Notation:

- $a \in \{\text{cloudy}, \text{clear}, \text{partial}, \text{haze}\}$ represents the possible atmospheric conditions.

- Vector $\boldsymbol{l} \in \left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \\ \vdots \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \\ 0 \\ \vdots \end{pmatrix}, \cdots \right\}$ represents whether the other labels are present, such that
  $l_i \in \{0, 1\}$ represents whether label $i$ is present.

We are interested in obtaining the joint posterior distribution, $P(a, \boldsymbol{l}|\boldsymbol{o}_t)$. With this posterior, we can do training using either CE,

$$\mathcal{F}_{\text{CE}} = -\sum_t \log P(a_t^*, \boldsymbol{l}_t^*|\boldsymbol{o}_t),$$

or maximum expected F-score,

$$\mathcal{F}_{\text{F}} = -\sum_t \sum_{a, \boldsymbol{l}} F(a, \boldsymbol{l}, a_t^*, \boldsymbol{l}_t^*) P(a, \boldsymbol{l}|\boldsymbol{o}_t).$$

With this posterior, we can also do decoding either MAP,

$$a^*, \boldsymbol{l}^* = \arg\max_{a, \boldsymbol{l}} P(a, \boldsymbol{l}|\boldsymbol{o}_t),$$

or maximum expected F-score,

$$a^*, \boldsymbol{l}^* = \arg\max_{a, \boldsymbol{l}} \sum_{a', \boldsymbol{l}'} F(a', \boldsymbol{l}', a, \boldsymbol{l}) P(a', \boldsymbol{l}'|\boldsymbol{o}_t).$$

# 1 Obtaining the joint posterior distribution

There are several possible ways that we can choose to formulate the posterior:

1. Sigmoid for $c \Rightarrow$ cloudy, softmax for $b \Rightarrow \{\text{clear, partial, haze}\}$, and separate sigmoids for $\boldsymbol{l} \Rightarrow \{\text{other labels}\}$. $b$ and $\boldsymbol{l}$ are conditionally independent given $c$.

$$P(c, b, \boldsymbol{l}|\boldsymbol{o}_t) = P(c|\boldsymbol{o}_t) P(b|c, \boldsymbol{o}_t) P(\boldsymbol{l}|c, \boldsymbol{o}_t)$$

2. Softmax for $a \Rightarrow \{\text{cloudy, clear, partial, haze}\}$ and separate sigmoids for $\boldsymbol{l} \Rightarrow \{\text{other labels}\}$.

$$P(a, \boldsymbol{l}|\boldsymbol{o}_t) = P(a|\boldsymbol{o}_t) P(\boldsymbol{l}|a, \boldsymbol{o}_t)$$

3. Other labels are conditionally independent of the atmospheric condition.

$$P(a, \boldsymbol{l}|\boldsymbol{o}_t) = P(a|\boldsymbol{o}_t) P(\boldsymbol{l}|\boldsymbol{o}_t)$$

## 1.1 Group outputs into cloudy, {clear, partial, haze}, {other labels}

We need to formulate the joint label probability, $P(c, b, \boldsymbol{l}|\boldsymbol{o}_t)$. We can factorise the distribution into

$$P(c, b, \boldsymbol{l}|\boldsymbol{o}_t) = P(c|\boldsymbol{o}_t) P(b|c, \boldsymbol{o}_t) P(\boldsymbol{l}|c, \boldsymbol{o}_t).$$

Here, we have assumed that $b$ and $\boldsymbol{l}$ are conditionally independent.

We can then model $P(c|\boldsymbol{o}_t)$ as a sigmoid output, $f(\boldsymbol{o}_t)$,

$$P(c|\boldsymbol{o}_t) = f^c(\boldsymbol{o}_t)[1 - f(\boldsymbol{o}_t)]^{1-c}.$$

We can model $P(a|c = 1, \boldsymbol{o}_t)$ as a softmax output, $\boldsymbol{g}(\boldsymbol{o}_t)$, then we have

$$P(b|c, \boldsymbol{o}_t) = [\delta(b, 0)]^c g_b^{1-c}(\boldsymbol{o}_t).$$

It can be seen that there is a problem, because 0 is not in the set of values that $b$ can take.

We assume that all other labels are independent of each other,

$$P(\boldsymbol{l}|c, \boldsymbol{o}_t) = \prod_i P(l_i|c, \boldsymbol{o}_t).$$

We can model $P(l_i|c = 1, \boldsymbol{o}_t)$ as a sigmoid, $h_i(\boldsymbol{o}_t)$, then we have

$$P(l_i|c, \boldsymbol{o}_t) = [\delta(l_i, 0)]^c \left[ h_i^{l_i}(\boldsymbol{o}_t) \{1 - h_i(\boldsymbol{o}_t)\}^{1-l_i} \right]^{1-c}.$$

The joint distribution is

$$P(c, b, \boldsymbol{l}|\boldsymbol{o}_t) = f^c(\boldsymbol{o}_t)[1 - f(\boldsymbol{o}_t)]^{1-c}[\delta(b, 0)]^c g_b^{1-c}(\boldsymbol{o}_t) \prod_i [\delta(l_i, 0)]^c \left[ h_i^{l_i}(\boldsymbol{o}_t) \{1 - h_i(\boldsymbol{o}_t)\}^{1-l_i} \right]^{1-c}.$$

The CE criterion for frame $t$ is

$$\begin{aligned}
\mathcal{F}_{\mathrm{CE}} &= -\log P(c_t^*, b_t^*, \boldsymbol{l}_t^*|\boldsymbol{o}_t) \\
&= -c_t^* \log f(\boldsymbol{o}_t) - (1 - c_t^*) \log[1 - f(\boldsymbol{o}_t)] - c_t^* \log \delta(b_t^*, 0) - (1 - c_t^*) g_{b_t^*}(\boldsymbol{o}_t) \\
&\quad - \sum_i \{c_t^* \log \delta(l_{it}^*, 0) + (1 - c_t^*)[l_{it}^* \log h_i(\boldsymbol{o}_t) + (1 - l_{it}^*) \log(1 - h_i(\boldsymbol{o}_t))]\}
\end{aligned}$$

There is a problem, because $b$ cannot take the value 0. So, $\log \delta(b, 0)$ is always infinite. However, the gradient is still finite, as $\log \delta(b, 0)$ is not involved in the gradient. When we discard all the terms that do not depend on the model parameters, the criterion reduces to

$$\mathcal{F}_{\mathrm{CE}} \Rightarrow -c_t^* \log f(\boldsymbol{o}_t) - (1 - c_t^*) \log[1 - f(\boldsymbol{o}_t)] - (1 - c_t^*) \sum_i [l_{it}^* \log h_i(\boldsymbol{o}_t) + (1 - l_{it}^*) \log[1 - h_i(\boldsymbol{o}_t)]].$$

## 1.2 Group outputs into {atmophere}, {other labels}

To avoid the infinite CE cost, we can instead group together all four atmospheric conditions into a single softmax output. The joint posterior can then be factorised as

$$P\left(a, \boldsymbol{l} | \boldsymbol{o}_t\right) = P\left(a | \boldsymbol{o}_t\right) P\left(\boldsymbol{l} | a, \boldsymbol{o}_t\right).$$

We can model $P\left(a | \boldsymbol{o}_t\right)$ as a softmax, $\boldsymbol{g}\left(\boldsymbol{o}_t\right)$,

$$P\left(a | \boldsymbol{o}_t\right) = g_a\left(\boldsymbol{o}_t\right).$$

We can again model $P\left(\boldsymbol{l} | a, \boldsymbol{o}_t\right) = \prod_i P\left(l_i | a, \boldsymbol{o}_t\right)$ as separate sigmoids, $h_i\left(\boldsymbol{o}_t\right)$, then we have

$$P\left(l_i | a, \boldsymbol{o}_t\right) = \left[\delta\left(l_i, 0\right)\right]^{\delta(a, \text{cloudy})} \left[h_i^{l_i}\left(\boldsymbol{o}_t\right) \left\{1 - h_i\left(\boldsymbol{o}_t\right)\right\}^{1-l_i}\right]^{1-\delta(a, \text{cloudy})}.$$

The joint distribution is

$$P\left(a, \boldsymbol{l} | \boldsymbol{o}_t\right) = g_a\left(\boldsymbol{o}_t\right) \prod_i \left[\delta\left(l_i, 0\right)\right]^{\delta(a, \text{cloudy})} \left[h_i^{l_i}\left(\boldsymbol{o}_t\right) \left\{1 - h_i\left(\boldsymbol{o}_t\right)\right\}^{1-l_i}\right]^{1-\delta(a, \text{cloudy})}.$$

The CE criterion for frame $t$ is

$$\mathcal{F}_{\text{CE}} = -\log g_{a_t^*}\left(\boldsymbol{o}_t\right) - \sum_i \left\{\delta\left(a_t^*, \text{cloudy}\right) \log \delta\left(l_{it}^*, 0\right) + \left[1 - \delta\left(a_t^*, \text{cloudy}\right)\right] \left[l_{it}^* \log h_i\left(\boldsymbol{o}_t\right) + \left(1 - l_{it}^*\right) \log \left(1 - h_i\left(\boldsymbol{o}_t\right)\right)\right]\right\}$$

By discarding all the terms that do not depend on the model parameters, the criterion reduces to

$$\mathcal{F}_{\text{CE}} \Rightarrow -\log g_{a_t^*}\left(\boldsymbol{o}_t\right) - \left[1 - \delta\left(a_t^*, \text{cloudy}\right)\right] \sum_i \left[l_{it}^* \log h_i\left(\boldsymbol{o}_t\right) + \left(1 - l_{it}^*\right) \log \left(1 - h_i\left(\boldsymbol{o}_t\right)\right)\right].$$

## 1.3  Conditionally independent $a$ and $l$

Although the task instructions mentions that images labelled with cloudy should not have any other labels, labelling errors do exists in the training set. It may be possible that labelling errors also exist in the test set. If this is the case it may be better not to force $P(l = \mathbf{0}|a = \text{cloudy}, \mathbf{o}_t) = 1$. We can remove this constraint by assuming that $a$ and $l$ are conditionally independent, so that the joint posterior factorises into

$$P(a, l|\mathbf{o}_t) = P(a|\mathbf{o}_t) P(l|\mathbf{o} - t).$$

We can again model $P(a|\mathbf{o}_t)$ as a softmax,

$$P(a|\mathbf{o}_t) = g_a(\mathbf{o}_t).$$

We can model $P(l|\mathbf{o} - t) = \prod_i P(l_i|\mathbf{o} - t)$ as separate sigmoids, regardless of the atmospheric condition,

$$P(l_i|\mathbf{o}_t) = h_i^{l_i}(\mathbf{o}_t) \{1 - h_i(\mathbf{o}_t)\}^{1-l_i}.$$

The joint posterior is then

$$P(a, l|\mathbf{o}_t) = g_a(\mathbf{o}_t) \prod_i h_i^{l_i}(\mathbf{o}_t) \{1 - h_i(\mathbf{o}_t)\}^{1-l_i}.$$

The CE criterion for frame $t$ is

$$\mathcal{F}_{\text{CE}} = -\log g_{a_t^*}(\mathbf{o}_t) - \sum_i \{l_{it}^* \log h_i(\mathbf{o}_t) + (1 - l_{it}^*) \log[1 - h_i(\mathbf{o}_t)]\}.$$