# Neural networks

Training neural networks - regularization

# 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

# REGULARIZATION

**Topics:** L2 regularization

$$\Omega(\boldsymbol{\theta}) = \sum_k \sum_i \sum_j \left( W_{i,j}^{(k)} \right)^2 = \sum_k ||\mathbf{W}^{(k)}||_F^2$$

- Gradient: $\nabla_{\mathbf{W}^{(k)}} \Omega(\boldsymbol{\theta}) = 2\mathbf{W}^{(k)}$

- Only applied on weights, not on biases (weight decay)

- Can be interpreted as having a Gaussian prior over the weights

# REGULARIZATION

**Topics:** L1 regularization

$$\Omega(\boldsymbol{\theta}) = \sum_k \sum_i \sum_j |W_{i,j}^{(k)}|$$

- Gradient: $\nabla_{\mathbf{W}^{(k)}} \Omega(\boldsymbol{\theta}) = \text{sign}(\mathbf{W}^{(k)})$

  ‣ where $\text{sign}(\mathbf{W}^{(k)})_{i,j} = 1_{\mathbf{W}_{i,j}^{(k)} > 0} - 1_{\mathbf{W}_{i,j}^{(k)} < 0}$

- Also only applied on weights

- Unlike L2, L1 will push certain weights to be exactly 0

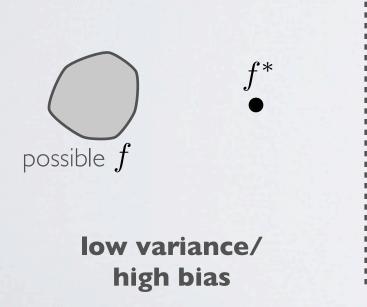- Can be interpreted as having a Laplacian prior over the weights
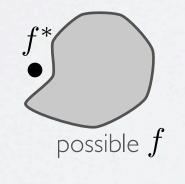
# MACHINE LEARNING

**Topics:** bias-variance trade-off

- Variance of trained model: does it vary a lot if the training set changes

- Bias of trained model: is the average model close to the true solution

- Generalization error can be seen as the sum of the (squared) bias and the variance



possible $f$

$f^*$

possible $f$

**low variance/
high bias**

**good trade-off**

possible $f$

$f^*$

$f^*$

possible $f$

**high variance/
low bias**