UNIVERSITY OF AMSTERDAM
Informatics Institute

AMLAB
Amsterdam
Machine Learning Lab

# Machine Learning 1

Lecture 7.3 - Supervised Learning
Classification - Logistic Regression:
Stochastic Gradient Descent

*Erik Bekkers*

*(Bishop 4.3.2)*

*Slide credits: Patrick Forré and
Rianne van den Berg*

# Logistic Regression for Two Classes

▸ Given: Dataset $\mathbf{X} = (\mathbf{x}_1, ..., \mathbf{x}_N)^T$ with binary targets

$$\mathbf{t} = (t_1, ..., t_N)^T \quad \text{with} \quad t_n \in \{\mathcal{C}_1, \mathcal{C}_2\} = \{1, 0\}$$

▸ Conditional likelihood function:

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}) = \prod_{n=1}^{N} p(\mathbf{t}_n|\mathbf{x}_n, \mathbf{w}) = \prod_{n=1}^{N} y_n^{t_n} (1 - y_n)^{1-t_n}$$

$$y_n = p(C_1|\phi_n) = \sigma(w^T \phi_n) \qquad \phi_n = \phi(\mathbf{x}_n)$$

▸ Maximizing the conditional likelihood/minimizing the cross-entropy

$$E(\mathbf{w}) = -\ln p(\mathbf{t}, \mathbf{X}, \mathbf{w}) = -\sum_{n=1}^{N} t_n \ln y_n + (1 - t_n)\ln(1 - y_n)$$

▸ E($\mathbf{w}$): convex, but no closed form solution!

$$y_n = \sigma(\mathbf{w}^T \phi_n) \text{ is nonlinear in } \mathbf{w}$$

# Logistic Regression (K=2): SGD

$y_n = \sigma(\mathbf{w}^T \phi_n)$

‣ Stochastic Gradient Descent for cross-entropy:

$$E(\mathbf{w}) = -\sum_{n=1}^{N} t_n \ln y_n + (1 - t_n) \ln(1 - y_n) = \sum_{n=1}^{N} E_n(\mathbf{w})$$

‣ Update rule given a random data point $(\mathbf{x}_n, t_n)$

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w}^{(\tau)})^T$$

‣ Gradient: $\nabla E_n(\mathbf{w}) = \left( \dfrac{\partial E_n(\mathbf{w})}{\partial w_0}, ..., \dfrac{\partial E_n(\mathbf{w})}{\partial w_{M-1}} \right)$

‣ $\dfrac{\partial E_n(\mathbf{w})}{\partial w_j} = \dfrac{\partial E_n(\mathbf{w})}{\partial y_n} \dfrac{\partial y_n}{\partial w_j} = \left( -\dfrac{t_n}{y_n} + \dfrac{1-t_n}{1-y_n} \cdot \right) \cdot \dfrac{\partial y_n}{\partial w_j}$

‣ $\dfrac{\partial y_n}{\partial w_j} = \dfrac{\partial}{\partial w_j} \sigma(\mathbf{w}^T \phi_n)$

# Logistic Regression (K=2): SGD

- $\dfrac{\partial y_n}{\partial w_j} = \dfrac{\partial}{\partial w_j}\sigma(\mathbf{w}^T\boldsymbol{\phi}_n)$

$$\frac{\partial \underline{w}^T \phi_n}{\partial w_j} = \frac{\partial}{\partial w_j}\sum_{i=0}^{M-1} w_i \phi_{ni} = \phi_{nj}$$

- Use $\dfrac{\partial\sigma(a)}{\partial a} = \sigma(a)(1-\sigma(a))$

- $\dfrac{\partial}{\partial w_j}\sigma(\mathbf{w}^T\boldsymbol{\phi}_n) = \underbrace{\sigma(\underline{w}^T\phi_n)}_{y_n}(1 - \underbrace{\sigma(\underline{w}^T\phi_n)}_{y_n}) \cdot \dfrac{\partial \underline{w}^T \phi_n}{\partial w_j}$

$$= y_n(1-y_n)\phi_{nj}$$

- $\dfrac{\partial E_n(\mathbf{w})}{\partial w_j} = -\dfrac{t_n}{y_n}\dfrac{\partial y_n}{\partial w_j} + \dfrac{1-t_n}{1-y_n}\dfrac{\partial y_n}{\partial w_j}$

(activation fn = link fn)

Bishop 4.3.6

$$= -\frac{t_n}{y_n}\cdot y_n(1-y_n)\phi_{nj} + \frac{1-t_n}{1-y_n}\cdot y_n(1-y_n)\phi_{nj}$$

$$= -t_n\phi_{nj} + t_n y_n \phi_{nj} + y_n \phi_{nj} - t_n y_n \phi_{nj} = (y_n - t_n)\phi_{nj}$$

# Logistic Regression (K=2): SGD

‣ Stochastic Gradient Descent for cross-entropy:

$$E(\mathbf{w}) = -\sum_{n=1}^{N} t_n \ln y_n + (1 - t_n)\ln(1 - y_n) = \sum_{n=1}^{N} E_n(\mathbf{w})$$

‣ Update rule given a random data point $(\mathbf{x}_n, t_n)$

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w}^{(\tau)})^T$$

‣ $$\frac{\partial E_n(\mathbf{w})}{\partial w_j} = (y_n - t_n)\phi_j(\mathbf{x}_n)$$

‣ Gradient: $\nabla E_n(\mathbf{w})^T = \left( \dfrac{\partial E_n(\mathbf{w})}{\partial w_0}, ..., \dfrac{\partial E_n(\mathbf{w})}{\partial w_{M-1}} \right)^T = (y_n - t_n)\phi_n$

‣ Update rule:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \, (y_n - t_n)\phi_n$$

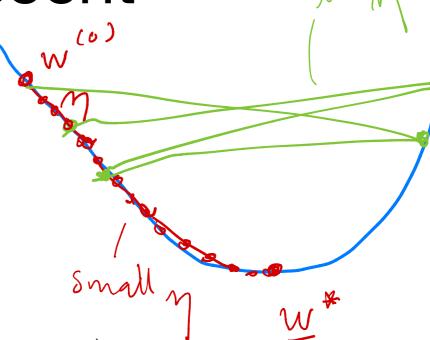*perceptron*

# Stochastic Gradient Descent

1. Initialize $\mathbf{w}^{(0)}$

2. Choose a learning rate $\eta$

3. While $||\mathbf{w}^{(\tau+1)} - \mathbf{w}^{(\tau)}|| > \varepsilon$

   I. Choose a random data point $(\mathbf{x}_n, t_n)$

   II. Update **w:**

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta\,(y_n^{(\tau)} - t_n)\boldsymbol{\phi}(\mathbf{x}_n)$$

‣ If $\eta$ too large: no convergence

‣ If $\eta$ too small: very slow convergence

‣ Converged **w\***: estimate of minimizer of E(**w**)!

*large $\eta$*

$w^{(0)}$

$\eta$

*small $\eta$*

$w^*$