

CS221 Problem Workout Solutions

Sept 23

1) [CA session] Problem 1: Backpropagation

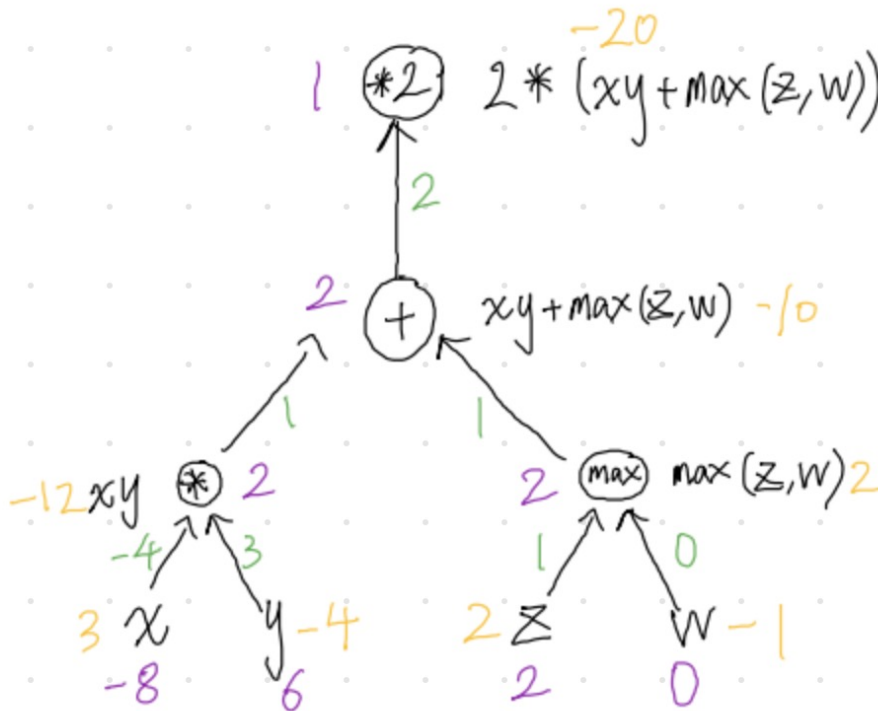
Consider the following function

$$\text{Loss}(x, y, z, w) = 2(xy + \max\{w, z\})$$

Run the backpropagation algorithm to compute the gradient $\nabla_w \text{Loss}(x, y, z, w)$ and $\nabla_z \text{Loss}(x, y, z, w)$ at $x = 3$, $y = -4$, $z = 2$ and $w = -1$. Use the following nodes: addition, multiplication, max, multiplication by a constant.

Solution When calculating the gradients, we run backpropagation from the root node to the leaves nodes. As shown on the computation graph below, the purple values are the gradients of Loss with respect to each node.

We have $\nabla_w \text{Loss}(x, y, z, w) = 0$ and $\nabla_z \text{Loss}(x, y, z, w) = 2$.



2) [CA session] Problem 2: Non-linear features

Consider the following two training datasets of (x, y) pairs:

- $\mathcal{D}_1 = \{(-1, +1), (0, -1), (1, +1)\}$.
- $\mathcal{D}_2 = \{(-1, -1), (0, +1), (1, -1)\}$.

Observe that neither dataset is linearly separable if we use $\phi(x) = x$, so let's fix that.

Define a two-dimensional feature function $\phi(x)$ such that:

- There exists a weight vector \mathbf{w}_1 that classifies \mathcal{D}_1 perfectly (meaning that $\mathbf{w}_1 \cdot \phi(x) > 0$ if x is labeled $+1$ and $\mathbf{w}_1 \cdot \phi(x) < 0$ if x is labeled -1); and
- There exists a weight vector \mathbf{w}_2 that classifies \mathcal{D}_2 perfectly.

Note that the weight vectors can be different for the two datasets, but the features $\phi(x)$ must be the same.

Solution One option is $\phi(x) = [1, x^2]$, and using $\mathbf{w}_1 = [-1, 2]$ and $\mathbf{w}_2 = [1, -2]$.

Then in \mathcal{D}_1 :

- For $x = -1$, $\mathbf{w}_1 \cdot \phi(x) = [-1, 2] \cdot [1, 1] = 1 > 0$
- For $x = 0$, $\mathbf{w}_1 \cdot \phi(x) = [-1, 2] \cdot [1, 0] = -1 < 0$
- For $x = 1$, $\mathbf{w}_1 \cdot \phi(x) = [-1, 2] \cdot [1, 1] = 1 > 0$

In \mathcal{D}_2 :

- For $x = -1$, $\mathbf{w}_2 \cdot \phi(x) = [1, -2] \cdot [1, 1] = -1 < 0$
- For $x = 0$, $\mathbf{w}_2 \cdot \phi(x) = [1, -2] \cdot [1, 0] = 1 > 0$
- For $x = 1$, $\mathbf{w}_2 \cdot \phi(x) = [1, -2] \cdot [1, 1] = -1 < 0$

Note that there are many options that work, so long as -1 and 1 are separated from 0 .

3) [breakout] Problem 3: Non-linear decision boundaries

Suppose we are performing classification where the input points are of the form $(x_1, x_2) \in \mathbb{R}^2$. We can choose any subset of the following set of features:

$$\mathcal{F} = \left\{ x_1^2, x_2^2, x_1x_2, x_1, x_2, \frac{1}{x_1}, \frac{1}{x_2}, 1, \mathbf{1}[x_1 \geq 0], \mathbf{1}[x_2 \geq 0] \right\} \quad (1)$$

For each subset of features $F \subseteq \mathcal{F}$, let $D(F)$ be the set of all decision boundaries corresponding to linear classifiers that use features F .

For each of the following sets of decision boundaries E , provide the minimal F such that $D(F) \supseteq E$. If no such F exists, write ‘none’.

- E is all lines [CA hint]:

 (2)

- E is all circles centered at the origin:

 (3)

- E is all circles:

 (4)

- E is all axis-aligned rectangles:

 (5)

- E is all axis-aligned rectangles whose lower-right corner is at $(0, 0)$:

 (6)

Solution

- Lines: $x_1, x_2, 1$ ($ax_1 + bx_2 + c = 0$)
- Circles centered at the origin: $x_1^2, x_2^2, 1$ ($x_1^2 + x_2^2 = r^2$)
- Circles centered anywhere in the plane: $x_1^2, x_2^2, x_1, x_2, 1$ ($(x_1 - a)^2 + (x_2 - b)^2 = r^2$)
- Axis aligned rectangles: not possible (need features of the form $\mathbf{1}[x_1 \leq a]$)
- Axis aligned rectangles with lower right corner at $(0, 0)$: not possible

4) [breakout, optional] Problem 4: K-means

Consider doing ordinary K -means clustering with $K = 2$ clusters on the following set of 3 one-dimensional points:

$$\{-2, 0, 10\}. \tag{7}$$

Recall that K -means can get stuck in local optima. Describe the precise conditions on the initialization $\mu_1 \in \mathbb{R}$ and $\mu_2 \in \mathbb{R}$ such that running K -means will yield the global optimum of the objective function. Notes:

- Assume that $\mu_1 < \mu_2$.
- Assume that if in step 1 of K -means, no points are assigned to some cluster j , then in step 2, that centroid μ_j is set to ∞ .
- Hint: try running K -means from various initializations μ_1, μ_2 to get some intuition; for example, if we initialize $\mu_1 = 1$ and $\mu_2 = 9$, then we converge to $\mu_1 = -1$ and $\mu_2 = 10$.

Solution The objective is minimized for $\mu_1 = -1$ and $\mu_2 = 10$. First, note that if all three points end up in one cluster, K -means definitely fails to recover the global optimum. Therefore, -2 must be assigned to the first cluster, and 10 must be assigned to the second cluster. 0 can be assigned to either: If 0 is assigned to cluster 1, then we're done. If it is assigned to cluster 2, then we have $\mu_1 = -2, \mu_2 = 5$; in the next iteration, 0 will be assigned to cluster 1 since it's closer. Therefore, the condition on the initialization written formally is $|-2 - \mu_1| < |-2 - \mu_2|$ and $|10 - \mu_1| > |10 - \mu_2|$.