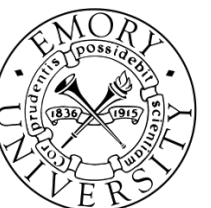


Gradient Descent

Natural Language Processing

Emory University

Jinho D. Choi



Supervised Learning

input output $y = \pm 1$ ← binomial distribution

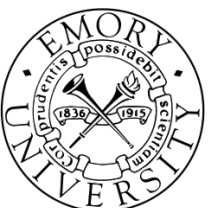
$$(X, Y) = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

prediction → $\hat{y} = f(x)$ ← predicts the output of x

Expected risk → $E(f) = \int \ell(\hat{y}; y) \cdot P(x, y)$

loss function joint distribution **unknown!**

Empirical risk → $\hat{E}(f) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{y}_i; y_i)$ ← minimize!



Linear Prediction

$$\hat{E}(f) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{y}_i; y_i)$$

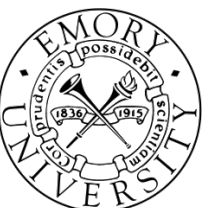
$$\ell(\hat{y}; y) = \frac{1}{2} (\hat{y} - y)^2 \leftarrow \text{least squares}$$

$$\hat{y} = f(x) = w^T \Phi(x) = w^T x \leftarrow \text{linear function}$$

↑
feature vector

$$\ell(w, x; y) = \frac{1}{2} (w^T x - y)^2$$

Find a **weight vector** that minimizes the **loss**.

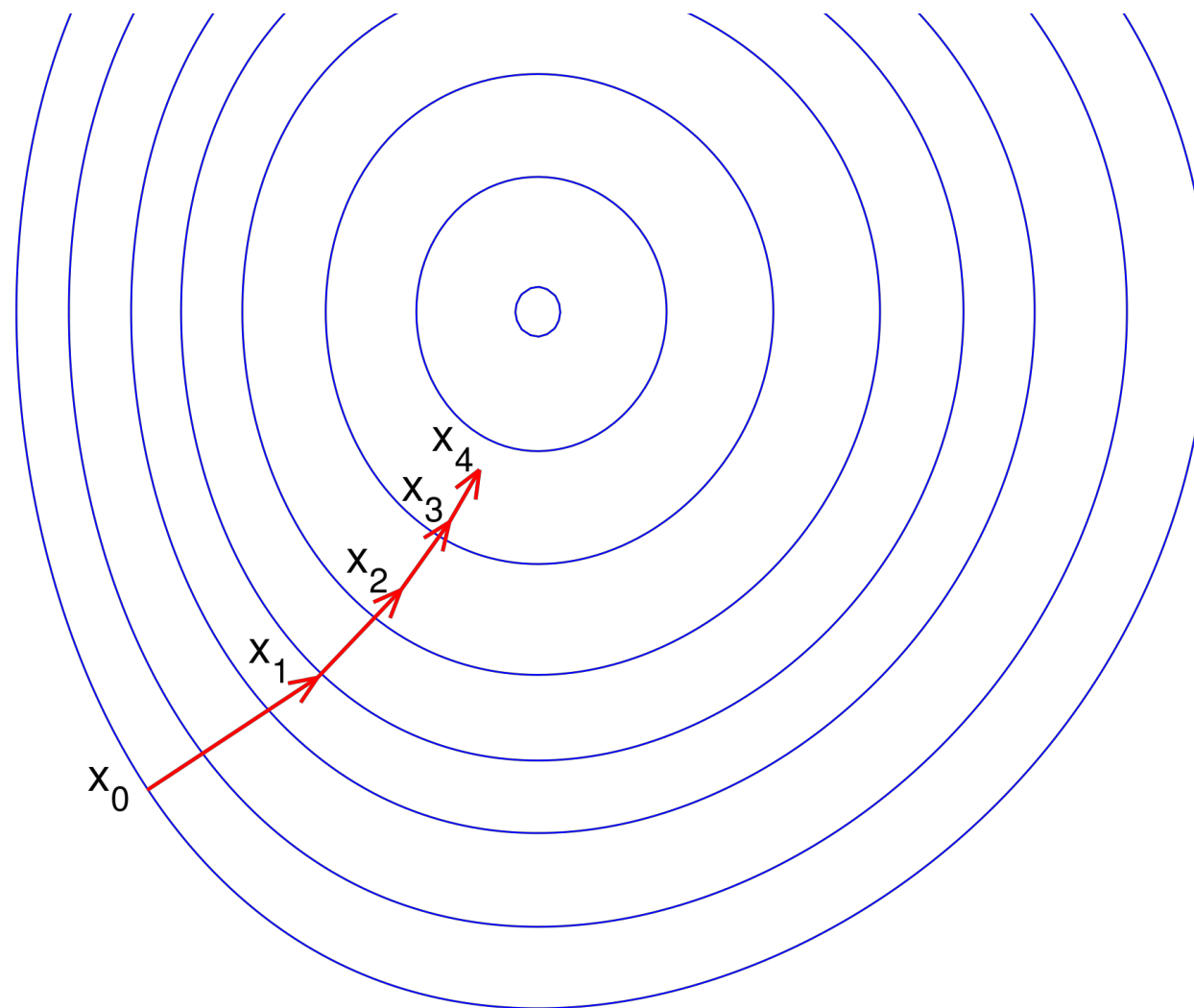


Gradient Descent

learning rate

derivative of the loss

$$w_{t+1} \leftarrow w_t - \eta_t \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial w} \ell(w_t, x_i; y_i)$$



Minimize **loss**

Derivative $\rightarrow 0$

Convex optimization

Global **optimum**?

Gradient Descent

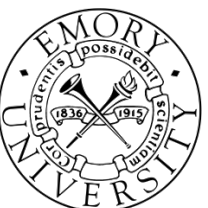
$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \mathbf{w}} \ell(\mathbf{w}_t, \mathbf{x}_i; y_i)$$

$$\ell(\mathbf{w}, \mathbf{x}; y) = \frac{1}{2} (\mathbf{w}^T \mathbf{x} - y)^2$$

$$\frac{\partial}{\partial \mathbf{w}} \ell(\mathbf{w}, \mathbf{x}; y) = \frac{\partial}{\partial \mathbf{w}} \frac{1}{2} (\mathbf{w}^T \mathbf{x} - y)^2 = (\mathbf{w}^T \mathbf{x} - y) \mathbf{x}$$

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i) \mathbf{x}_i$$

How often is the **weight vector** updated?



Stochastic Gradient Descent

$$w_{t+1} \leftarrow w_t - \eta_t \frac{1}{n} \sum_{i=1}^n (w^T x_i - y_i) x_i$$

$$w_{t+1} \leftarrow w_t - \eta_t (w_t^T x_i - y_i) x_i$$

updated for **every** instance

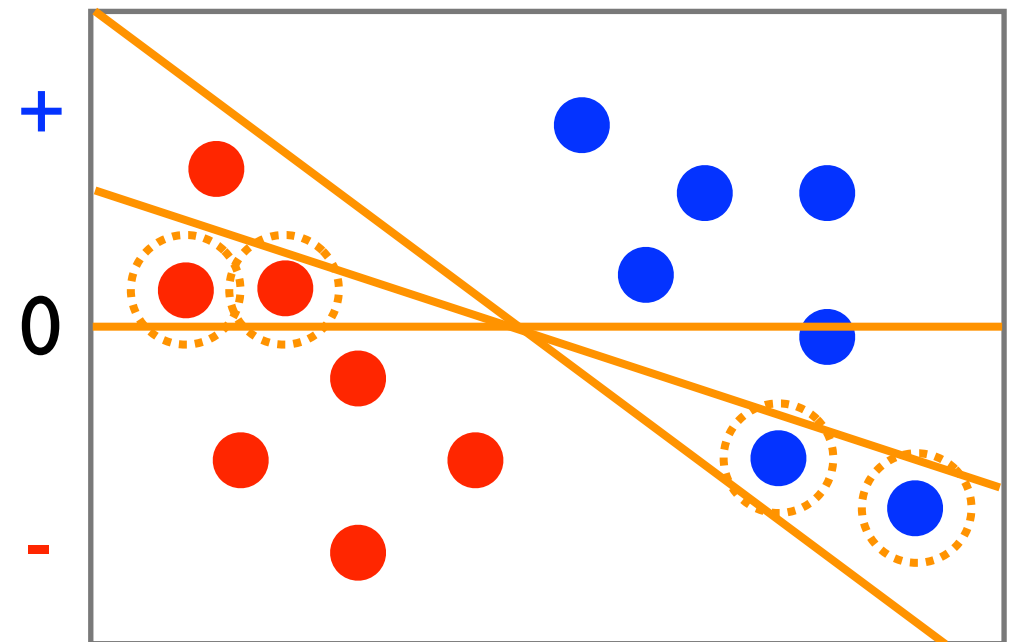
$$w_0 \leftarrow 0$$

$$w_0^T x_1 > 0 \quad w_1 \leftarrow w_0 - \eta(\oplus + 1)x_1$$

$$w_1^T x_2 < 0 \quad w_2 \leftarrow w_1 - \eta(\ominus + 1)x_2$$

$$w_2^T x_3 < 0 \quad w_3 \leftarrow w_2 - \eta(\ominus - 1)x_3$$

$$w_3^T x_4 > 0 \quad w_4 \leftarrow w_3 - \eta(\oplus - 1)x_4$$



Perceptron

Stochastic gradient descent

$$w_{t+1} \leftarrow w_t - \eta_t \Delta \ell$$

Least squares

$$\ell(w, x; y) = \frac{1}{2} (w^T x - y)^2$$

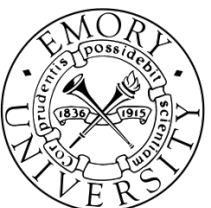
$$\Delta \ell = (w^T x - y) x$$

Perceptron

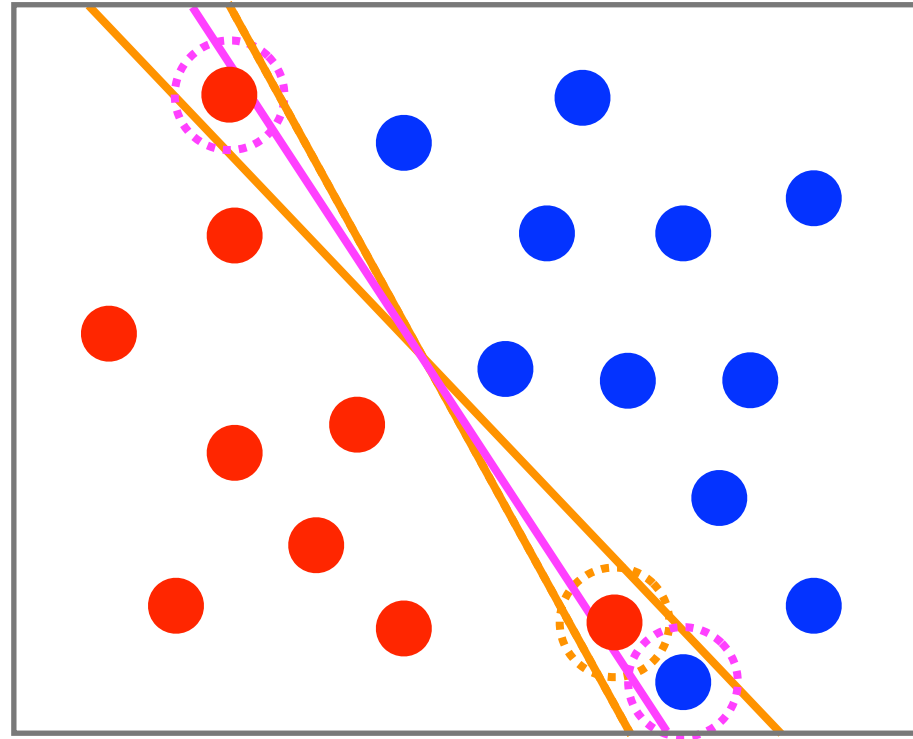
$$\ell(w, x; y) = \max\{0, -w^T x \cdot y\}$$

$$\Delta \ell = \begin{cases} -x \cdot y & w^T x \cdot y < 0 \\ 0 & \text{otherwise} \end{cases}$$

$$w_{t+1} \leftarrow w_t + \eta_t \begin{cases} x \cdot y & w_t^T x \cdot y < 0 \\ 0 & \text{otherwise} \end{cases}$$



Averaged Perceptron



The final hyperplane may be **overfitted** to later instances.

Take the **average** of all hyperplanes including ones that are not updated.

Averaged Perceptron

$$w_{t+1} \leftarrow w_t + \eta_t(x \cdot y) \quad \text{if } w_t^T \boxed{x} \cdot y < 0$$

$$\bar{w} \leftarrow \frac{1}{c} \sum_{t=0}^{c-1} w_t$$

sparse vector?

$$w_{t+1} \leftarrow w_t + \eta_t(x \cdot y)$$

$$v_{t+1} \leftarrow v_t + \eta_t \cdot \boxed{c}(x \cdot y)$$

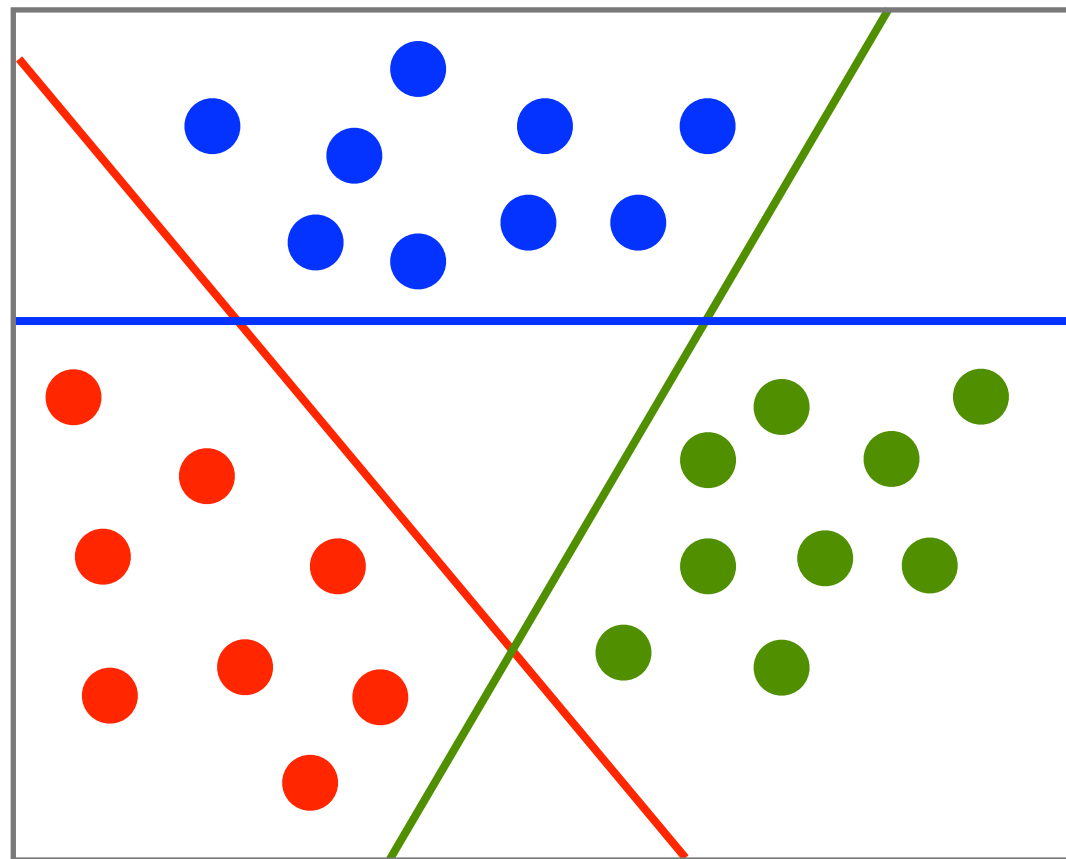
Initialization: $c \leftarrow 1$

Update rule: $c \leftarrow c + 1$ for every instance

$$\bar{w} \leftarrow w - \boxed{\frac{1}{c} \cdot v}$$

Multinomial Perceptron

Binomial distribution requires
1 hyperplane to separate 2 classes.



How many for
 m classes?

Multinomial distribution requires
 m hyperplanes to separate m classes.



Multinomial Perceptron

$$\mathbf{x} = \begin{bmatrix} 1 & 0 & 0 & 1 & 0 \end{bmatrix} \quad \text{5 features (including bias)}$$

Binomial

$$y = \{-1, 1\}$$

$$\mathbf{w} = \begin{bmatrix} a & b & c & d & e \end{bmatrix}$$

$$w^T \mathbf{x} = a + d \quad \hat{y} = \begin{cases} 1 & w^T \mathbf{x} \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

Multinomial

$$y = \{0, 1, 2, 3\}$$

$$\mathbf{w} = \begin{bmatrix} a_0 & a_1 & a_2 & a_3 & b_0 & b_1 & b_2 & b_3 & c_0 & c_1 & c_2 & c_3 & d_0 & d_1 & d_2 & d_3 & e_0 & e_1 & e_2 & e_3 \end{bmatrix}$$

$$w_y^T \mathbf{x} = a_y + d_y$$

$$\hat{y} = \arg \max_y w_y^T \mathbf{x}$$



Binomial vs. Multinomial Perceptron

$$\text{if } w_t^T x \cdot y < 0 \Leftrightarrow y \neq \hat{y}$$

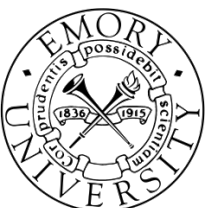
Binomial

$$w_{t+1} \leftarrow w_t + \eta_t (x \cdot y)$$

Multinomial

$$w_{y,t+1} \leftarrow w_{y,t} + \eta_t \cdot x$$

$$w_{\hat{y},t+1} \leftarrow w_{\hat{y},t} - \eta_t \cdot x$$



Hinge Loss

Perceptron

$$\ell(w, x; y) = \max\{0, -w^T x \cdot y\}$$

$$\Delta\ell = \begin{cases} -x \cdot y & w^T x \cdot y < 0 \\ 0 & \text{otherwise} \end{cases}$$

Hinge loss

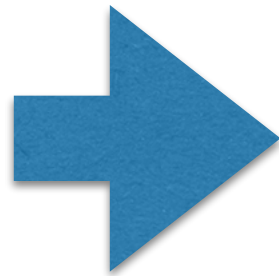
$$\ell(w, x; y) = \max\{0, 1 - w^T x \cdot y\}$$

$$\Delta\ell = \begin{cases} -x \cdot y & w^T x \cdot y < 1 \\ 0 & \text{otherwise} \end{cases}$$

Adaptive Gradient Descent

Perceptron

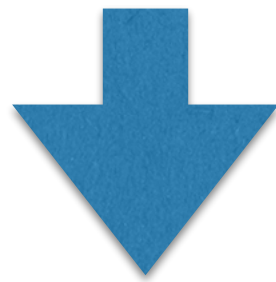
if $w_t^T x \cdot y < 0$



Hinge loss

if $w_t^T \cdot y < 1$

$$w_{t+1} \leftarrow w_t + \boxed{\eta_t} (x \cdot y)$$



$$g_{t+1} \leftarrow g_t + x \circ x$$

$$w_{t+1} \leftarrow w_t + \boxed{\frac{\eta}{\rho + \sqrt{g_{t+1}}}} \cdot (x \cdot y)$$