# Neural networks

Training neural networks - loss function

# MACHINE LEARNING

**Topics:** stochastic gradient descent (SGD)

- Algorithm that performs updates after each example
  - ‣ initialize $\boldsymbol{\theta}$    ( $\boldsymbol{\theta} \equiv \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \ldots, \mathbf{W}^{(L+1)}, \mathbf{b}^{(L+1)}\}$ )
  - ‣ for N iterations
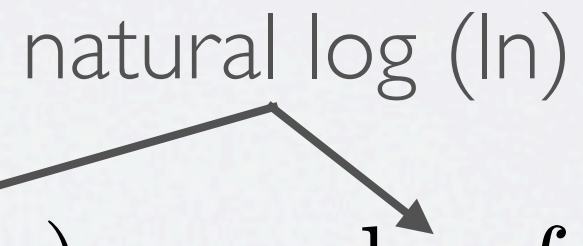    - for each training example $(\mathbf{x}^{(t)}, y^{(t)})$
      - ✓ $\Delta = -\nabla_{\boldsymbol{\theta}} l(f(\mathbf{x}^{(t)}; \boldsymbol{\theta}), y^{(t)}) - \lambda \nabla_{\boldsymbol{\theta}} \Omega(\boldsymbol{\theta})$
      - ✓ $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \, \Delta$

$$\left.\begin{array}{c} \\ \\ \end{array}\right\}$$
training epoch
=
iteration over **all** examples

- To apply this algorithm to neural network training, we need
  - ‣ the loss function $l(\mathbf{f}(\mathbf{x}^{(t)}; \boldsymbol{\theta}), y^{(t)})$
  - ‣ a procedure to compute the parameter gradients $\nabla_{\boldsymbol{\theta}} l(\mathbf{f}(\mathbf{x}^{(t)}; \boldsymbol{\theta}), y^{(t)})$
  - ‣ the regularizer $\Omega(\boldsymbol{\theta})$  (and the gradient $\nabla_{\boldsymbol{\theta}} \Omega(\boldsymbol{\theta})$ )
  - ‣ initialization method

# LOSS FUNCTION

**Topics:** loss function for classification

- Neural network estimates $f(\mathbf{x})_c = p(y = c|\mathbf{x})$

  ‣ we could maximize the probabilities of $y^{(t)}$ given $\mathbf{x}^{(t)}$ in the training set

- To frame as minimization, we minimize the negative log-likelihood

natural log (ln)

$$l(\mathbf{f}(\mathbf{x}), y) = -\sum_c 1_{(y=c)} \log f(\mathbf{x})_c = -\log f(\mathbf{x})_y$$

  ‣ we take the log to simplify for numerical stability and math simplicity

  ‣ sometimes referred to as cross-entropy