# Semi-Supervised Learning

Eugene Lim
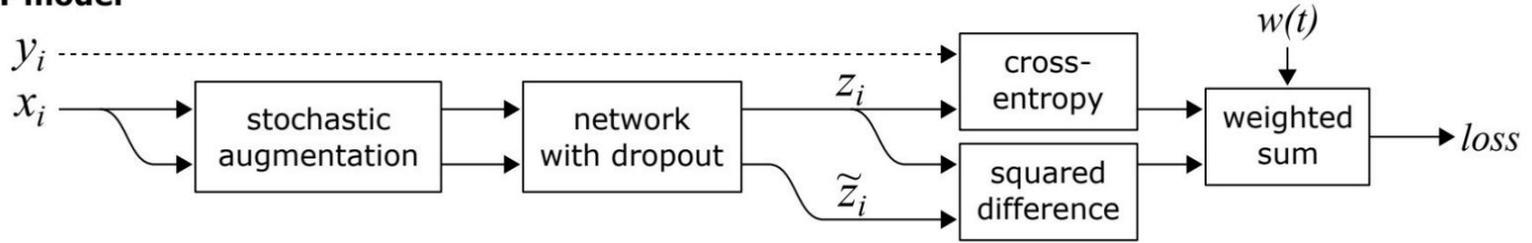
$$(x, y) \sim p(x, y)$$
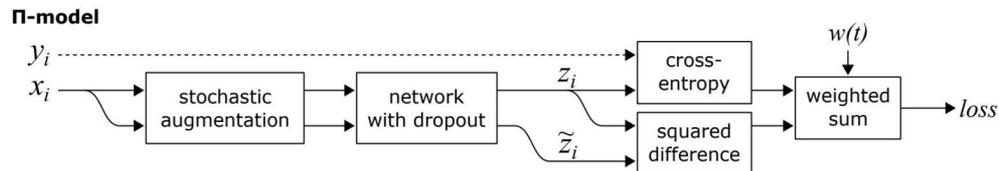
$$(x, y) \sim p(x, y)$$

$$x \sim p(x)$$

# Pi Model

**Π-model**

Pi Model

**Π-model**



---

**Algorithm 1** Π-model pseudocode.

---

**Require:** $x_i$ = training stimuli
**Require:** $L$ = set of training input indices with known labels
**Require:** $y_i$ = labels for labeled inputs $i \in L$
**Require:** $w(t)$ = unsupervised weight ramp-up function
**Require:** $f_\theta(x)$ = stochastic neural network with trainable parameters $\theta$
**Require:** $g(x)$ = stochastic input augmentation function
  **for** $t$ in $[1, num\_epochs]$ **do**
    **for** each minibatch $B$ **do**
      $z_{i \in B} \leftarrow f_\theta(g(x_{i \in B}))$          ▷ evaluate network outputs for augmented inputs
      $\tilde{z}_{i \in B} \leftarrow f_\theta(g(x_{i \in B}))$          ▷ again, with different dropout and augmentation
      $loss \leftarrow -\frac{1}{|B|} \sum_{i \in (B \cap L)} \log z_i[y_i]$          ▷ supervised loss component
           $+ w(t)\frac{1}{C|B|} \sum_{i \in B} \|z_i - \tilde{z}_i\|^2$          ▷ unsupervised loss component
      update $\theta$ using, e.g., ADAM          ▷ update network parameters
    **end for**
  **end for**
  **return** $\theta$

---

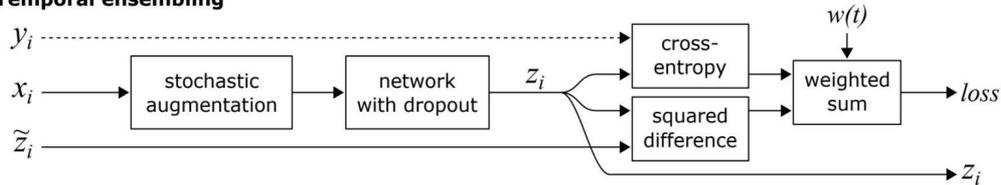# Temporal Ensembling

**Temporal ensembling**

**Temporal ensembling**



---

**Algorithm 2**  Temporal ensembling pseudocode. Note that the updates of $Z$ and $\tilde{z}$ could equally well be done inside the minibatch loop; in this pseudocode they occur between epochs for clarity.

---

**Require:** $x_i$ = training stimuli
**Require:** $L$ = set of training input indices with known labels
**Require:** $y_i$ = labels for labeled inputs $i \in L$
**Require:** $\alpha$ = ensembling momentum, $0 \le \alpha < 1$
**Require:** $w(t)$ = unsupervised weight ramp-up function
**Require:** $f_\theta(x)$ = stochastic neural network with trainable parameters $\theta$
**Require:** $g(x)$ = stochastic input augmentation function
$\quad Z \leftarrow \mathbf{0}_{[N \times C]}$ $\qquad\qquad\qquad\qquad\qquad\quad$ ▷ initialize ensemble predictions
$\quad \tilde{z} \leftarrow \mathbf{0}_{[N \times C]}$ $\qquad\qquad\qquad\qquad\qquad\quad$ ▷ initialize target vectors
$\quad$**for** $t$ in $[1, num\_epochs]$ **do**
$\qquad$**for** each minibatch $B$ **do**
$\qquad\quad z_{i \in B} \leftarrow f_\theta(g(x_{i \in B}, t))$ $\qquad\qquad\quad$ ▷ evaluate network outputs for augmented inputs
$\qquad\quad loss \leftarrow -\frac{1}{|B|} \sum_{i \in (B \cap L)} \log z_i[y_i]$ $\qquad$ ▷ supervised loss component
$\qquad\qquad + w(t) \frac{1}{C|B|} \sum_{i \in B} \|z_i - \tilde{z}_i\|^2$ $\qquad$ ▷ unsupervised loss component
$\qquad\quad$ update $\theta$ using, e.g., ADAM $\qquad\qquad\quad$ ▷ update network parameters
$\qquad$**end for**
$\qquad Z \leftarrow \alpha Z + (1 - \alpha)z$ $\qquad\qquad\qquad\quad$ ▷ accumulate ensemble predictions
$\qquad \tilde{z} \leftarrow Z/(1 - \alpha^t)$ $\qquad\qquad\qquad\qquad$ ▷ construct target vectors by bias correction
$\quad$**end for**
$\quad$**return** $\theta$

---

# How to do better?

**Temporal ensembling**

$y_i$ ---- → cross-entropy

$x_i$ → stochastic augmentation → network with dropout → $z_i$ → cross-entropy / squared difference → weighted sum ← $w(t)$ → $loss$

$\widetilde{z}_i$ → squared difference

$z_i$
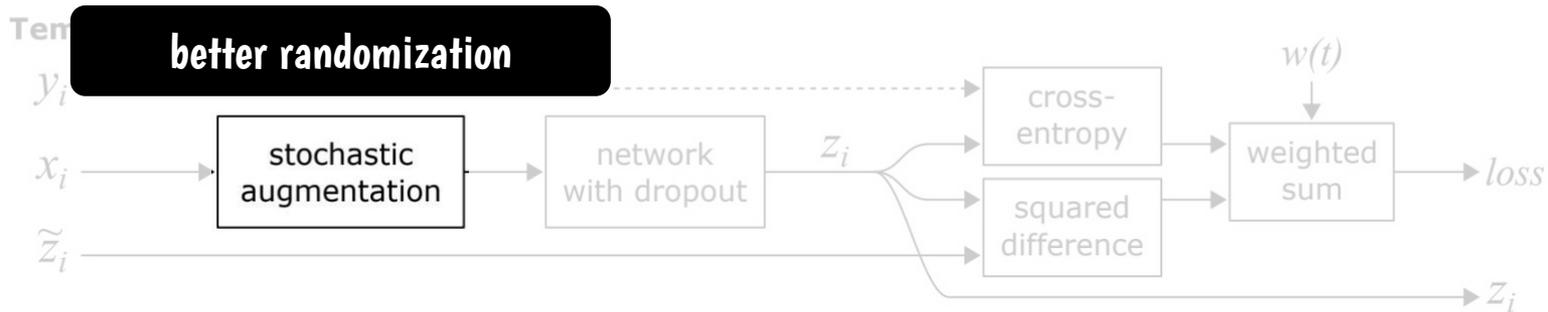
How to do better?

Temporal ensembling

How to do better?

Ten

**better randomization**

$y_i$

$x_i$ ────────▶ ┌─────────────────┐
                │   stochastic    │
                │  augmentation   │
                └─────────────────┘

$\widetilde{z}_i$ ──────────────────────

```
          network
        with dropout            $z_i$
```

$w(t)$

cross-
entropy

squared
difference

weighted
sum ────▶ $loss$

$z_i$

# Mean Teacher

Temporal ensembling

$y_i$

$x_i$ → stochastic augmentation → network with dropout → $z_i$ → cross-entropy

$w(t)$

weighted sum → $loss$

squared difference

$\widetilde{z}_i$

$z_i$

**better "targets"**

Mean Teacher

I'M A ROLE MODEL

imgflip.com

Mean Teacher

I'M A ROLE MODEL

BECAUSE I'M AN AVERAGE OF ALL MY PAST STUDENTS

imgflip.com

Mean Teacher

Mean Teacher

# Virtual Adversarial Training

**better randomization**

Te... $y_i$

$x_i$ → stochastic augmentation → network with dropout → $z_i$ → cross-entropy → weighted sum → *loss*

$\widetilde{z}_i$ → squared difference
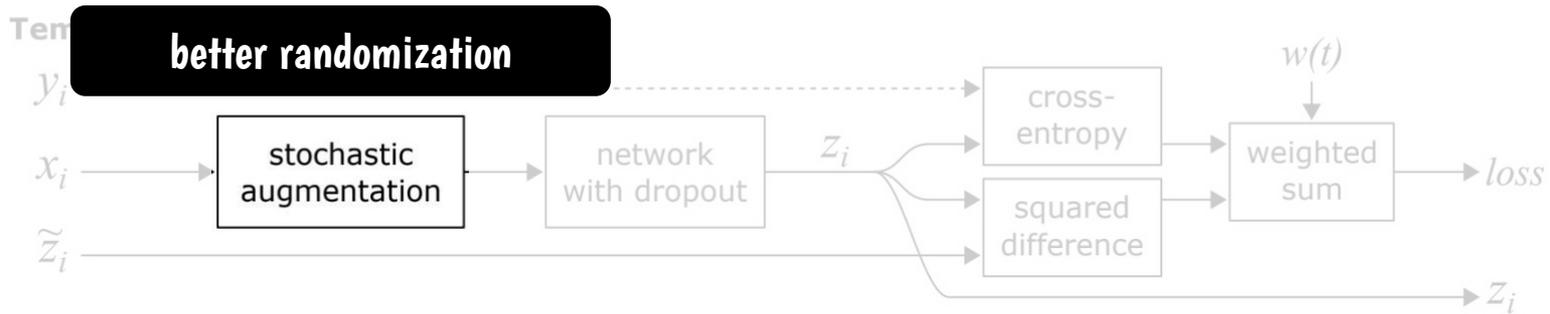
$w(t)$

$z_i$

Virtual Adversarial Training

Virtual Adversarial Training

$$L_{\text{adv}}(x_l, \theta) := D\left[q(y|x_l), p(y|x_l + r_{\text{adv}}, \theta)\right]$$

$$\text{where } r_{\text{adv}} := \arg\max_{r; \|r\| \leq \epsilon} D\left[q(y|x_l), p(y|x_l + r, \theta)\right],$$

$$L_{\text{adv}}(x_l, \theta) \approx [r(y|x), p(y|x_l + r_{\text{adv}}, \theta)]$$

$$\text{where } r_{\text{adv}} := \arg\max_{r;\|r\|\leq\epsilon} D\left[q(y|x_l), p(y|x_l + r, \theta)\right],$$

😢 **no closed form** 😢

$$L_{\text{adv}}(x_l, \theta) := D\left[q(y|x_l), p(y|x_l + r_{\text{adv}}, \theta)\right]$$

$$\text{where } r_{\text{adv}} := \underset{r; \|r\| \le \epsilon}{\arg\max} D\left[q(y|x_l), p(y|x_l + r, \theta)\right],$$

$$r_{\text{adv}} \quad \approx \quad \epsilon \frac{g}{\|g\|_2}, \text{ where } g = \nabla_{x_l} D\left[h(y; y_l), p(y|x_l, \theta)\right]$$

Adversarial training is a successful method that works for many supervised problems. However, full label information is not available at all times. Let $x_*$ represent either $x_l$ or $x_{ul}$. Our objective function is now given by

$$D\left[q(y|x_*), p(y|x_* + r_{\text{qadv}}, \theta)\right]$$

$$\text{where } r_{\text{qadv}} := \arg\max_{r; \|r\| \leq \epsilon} D\left[q(y|x_*), p(y|x_* + r, \theta)\right],$$

Adversarial training is a successful method that works for many supervised problems. However, full label information is not available at all times. Let $x_*$ represent either $x_l$ or $x_{ul}$. 😢 **we don't have this** 😢 n is now given by
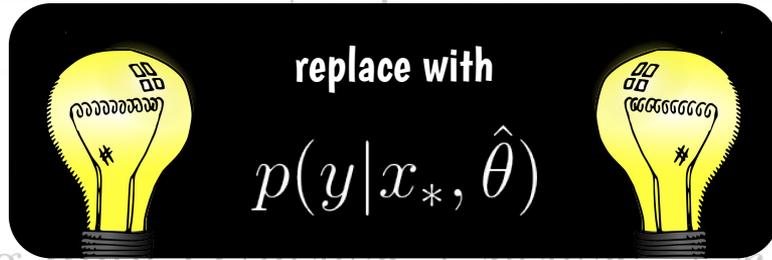
$$D\left[q(y|x_*), p(y|x_* + r_{\text{qadv}}, \theta)\right]$$

$$\text{where } r_{\text{qadv}} := \arg\max_{r; \|r\| \leq \epsilon} D\left[q(y|x_*), p(y|x_* + r, \theta)\right],$$

Adversarial training is a successful method that works for many supervised problems. However, full label information is not available at all times. Let $x_*$ represent either $x_l$ or $x_{ul}$.



$$D\left[q(y|x_*), p(y|x_*\right.$$

$$p(y|x_*, \hat{\theta})$$

$$\text{where } r_{\text{qadv}} := \arg\max_{r; \|r\| \leq \epsilon} D\left[q(y|x_*), p(y|x_*+r, \theta)\right],$$

$$\text{LDS}(x_*, \theta) := D\left[p(y|x_*, \hat{\theta}), p(y|x_* + r_{\text{vadv}}, \theta)\right]$$

$$r_{\text{vadv}} := \arg\max_{r; \|r\|_2 \leq \epsilon} D\left[p(y|x_*, \hat{\theta}), p(y|x_* + r)\right],$$

$$\mathcal{R}_{\text{vadv}}(\mathcal{D}_l, \mathcal{D}_{ul}, \theta) := \frac{1}{N_l + N_{ul}} \sum_{x_* \in \mathcal{D}_l, \mathcal{D}_{ul}} \text{LDS}(x_*, \theta)$$

$$\ell(\mathcal{D}_l, \theta) + \alpha \mathcal{R}_{\text{vadv}}(\mathcal{D}_l, \mathcal{D}_{ul}, \theta)$$