we need to ask a user to get his/her utility of a given alternative with accuracy 1%?
- b) Assuming that utility is proportional to the square root of money amount, would a person prefer \$9 without any condition or \$100 with probability 0.1?

Answer.

a) Before we get any answers, all we know is that the utility is somewhere on the interval $[0, 1]$. The width of this interval is 1. After each answer, the width decreases by half. So, after k questions, the width becomes 2^{-k} . The smallest k for which 2^{-k} is smaller than 0.01 is $k = 7$ – then $2^{-k} = 1/128 < 1/20$. So, we need to ask 7 questions.

b) According to decision theory, a person prefers an alternative with the largest utility.

- For the first alternative, the utility is $u_1 = \sqrt{9} = 3$ units.
- For the second alternative, the utility is

$$u_1 = 0.1 \cdot u(100) + 0.99 \cdot u(0) = 0.1 \cdot \sqrt{100} + 0.9 \cdot \sqrt{0} = 0.1 \cdot 10 + 0.9 \cdot 0 = 1.$$

The first utility is larger, so the person will prefer to get \$9 without any condition.

7. Suppose that we have two alternatives, with gains $[2, 5]$ and $[3, 4]$.

- a) Which of them are possibly optimal? definitely optimal?
- b) Which of the alternatives should we choose if Hurwicz coefficient α_H is 0.6?

Answer.

a) An alternative is definitely optimal if its lower endpoint is larger than or equal to all other upper endpoints. One can check that in this case, there is no such alternative:

- for Alternative 1, its lower endpoint 2 is not larger than the upper endpoint 4 of the second alternative;
- for Alternative 2, its lower endpoint 3 is not larger than the upper endpoint 5 of the first alternative.

An alternative is possibly optimal if its upper endpoint is larger than or equal to all the lower endpoints – i.e., equivalently, to the maximum of the lower endpoints. In our case, this maximum is equal to $\max(2, 3) = 3$. So, an alternative is possible optimal if its upper endpoint is larger than or equal to 3. This property is true for both Alternatives 1 and 2, so they are both possible optimal.

b) In the Hurwicz approach, we replace each interval $[\underline{x}, \bar{x}]$ with the value $x = \alpha_H \cdot \bar{x} + (1 - \alpha_H) \cdot \underline{x}$, and select the alternative for which this number is the largest. So, for $\alpha_H = 0.6$, we have:

$$x_1 = 0.6 \cdot 5 + (1 - 0.6) \cdot 2 = 3.0 + 0.4 \cdot 2 = 3.0 + 0.8 = 3.8,$$

$$x_2 = 0.6 \cdot 4 + (1 - 0.6) \cdot 3 = 2.4 + 0.4 \cdot 3 = 1.6 + 1.2 = 2.8.$$

so we select Alternative 1.

8–10.

8. Prove that for each final invariant optimality criterion, the optimal alternative is itself optimal.
9. Use this result to explain why ReLU is the optimal activation function.
10. Use the result from Problem 8 to explain why we should use $1 - x$ for negation, $a \cdot b$ for “and”, and $a + b - a \cdot b$ for “or”.

Answer.

8. An optimality criterion is a pair of relations ($>$, \sim), where $a > b$ means that a is better than b , and \sim means that a and b are of the same quality. An optimality criterion is final if there is exactly one optimal alternative. The criterion is invariant with respect to some transformation T if $a > b$ implies $T(a) > T(b)$ and $a \sim b$ implies $T(a) \sim T(b)$. Let a_{opt} be an optimal alternative, this means that for every alternative a we have either $a_{\text{opt}} > a$ or $a_{\text{opt}} \sim a$. In particular, for every a , we have either $a_{\text{opt}} > T^{-1}(a)$ or $a_{\text{opt}} \sim T^{-1}(a)$. Since the optimality criterion is invariant, this implies that either $T(a_{\text{opt}}) > T(T^{-1}(a)) = a$ or $T(a_{\text{opt}}) \sim a$. This is true for all a . By definition of an optimal alternative, this means that $T(a_{\text{opt}})$ is optimal. Since the optimality criterion is final, there is only one optimal alternative, so indeed $T(a_{\text{opt}}) = a_{\text{opt}}$.

9. We want to select an activation function $s(x)$ that is optimal, and it is reasonable to require that the optimality criterion should be invariant with respect to the change of a measuring unit. Thus, the optimal activation function should be thus invariant. In other words, we want to make sure that if $y = s(x)$ and we select a new measuring unit, i.e., switch to new numerical values $x' = \lambda \cdot x$ and $y' = \lambda \cdot y$, then for these new values x' and y' , we will have the exact same dependence:

$$y' = s(x').$$

Substituting the expressions $x' = \lambda \cdot x$ and $y' = \lambda \cdot y$ into this formula, we conclude that $\lambda \cdot y = s(\lambda \cdot x)$. Here, $y = s(x)$, so we conclude that $s(\lambda \cdot x) = \lambda \cdot s(x)$ for all possible x and $\lambda > 0$.

For $x = 1$, we conclude that $s(\lambda) = \lambda \cdot s(1)$. If we denote $s(1)$ by c_+ , and rename λ into z , we conclude that for all $z > 0$, we get

$$s(z) = c_+ \cdot z.$$

For $x = -1$, we conclude that $s(-\lambda) = \lambda \cdot s(-1)$. If we denote $-s(-1)$ by c_- (so that $s(-1) = -c_-$) and denote $-\lambda$ by z (so that $\lambda = -z$), we conclude that for all negative values z , we have

$$s(z) = (-c_-) \cdot (-z) = c_- \cdot z.$$

Thus, we conclude that the activation function $s(z)$ should have the following *piecewise linear* form:

- for $z > 0$, we have $s(z) = c_+ \cdot z$;
- for $z < 0$, we have $s(z) = c_- \cdot z$.

In particular, for $c_- = 0$ and $c_+ = 1$, we get ReLU. One can show that because of the linear layers, the use of any other such piecewise linear activation function is equivalent to using ReLU.

10. We use linear interpolation because it is the only one that is invariant with respect to change of measuring unit and starting point and is, therefore optimal with respect to any invariant optimality criterion. The general formula for linear interpolation has the following form: if we know the values $y_1 = f(x_1)$ and $y_2 = f(x_2)$, then for every other x , we have

$$f(x) = y_1 + \frac{y_2 - y_1}{x_2 - x_1} \cdot (x - x_1).$$

For negation, for $x_1 = 0$, we have $y_1 = f_-(0) = 1$, and for $x_2 = 1$, we have $y_2 = f_-(1) = 0$. Thus, we get

$$f_-(x) = 1 + \frac{0 - 1}{1 - 0} \cdot (x - 0) = 1 - x.$$

For “and”, let us first apply linear interpolation to find the values $f_{\&}(0, x)$ for different x . For $x_1 = 0$, we have $y_1 = f_{\&}(0, 0) = 0$, while for $x_2 = 1$, we have $y_2 = f_{\&}(0, 1) = 0$. Thus, we have

$$f_{\&}(0, x) = 0 + \frac{0 - 0}{1 - 0} \cdot (x - 0) = 0.$$

Let us now apply linear interpolation to find the values $f_{\&}(1, x)$ for different x . For $x_1 = 0$, we have $y_1 = f_{\&}(1, 0) = 0$, while for $x_2 = 1$, we have $y_2 = f_{\&}(1, 1) = 1$. Thus, we have

$$f_{\&}(1, x) = 0 + \frac{1 - 0}{1 - 0} \cdot (x - 0) = x.$$

Finally, let us apply linear interpolation to find the value $f_{\&}(x, b)$. For $x_1 = 0$, we have $y_1 = f_{\&}(0, b) = 0$, while for $x_2 = 1$, we have $y_2 = f_{\&}(1, b) = b$. Thus, we have

$$f_{\&}(x, b) = 0 + \frac{b - 0}{1 - 0} \cdot (x - 0) = b \cdot x.$$

For “or”, let us first apply linear interpolation to find the values $f_{\vee}(0, x)$ for different x . For $x_1 = 0$, we have $y_1 = f_{\vee}(0, 0) = 0$, while for $x_2 = 1$, we have $y_2 = f_{\vee}(0, 1) = 1$. Thus, we have

$$f_{\vee}(0, x) = 0 + \frac{1 - 0}{1 - 0} \cdot (x - 0) = x.$$

Let us now apply linear interpolation to find the values $f_{\vee}(1, x)$ for different x . For $x_1 = 0$, we have $y_1 = f_{\vee}(1, 0) = 1$, while for $x_2 = 1$, we have $y_2 = f_{\vee}(1, 1) = 1$. Thus, we have

$$f_{\vee}(1, x) = 1 + \frac{1-1}{1-0} \cdot (x-0) = 1.$$

Finally, let us apply linear interpolation to find the value $f_{\vee}(x, b)$. For $x_1 = 0$, we have $y_1 = f_{\vee}(0, b) = b$, while for $x_2 = 1$, we have $y_2 = f_{\vee}(1, b) = 1$. Thus, we have

$$f_{\vee}(x, b) = b + \frac{1-b}{1-0} \cdot (x-0) = b + (1-b) \cdot x = b + x - b \cdot x.$$