$P_\theta$: model output $\in \mathbb{R}^N$

$g$: one-hot ground truth $\in \mathbb{R}^N$

$u$: a prior distribution $\in \mathbb{R}^N$

- Smoothed target ($\varepsilon \in [0, 1]$)

$$t = (1 - \varepsilon) g + \varepsilon u \in \mathbb{R}^N$$

- loss (constant terms are ignored)

$$L = -t \cdot \log P_\theta \quad \leftarrow \text{dot product}$$

$$= (1 - \varepsilon) \underbrace{(-g \cdot \log P_\theta)}_{\text{cross entropy}} - \underbrace{\varepsilon (u \cdot \log P_\theta)}_{(*)}$$

$(*)$: if $u$ is the uniform distribution,

$$u = [\tfrac{1}{N} \ \tfrac{1}{N} \ \cdots ]$$

$$(u \cdot \log P_\theta) = \text{mean} (\log P_\theta)$$

Note: we can use whatever $u$