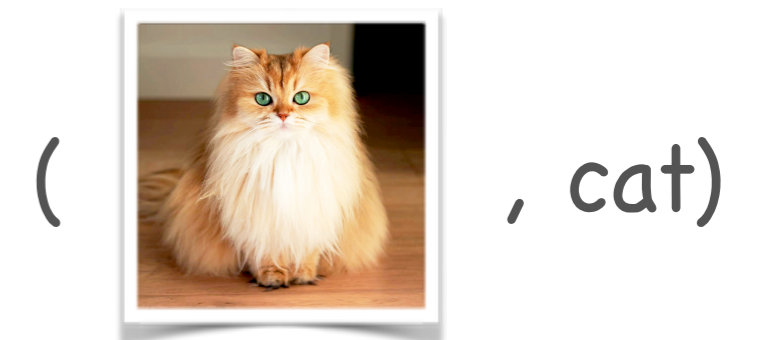


Optimization of deep networks

© 2019 Philipp Krähenbühl and Chao-Yuan Wu

Data

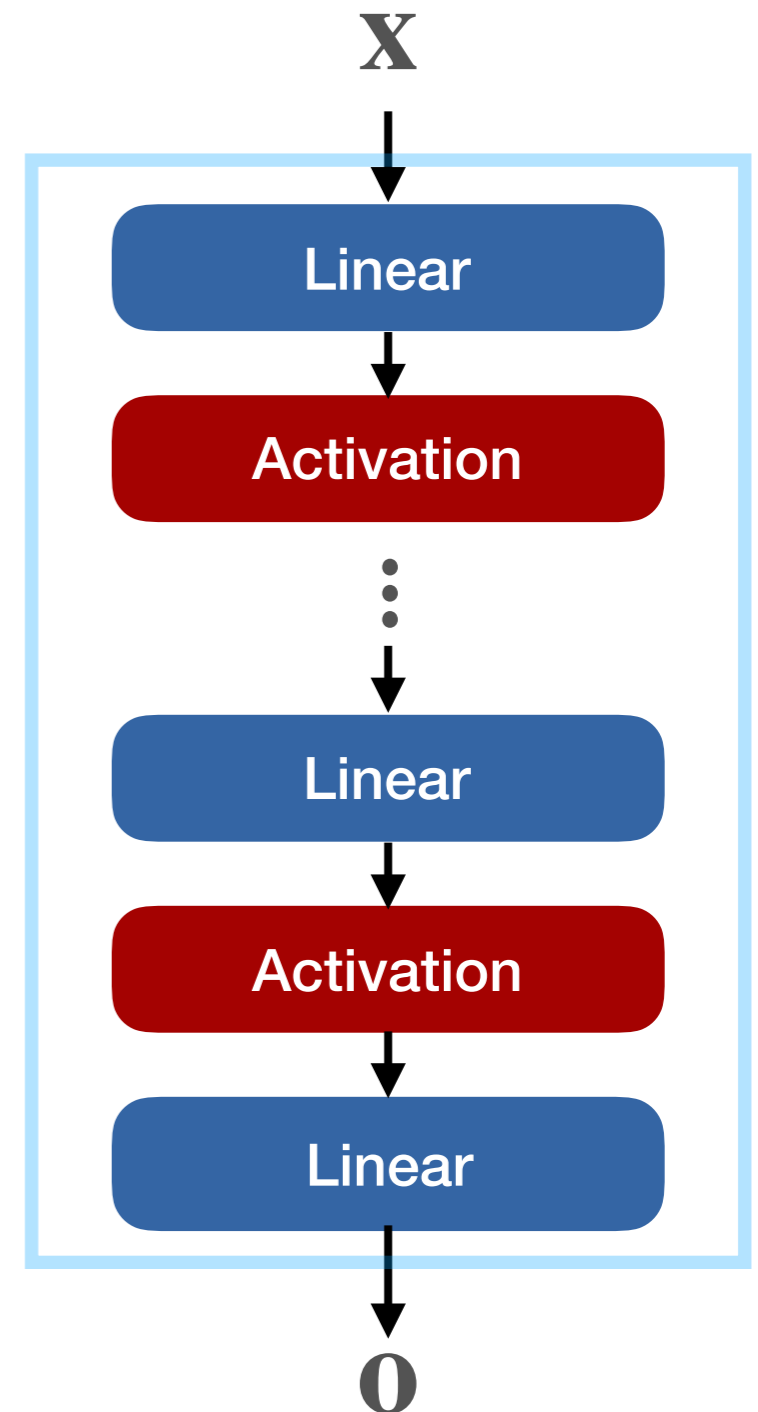
- Input: $\{\mathbf{x}_0, \dots, \mathbf{x}_{N-1}\}$
- Label: $\{y_0, \dots, y_{N-1}\}$
- Dataset: $D = \{(\mathbf{x}_0, y_0), \dots, (\mathbf{x}_{N-1}, y_{N-1})\}$



⋮

Model

- Deep network $f: (\mathbf{x}, \theta) \rightarrow \mathbf{o}$
- Layers of computation
- Parameters θ
- Differentiable computation graph



Loss

- Differentiable $\ell(\mathbf{o}, \mathbf{y})$

- Regression

- Distance norm

$$\ell(\mathbf{o}, \mathbf{y}) = \|\mathbf{o} - \mathbf{y}\|$$

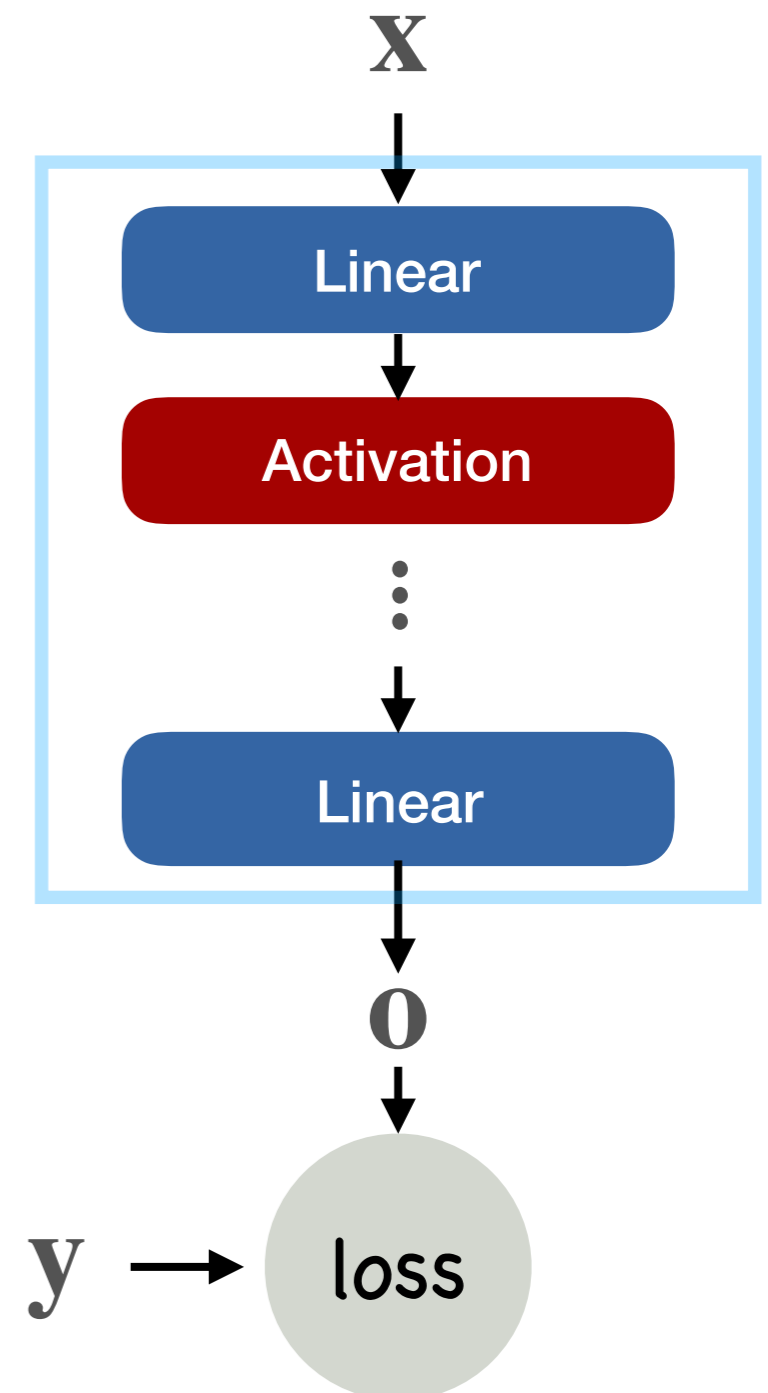
- Classification

- Cross Entropy

$$\ell(\mathbf{o}, \mathbf{y}) = -\log p(y)$$

- Over training dataset

- $L(\theta) = \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim D}[\ell(f(\mathbf{x}, \theta), \mathbf{y})]$



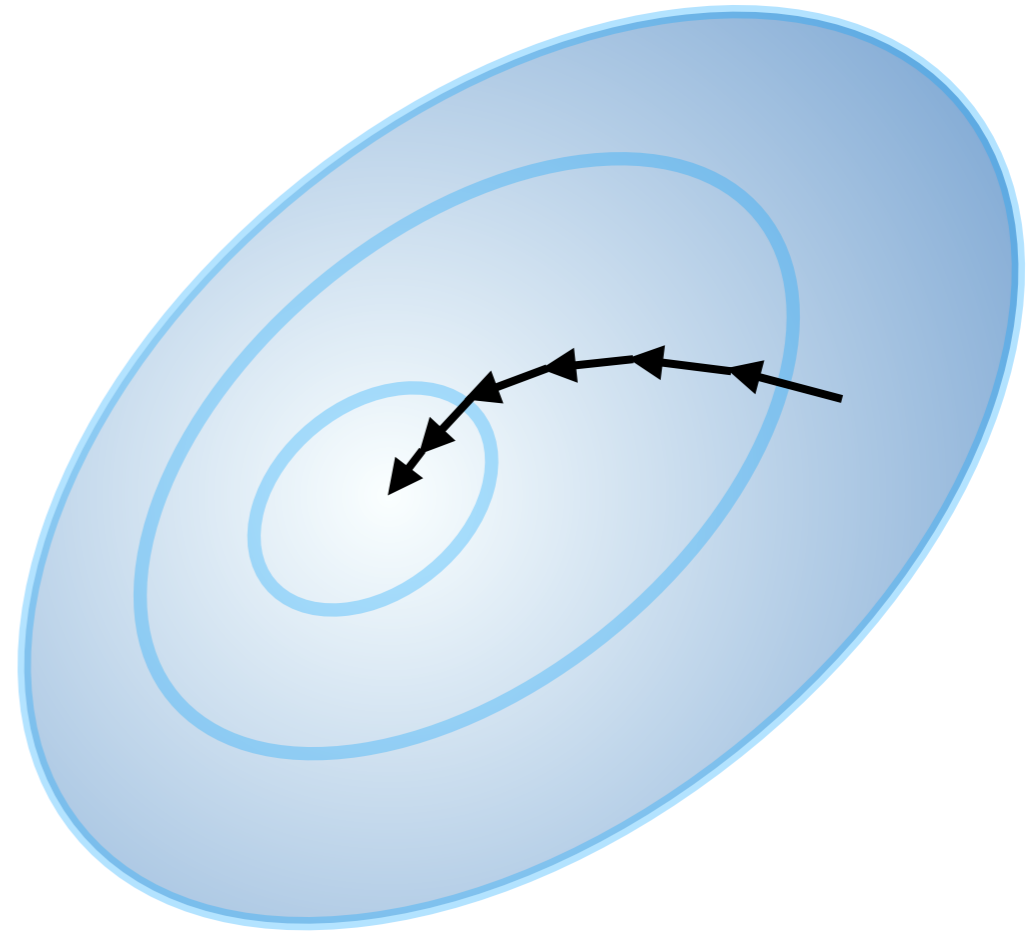
Optimization

- Minimize $L(\theta)$

Gradient Descent

- Repeat until convergence:

- $\theta := \theta - \epsilon \frac{dL(\theta)}{d\theta}$



Issue with Gradient Descent

- Slow to compute gradient

- $$\frac{dL(\theta)}{d\theta} = \mathbb{E}_{\mathbf{x}, \mathbf{y} \in D} \left[\frac{d\ell(f(\mathbf{x}, \theta), \mathbf{y})}{d\theta} \right]$$