# Training NNs via the Augmented Lagrangian Method (ALM)

NN model:

$$\Phi(W_0, \ldots, W_N; z) = W_N \Phi(W_{p-1} \cdots \Phi(W_0 x) \cdots)$$

NN training: Given data $(z_0^j, y^j)$, $j = 1, \ldots, m$, solve

$$\min_{W_0, \ldots, W_N} \frac{1}{2} \sum_{j=1}^{m} \| \Phi(W_0, \ldots, W_N; z_0^j) - y^j \|_2^2$$

Introducing dynamics $z_{a+1} = \Phi(W_a z_a)$ we can express it as an optimal control problem

$$\min \frac{1}{2} \sum_{j=1}^{m} \| W_N z_N^j - y^j \|_2^2$$

$$\text{s.t.} \quad z_{a+1}^j = \Phi(W_a z_a^j), \quad a = 0, \ldots, N-1, \quad j = 1, \ldots, m$$

In abstract terms

$$\min \frac{1}{2} \| F(x) \|_2^2$$

$$\text{s.t.} \quad h(x) = 0$$

where $F(x) = (W_N z_N^1 - y^1, \ldots, W_N z_N^m - y^m)$

$$h(x) = (h^1(x), \ldots, h^m(x))$$

$$h^j(x) = (z_1^j - \Phi(W_0 z_0^j), \ldots, z_N^j - \Phi(W_{N-1} z_{N-1}^j))$$

$$x = (W_0, \ldots, W_N, (z_1^j, \ldots, z_N^j)_{j=1}^m)$$

$$\min \frac{1}{2}\|F(x)\|_2^2$$
$$\text{s.t.} \quad h(x) = 0$$

## Augmented Lagrangian

$$\mathcal{L}_\beta(x,y) = \frac{1}{2}\|F(x)\|_2^2 + \langle y, h(x)\rangle + \frac{\beta}{2}\|h(x)\|_2^2$$

$$= \frac{1}{2}\|F(x)\|_2^2 + \frac{\beta}{2}\|h(x) + y/\beta\|_2^2 - \frac{1}{2\beta}\|y\|_2^2$$

$$= \frac{\beta}{2}\left\|\begin{bmatrix} F(x)/\sqrt{\beta} \\ h(x) + y/\beta \end{bmatrix}\right\|^2$$

## ALM: Given $y^0$, iterate

1. Find $x^a$ s.t. $\|\nabla_x \mathcal{L}_{\beta_a}(x^a, y^a)\|_2 \le \epsilon_a$

2. Update multipliers $y^{a+1} = y^a + \beta_a h(x^a)$

Step 1 amounts to solving a nonlinear least-squares problem:
$$\min \frac{1}{2}\|F_{\beta_a}(x, y^a)\|_2^2$$
where $F_\beta(x,y) = \begin{bmatrix} F(x)/\sqrt{\beta} \\ h(x) + y/\beta \end{bmatrix}$

## Can use Levenberg-Marquardt (LM)

### Questions:

1. How can we solve the linear least-squares problem efficiently in the LM method? It has a lot of structure: exploit it!

2. How can we adjust penalty parameters $\beta_a$ and tolerances $\epsilon_a$? Check the literature!

3. Other methods for solving step 1?