# Meta Learning Framework.

A distribution over tasks $P(\Upsilon)$

$\Upsilon$ , training data : $D_\Upsilon^{tr} = \{ X_\Upsilon^{tr} \ Y_\Upsilon^{tr} \ L_\Upsilon \}$

validation data : $D_\Upsilon^{val} = \{ X_\Upsilon^{val} \ Y_\Upsilon^{tr} \ L_T \}$

meta-learning , inner loop : operates within a task $\Upsilon$

outer loop : operates across a task $\Upsilon$

w belongs to $\varphi$
not task adaptive

$\Uparrow$
Predictor
$P_{\theta, w}$
meta-parameters
$\varphi$

**Aim:** design a meta-learner that can generalize well on a new task by appropriately choosing the predictor's task adaptive parameters $\theta$ after observing $D_{\Upsilon i}^{tr}$

For each training task $\Upsilon_i$:

$$P_{\theta, w} : X_{\Upsilon_i}^{val} \to \hat{Y}_{\Upsilon_i}^{val} \quad \text{conditioned on} \quad D_{\Upsilon_i}^{tr}$$

The meta-parameter $\varphi$ are updated in the outer loop so as to obtain good generalization in the inner loop.

minimize $\mathbb{E}_\Upsilon \ L_\Upsilon ( \hat{Y}_\Upsilon^{val}, Y_\Upsilon^{val} )$

Training on multiple tasks enables the meta-learner to produce $P_{\theta, w}$ that generalize well on a set of unseen tasks $\{\Upsilon_i\}$ sampled from $p(\Upsilon)$

Meta-Learning Ingredients.

(1) $\varphi = (I_{t_0}, w, u)$   meta parameters

(2) $I_{t_0}$   initialization function

(3) $U_u$   update function.

task meta-data.

Initialization function   $I_{t_0}(D_{T_i}^{tr}, c_{T_i})$

defines. the Initial values. of $\theta$ for a given $T_i$

* task meta-data $c_{T_i}$ may have a form of task ID

or a texture description.

Update function   $U_u(\theta_{l-1}, D_{T_i}^{tr})$

defines an iterated update to predictor parameters $\theta$

at iteration $l$.

The initialization and update functions produces a sequence

of predictor parameter.   $\theta_{0:l} \equiv \{\theta_0 \text{ -- } \theta_{l-1}, \theta_l\}$

We let the final predictor be a function of the whole

sequence of parameters.   $P_{\theta_{0:l}, w}$

eq. ① $P_{\theta_{0:l}, w}(\cdot) = \sum_{j=0}^{l} w_j P_{\theta_j, w}(\cdot)$

② $P_{\theta_{0:l}, w}(\cdot) = P_{\theta_l, w}(\cdot)$

Parameter $\theta$; Meta-parameters $\varphi = (\theta_0, w, u)$

Inner Loop: $\theta_0 \leftarrow I_{t_0}(D_{T_i}^{tr}, c_{T_i})$

$\qquad\qquad \theta_l \leftarrow U_u(\theta_{l-1}, D_{T_i}^{tr}) \quad \forall l > 0$

Prediction at $x$: $P_{\theta:l,w}(x)$

Outer Loop: $\varphi \leftarrow \varphi - \eta \cdot \nabla_\varphi L_{T_i} [P_{\theta:l,w}(x_{T_i}^{val}), y_{T_i}^{val}]$


## N-Beats as a Meta-learning Algorithm.

$$\hat{y} = \sum_{l=1}^{\infty} \hat{y_l}$$

$$\hat{y_l} = g \circ f(x)$$

$$\hat{y_l} = g \circ f(x_{l-1} - \hat{x_{l-1}}) \qquad l > 1$$

$$\hat{x_{l-1}} = g \circ f(x_{l-1})$$

$$\hat{y} = g \circ f(x) + \sum_{l > 1} g \circ f(x_{l-1} - g \circ f(x_{l-1})) \qquad (7)$$

(i) each application of $g \circ f$ in (7) is a predictor and
(ii) each block of N-Beats is the iteration of the inner meta-learning loop.

$$P_{\theta, w}(\cdot) = g_{w_g} \circ f_{w_f, \theta}(\cdot)$$

$w = (w_g, w_f)$ learned across tasks in the outer loop.

Task-specific parameters $\theta$ consist of the sequence
of input shift vectors $\theta \equiv \{\mu_l\}_{l=0}^{L}$

defined such that the $l$-th block takes input $x_l = x - \mu_{l-1}$

$$\mu_{l-1} = x - x_l \qquad\qquad \text{for } N\text{-Beats.} \quad x_{l+1} = x_l - \hat{x}_l$$

$$\mu_l = x - x_{l+1} \qquad\qquad\qquad\qquad x_l - x_{l+1} = \hat{x}_l \quad ②$$

$$\mu_l = \mu_{l-1} + x_l - x_{l+1} \qquad ①$$

From ① ② : $\quad \mu_l = \mu_{l-1} + \hat{x}_l$

This yields a recursive expression.

$$\mu_l \leftarrow \mathcal{U}_u (\mu_{l-1}, D_{Ti}^{tr})$$
$$\equiv \mu_{l-1} + g_{wg} \circ f_{wf} (x - \mu_{l-1}) \qquad D_{Ti}^{tr} \equiv \{x\}$$

$\theta = \{\mu_i\}_{i=0}^{\perp}$ are combined in the final output.

$$P_{\mu_{0:L}, w}(\cdot) = \sum_{j=0}^{L} w_j \cdot P_{\mu_j, w}(\cdot) \qquad w_j = 1 \; \forall j.$$

Conclusion: Even if predictor parameters $w_g$ $w_f$ are shared across blocks and fixed. The behavior of. $P_{\mu_{0:L}, w}(\cdot) = g_{wg} \circ f_{wf, \mu_{0:L}}(\cdot)$ is governed by $(w, \mu_1, \mu_2 \cdots)$

Therefore, the expressive power of the architecture. can. be expected to grow with the growing number of blocks. in proportion to the growth of. the space spanned by $\mu_{0:L}$

I think this claim is wrong for the size of $\theta$ is nearly negligible compared with $w = (w_g, w_f)$