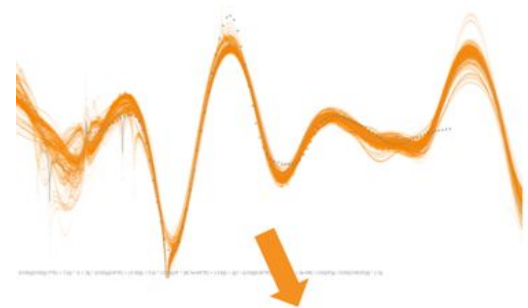# Meta-learning for Multi-task Symbolic Regression

6.883 Final Project
12/8/20
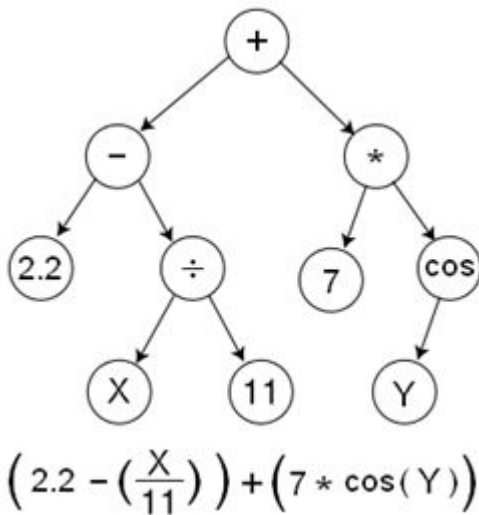Sam Kim and Ileana Rugina

# Symbolic Regression



$$e^{-x} x^3 \cos(x) \sin(x) (\cos(x) \sin(x)^2 - 1)$$

Offers interpretable results that can extrapolate



$$\left(2.2 - \left(\frac{X}{11}\right)\right) + \left(7 * \cos(Y)\right)$$

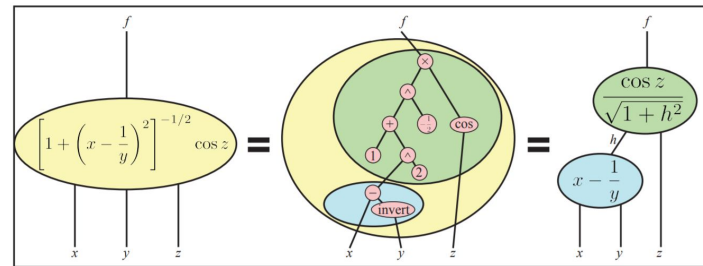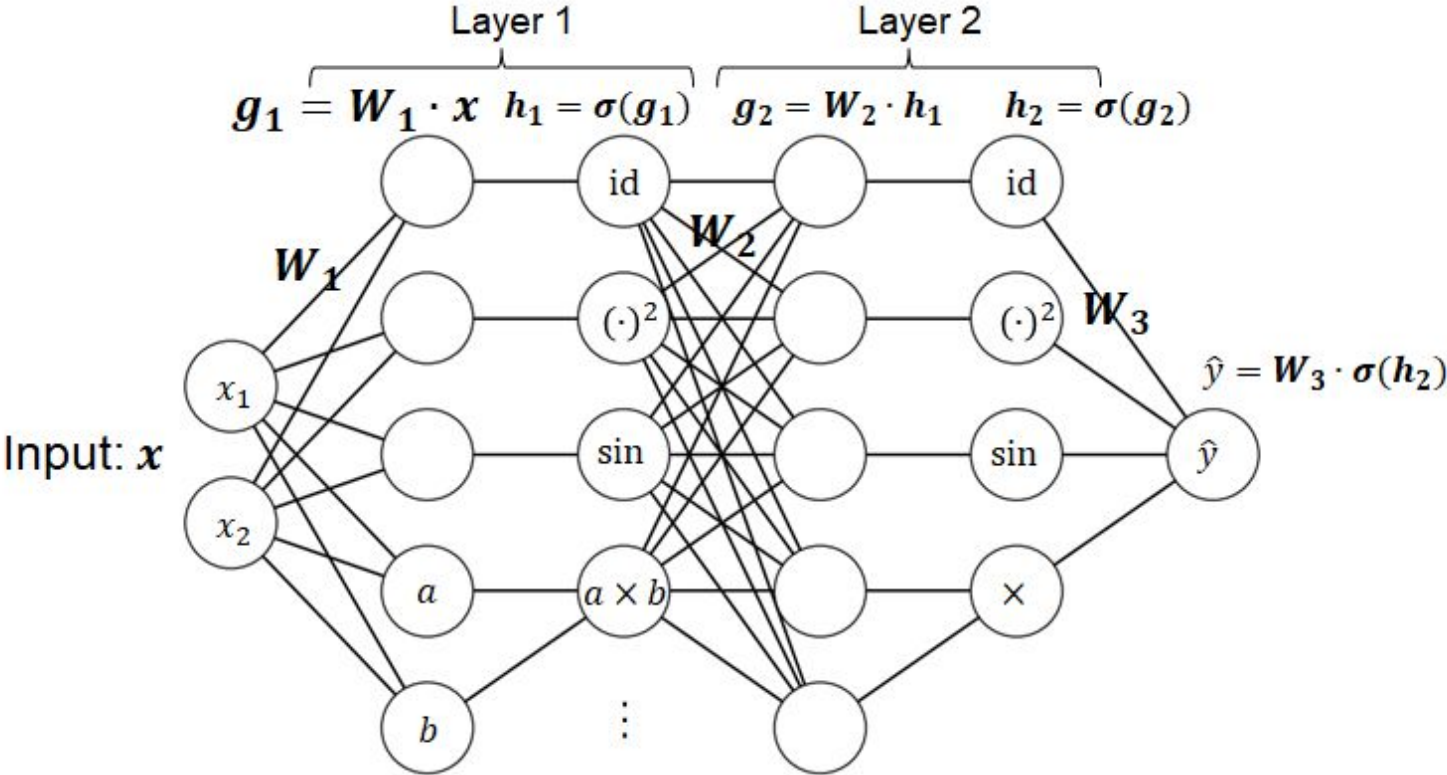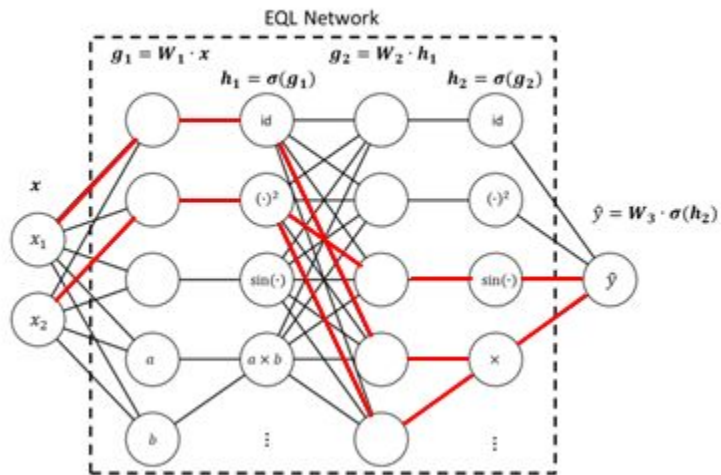Typically implemented using genetic programming



Figure 2: All functions can be represented as tree graphs whose nodes represent a set of basic functions (middle panel). Using a neural network trained to fit a mystery function (left panel), our algorithm seeks a decomposition of this function into others with fewer input variables (right panel), in this case of the form $f(x, y, z) = g[h(x, y), z]$.

AI Feynman 2.0
Uses neural networks to discover symmetries and simplify the problem

# Equation Learner (EQL) Network
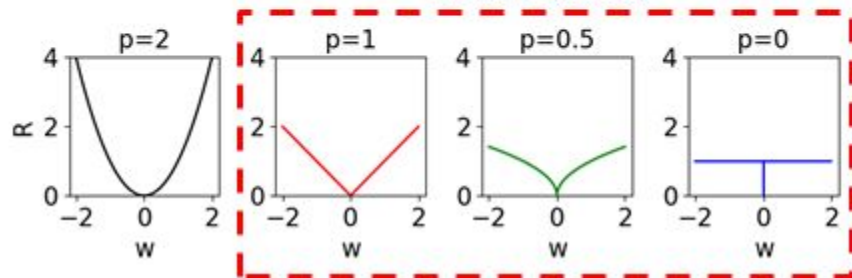
# Key Ingredient: Sparsity



EQL Network

$g_1 = W_1 \cdot x$  $g_2 = W_2 \cdot h_1$

$h_1 = \sigma(g_1)$  $h_2 = \sigma(g_2)$

id  id

$x$

$x_1$  $(\cdot)^2$  $(\cdot)^2$  $\hat{y} = W_3 \cdot \sigma(h_2)$

$x_2$  $\sin(\cdot)$  $\sin(\cdot)$  $\hat{y}$

$a$  $a \times b$  $\times$

$b$

## Sparsity through regularization

Loss function $= \text{MSE} + R$

$$R = \sum |w|^p$$
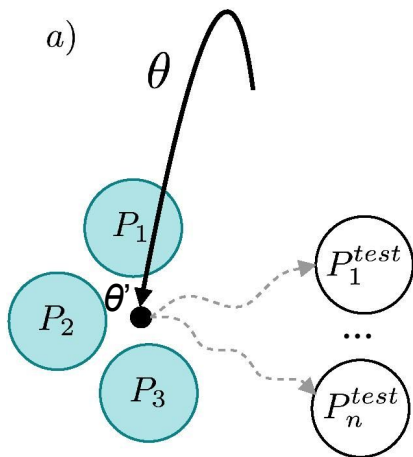
Promote sparsity

p=2    p=1    p=0.5    p=0

# Meta-Learning Methods:
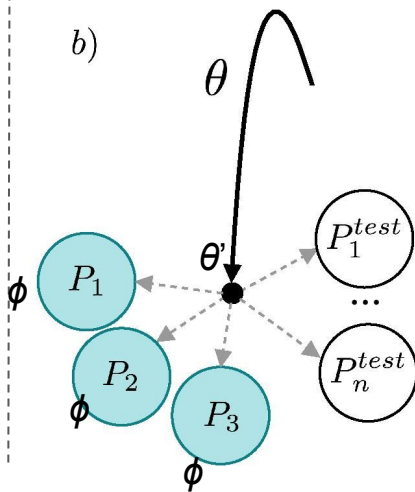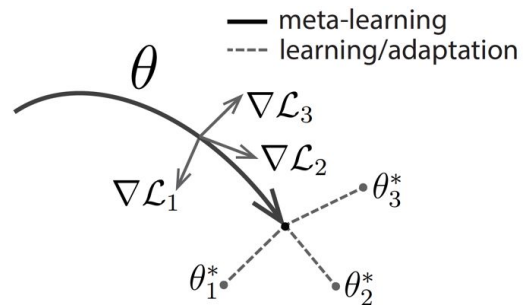
- ## Joint Training

Train on all functions "simultaneously" - pass in data from each function sequentially and perform parameter updates.

<u>No</u> second order terms or inductive bias towards good initialization

- ## Model-Agnostic Meta-Learning (MAML)

Find a good initialization by "fine-tuning" on each function during training to get θ'. Propagates second-order derivatives from θ' to the initial θ.



$a)$
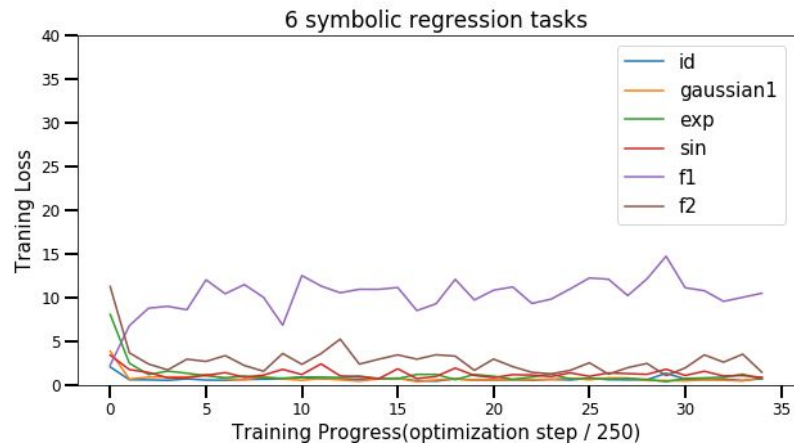


$b)$



**Algorithm 1** Model-Agnostic Meta-Learning

**Require:** $p(\mathcal{T})$: distribution over tasks
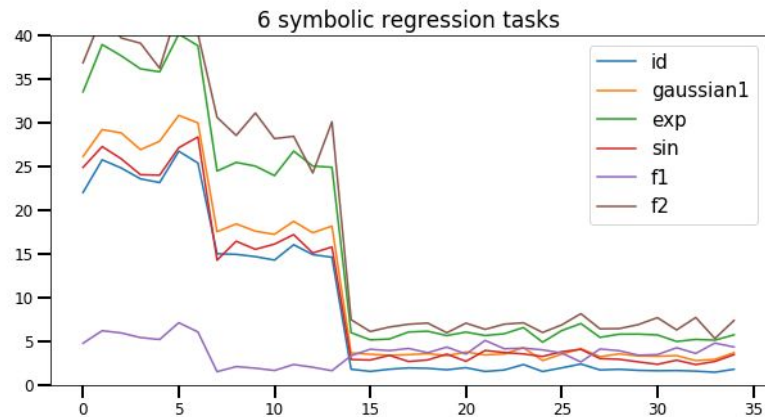**Require:** $\alpha, \beta$: step size hyperparameters
1: randomly initialize $\theta$
2: **while** not done **do**
3:     Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
4:     **for all** $\mathcal{T}_i$ **do**
5:         Evaluate $\nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$ with respect to $K$ examples
6:         Compute adapted parameters with gradient descent: $\theta_i' = \theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$
7:     **end for**
8:     Update $\theta \leftarrow \theta - \beta \nabla_\theta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta_i'})$
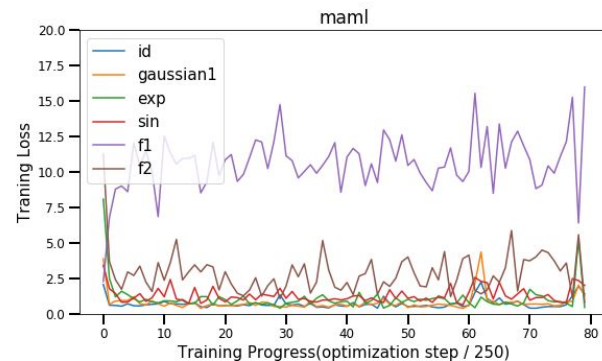9: **end while**

# Initial Results: Regularization Method
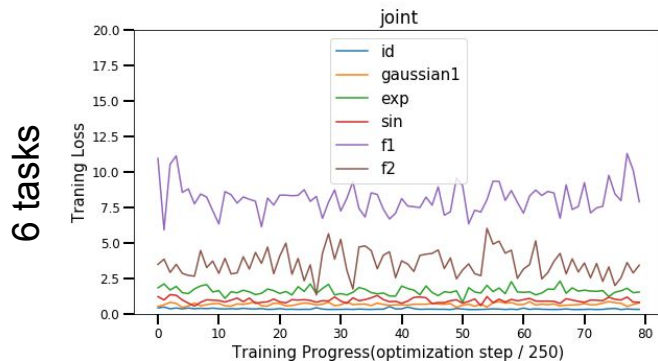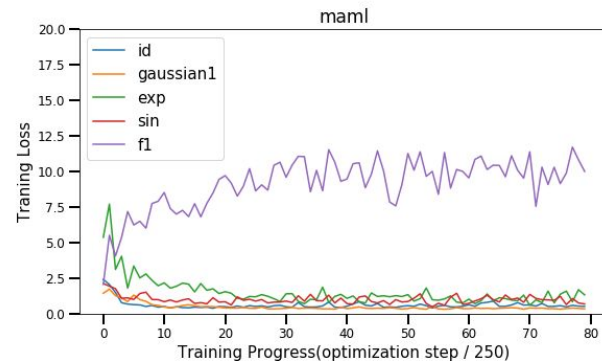
$L_0$ sparsity regularization

$L_{0.5}$ sparsity regularization

# Initial Results: Method Comparison

| ID | formula |
|---|---|
| "id" | $f(x) = x$ |
| "gaussian1" | $f(x) = \frac{\exp(-x^2/2)}{\sqrt{2*\pi}}$ |
| "exp" | $f(x) = \exp(x)$ |
| "sin" | $f(x) = \sin x$ |
| "f1" | $f(x) = x \cdot \sin x - 3$ |
| "f2" | $f(x) = x^2 + 3 \cdot x + 1$ |



$L_0$ regularization

# Distribution of Functions: Damped Harmonic Oscillator



Spring force

-kx

Equilibrium position

Linear damping force

-cv

viscous damping

Velocity

v

x

$$F = -kx - c\frac{\mathrm{d}x}{\mathrm{d}t} = m\frac{\mathrm{d}^2 x}{\mathrm{d}t^2},$$

$$\frac{\mathrm{d}^2 x}{\mathrm{d}t^2} + 2\zeta\omega_0\frac{\mathrm{d}x}{\mathrm{d}t} + \omega_0^2 x = 0,$$



The underdamped response of the oscillator is described by the equation:

$$x = e^{-\gamma t} a\cos\left[\omega_1 t - \alpha\right]$$

Overdamped

Critical Damping

One-half of critical damping

One-tenth of critical damping

Oscillator with resonant frequency 10 rad/s started from rest. After Barger & Olsson

# Results: Drawing from the Function Distribution

Function generation process:
Sample *C* ~ Bernoulli(0.5)
- If C = 0
  - Sample ω ~ Unif(0.5, 2)*2π
  - Sample φ ~ Unif(0, 2π)
  - Function:
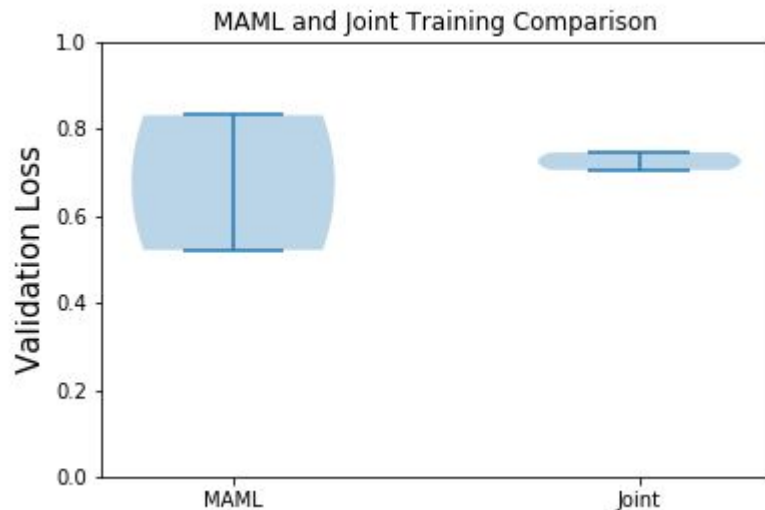    - $f(x) = \sin(\omega x + \phi)$

- Else
  - Sample a,b ~ Unif(-1, 1)
  - Function:
    - $f(x) = a\square e^x + b\square e^x$



MAML and Joint Training Comparison

| MAML | Joint Training |
|------|----------------|
| 0.68 | 0.72 |

# Out of Domain Experiments



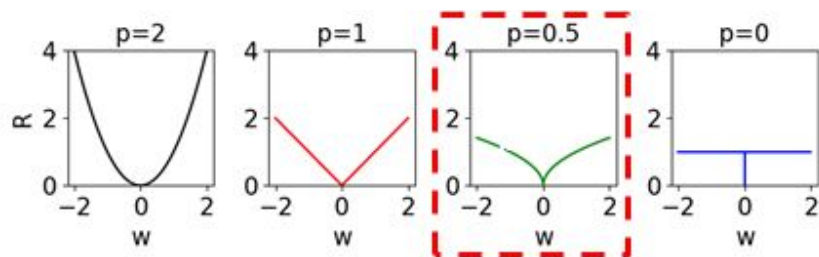| MAML | Joint Training |
|------|----------------|
| 8.76 | 6.75 |

# Conclusion

Demonstrate meta-learning on the EQL network for symbolic regression

When functions are drawn from the same distribution (which is often the case in certain fields of science and engineering), meta-learning improves performance and generalizability of the EQL network

MAML outperforms joint training on average for test functions drawn from same distribution as training
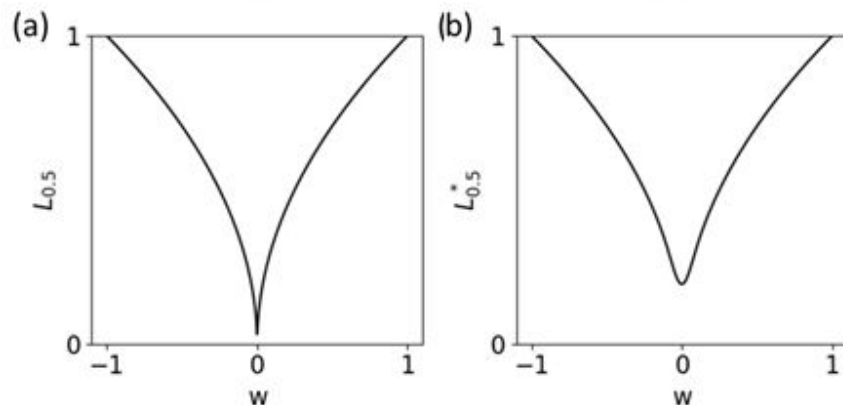
# Smoothed L0.5 Sparsity



$L_{0.5}$
- Promotes sparsity stronger than $L_1$
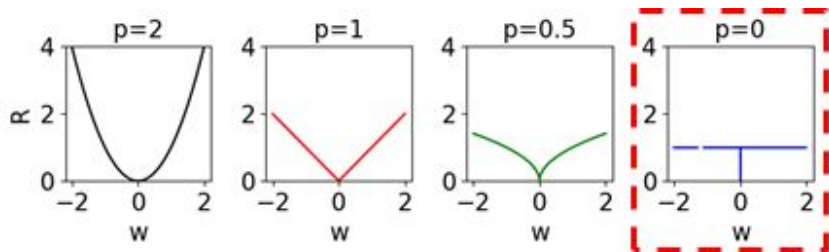- Downsides:
  - Non-convex
  - Infinite gradient

Smoothed $L_{0.5}$ regularization
- Avoids infinite gradient

$$L_{0.5}^*(w) = \begin{cases} |w|^{1/2} & w \geq a \\ \left(-\dfrac{w^4}{8a^3} + \dfrac{3w^2}{4a} + \dfrac{3a}{8}\right)^{1/2} & w < a \end{cases}$$

Fan, Qinweiet, et al.. "Convergence of online gradient method for feedforward neural networks with smoothing L1/2 regularization penalty." *Neurocomputing* 131 (2014): 208-216.

# Relaxed L0 Sparsity



$L_0$

- Promotes sparsity without penalizing magnitude
- Downside: non-differentiable
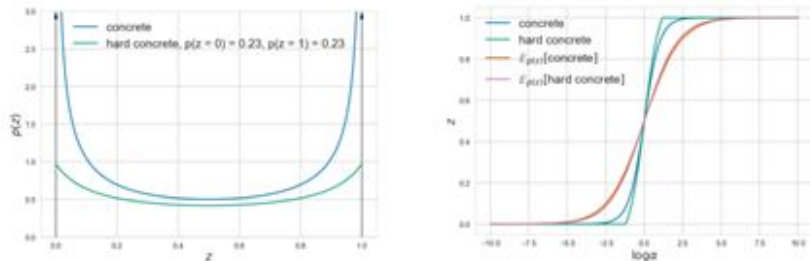
Relaxed $L_0$ regularization
- Can calculate gradients for backprop

Reparameterize weights as

$$\Theta = \tilde{\Theta} \odot z$$

where $z$ is a stochastic variable drawn from the Hard Concrete distribution

Louizos, Christos, Max Welling, and Diederik P. Kingma. "Learning Sparse Neural Networks through $ L_0 $ Regularization." *arXiv preprint arXiv:1712.01312* (2017).