$\theta_s \rightarrow$ shared weights (CNN)
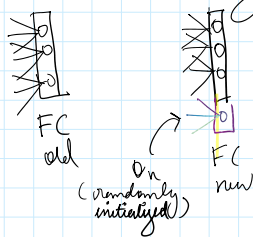
$\theta_o \rightarrow$ old weights (pre-trained)

$\theta_n \rightarrow$ new weights (fine tuning)

**Method**

1) Record responses $y_o$ on each new task image from the original network for outputs on the old tasks defined by $\theta_s$ & $\theta_o$.

$y_o \rightarrow$ set of class probabilities for each image.



FC old

$\theta_n$ (randomly initialized)

FC new

number of new parameters = no. of new classes
$\times$ no. of nodes in last shared layer

2) Train model to minimize loss for all tasks and Regularization R
using **SGD**

$R \rightarrow$ Weight decay of 0.0005

3) First freeze $\theta_s$ and $\theta_o$ & train $\theta_n$ to convergence. (warm up)

4) Jointly train all weights $\theta_s$, $\theta_o$, $\theta_n$ until convergence (joint optimize step)

☆ Warm up → enhances fine tuning's old task performance (not crucial)

Loss encourages predictions $\hat{y}_n$ to be consistent with the gnd. truths $y_n$

Multinomial logistic loss

$$L_{new}(y_n, \hat{y}_n) = -y_n . \log \hat{y}_n$$

$\hat{y}_n = $ softmax output

$y_n \rightarrow$ one hot encoded gnd truth vector

multilabel classification

multilabel classification

↓↓

true false prediction
for each label

for original task → output prob
use { as close as possible
knowledge { to output of
distillation { original network
loss.

$$\mathcal{L}_{old}(y_0, \hat{y}_0) = -H(y_0', \hat{y}_0')$$

$$= -\sum_{i=1}^{l} y_0'^{(i)} \log \hat{y}_0'^{(i)}$$

where $l \Rightarrow$ no. of labels
$y_0'^{(i)} \to$ recorded prob
$\hat{y}_0'^{(i)} \to$ current prob } Modified version of $y_0^{(i)}$ & $\hat{y}_0^{(i)}$

☆ Modified CE loss which increases the weight for smaller probabilities.

$$y_0'^{(i)} = \frac{\left(y_0^{(i)}\right)^{1/T}}{\sum_j \left(y_0^{(j)}\right)^{1/T}} \qquad \hat{y}_0'^{(i)} = \frac{\left(\hat{y}_0^{(i)}\right)^{1/T}}{\sum_j \left(\hat{y}_0^{(j)}\right)^{1/T}}$$

$T > 1 \to$ increases the weight of the smaller of logits.

$T = 2$ is used ( found by grid search)

$\lambda_0 \to 1$ { loss balance weight

larger values favour old task performance over new task

[1 = balance]

⊗ Instead of storing old task data [joint training] we use recorded responses (Yo) instead.

[Joint training] we use recorded responses $Y_o$ instead.

This is used for joint optimization of $\theta_s$

# Algorithm

start:

$\theta_s \rightarrow$ shared params

$\theta_o \rightarrow$ old task params

$X_n, Y_n \rightarrow$ training data & ground truth of the new task.

Initialize

$Y_o \leftarrow CNN(X_n, \theta_s, \theta_o)$

$\theta_n \leftarrow RandInt(|\theta_n|)$

Train

$\hat{Y}_o = CNN(X_n, \hat{\theta}_s, \hat{\theta}_o)$ // old task

$\hat{Y}_n = CNN(X_n, \hat{\theta}_s, \hat{\theta}_n)$ // new task

$\theta_s^*, \theta_o^*, \theta_n^* \leftarrow \underset{(\hat{\theta}_s, \hat{\theta}_o, \hat{\theta}_n)}{argmin} \left( \lambda_o L_{old}(Y_o, \hat{Y}_o) + L_{new}(Y_n, \hat{Y}_n) + R(\hat{\theta}_s, \hat{\theta}_o, \hat{\theta}_n) \right)$