# Introduction to Machine Learning
# Lecture 6

Mehryar Mohri

Courant Institute and Google Research
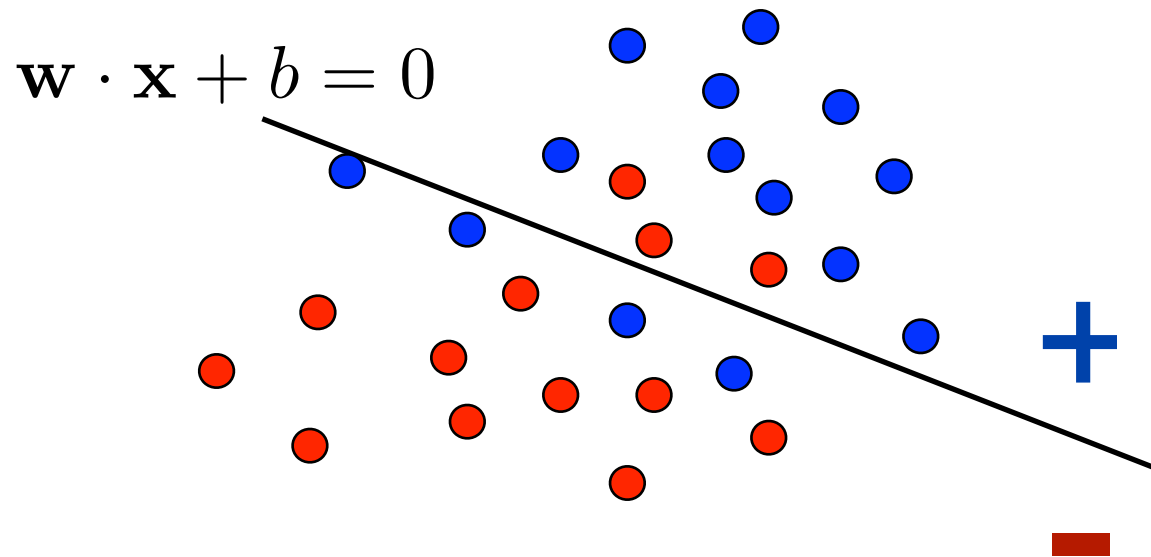mohri@cims.nyu.edu

# Perceptron and Winnow

# This Lecture

- On-Line linear classification: two algorithms.
  - Perceptron algorithm.
  - Winnow algorithm.

# Linear Classification

- Definition: a linear classifier is an algorithm that returns a hypothesis of the form

$$x \mapsto \operatorname{sgn}(\mathbf{w} \cdot \mathbf{x} + b),$$

with $\mathbf{w} \in \mathbb{R}^N, b \in \mathbb{R}.$

$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

# Margin Definitions

■ Definition: the (geometric) margin of a point $\mathbf{x}$ with label y for a linear classifier $h \colon \mathbf{x} \mapsto \mathbf{w} \cdot \mathbf{x} + b$ is its algebraic distance to the hyperplane $\mathbf{w} \cdot \mathbf{x} + b = 0$,

$$\rho(x) = \frac{y(\mathbf{w} \cdot \mathbf{x} + b)}{\|\mathbf{w}\|}.$$

■ Definition: the margin of a linear classifier $h$ for a sample $S = (x_1, \ldots, x_m)$ is the minimum margin of the points in that sample:

$$\rho = \min_{1 \leq i \leq m} \frac{y_i(\mathbf{w} \cdot \mathbf{x}_i + b)}{\|\mathbf{w}\|}.$$

# Perceptron Algorithm

(Rosenblatt, 1958)

$\text{Perceptron}(\mathbf{w}_0)$

1   $\mathbf{w}_1 \leftarrow \mathbf{w}_0$     $\triangleright$ typically $\mathbf{w}_0 = \mathbf{0}$

2   **for** $t \leftarrow 1$ **to** $T$ **do**

3      $\text{Receive}(\mathbf{x}_t)$

4      $\widehat{y}_t \leftarrow \text{sgn}(\mathbf{w}_t \cdot \mathbf{x}_t)$

5      $\text{Receive}(y_t)$

6      **if** $(\widehat{y}_t \neq y_t)$ **then**

7        $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_t\mathbf{x}_t$    $\triangleright$ more generally $\eta y_t\mathbf{x}_t, \eta > 0$

8      **else** $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t$

9   **return** $\mathbf{w}_{T+1}$

# Perceptron - Notes

- Update: if $y_t(\mathbf{w}_t \cdot \mathbf{x}_t) < 0$, then

$$y_t(\mathbf{w}_{t+1} \cdot \mathbf{x}_t) = y_t(\mathbf{w}_t \cdot \mathbf{x}_t) + \underbrace{\eta \|\mathbf{x}_t\|^2}_{\geq 0}.$$
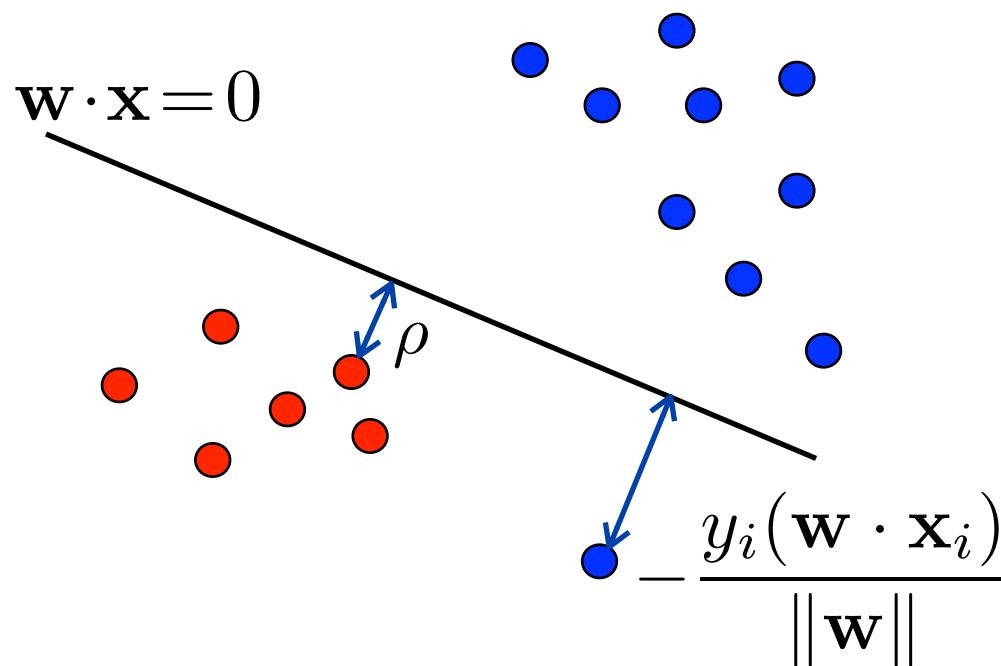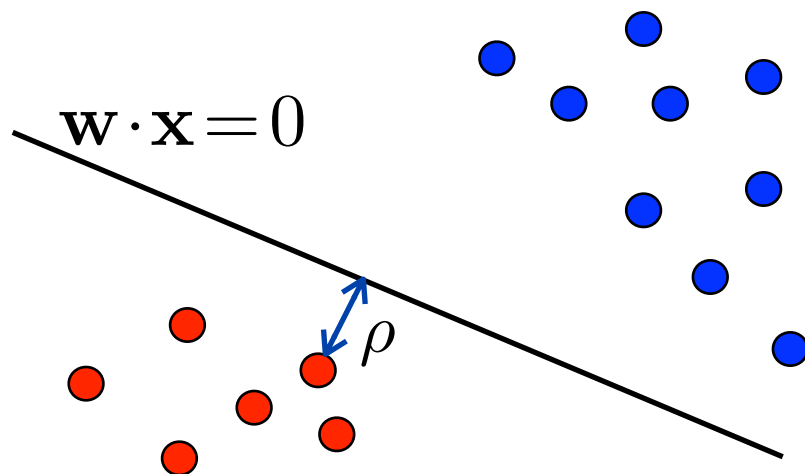
➡ change in the desired direction.

- Different modes of applications:

  - repeated passes over sample of size $m$ drawn according to some distribution $D$.

  - infinite sample drawn according to $D$.

  - no distributional assumption.

# Separating Hyperplane

■ Margin and errors

$$\mathbf{w} \cdot \mathbf{x} = 0 \qquad \rho \qquad\qquad \mathbf{w} \cdot \mathbf{x} = 0 \qquad \rho \qquad -\frac{y_i(\mathbf{w} \cdot \mathbf{x}_i)}{\|\mathbf{w}\|}$$

# Perceptron $=$ Stochastic Gradient Descent

- **Objective function**: convex but not differentiable.

$$F(\mathbf{w}) = \frac{1}{T} \sum_{t=1}^{T} \max\left(0, -y_t(\mathbf{w} \cdot \mathbf{x}_t)\right) = \mathop{\mathrm{E}}_{\mathbf{x} \sim \widehat{D}}[f(\mathbf{w}, \mathbf{x})]$$

with $f(\mathbf{w}, \mathbf{x}) = \max\left(0, -y(\mathbf{w} \cdot \mathbf{x})\right)$.

- **Stochastic gradient**: for each $\mathbf{x}_t$, the update is

$$\mathbf{w}_{t+1} \leftarrow \begin{cases} \mathbf{w}_t - \eta \nabla_{\mathbf{w}} f(\mathbf{w}_t, \mathbf{x}_t) & \text{if differentiable} \\ \mathbf{w}_t & \text{otherwise,} \end{cases}$$

where $\eta > 0$ is a learning rate parameter.

- **Here**:
$$\mathbf{w}_{t+1} \leftarrow \begin{cases} \mathbf{w}_t + \eta y_t \mathbf{x}_t & \text{if } y_t(\mathbf{w}_t \cdot \mathbf{x}_t) < 0 \\ \mathbf{w}_t & \text{otherwise.} \end{cases}$$

# Perceptron Algorithm - Bound

■ **Theorem:** Assume that $\|x_t\| \leq R$ for all $t \in [1, T]$ and that for some $\rho > 0$ and $\mathbf{v} \in \mathbb{R}^N$, for all $t \in [1, T]$,

$$\rho \leq \frac{y_t (\mathbf{v} \cdot \mathbf{x}_t)}{\|\mathbf{v}\|}.$$

Then, the number of mistakes made by the perceptron algorithm is bounded by $R^2 / \rho^2$.

■ **Proof:** Let $I$ be the set of $t$s at which there is an update and let $M$ be the total number of updates.

- Summing up the assumption inequalities gives:

$$M\rho \leq \frac{\mathbf{v} \cdot \sum_{t \in I} y_t \mathbf{x}_t}{\|\mathbf{v}\|}$$

$$= \frac{\mathbf{v} \cdot \sum_{t \in I} (\mathbf{w}_{t+1} - \mathbf{w}_t)}{\|\mathbf{v}\|} \quad \text{(definition of updates)}$$

$$= \frac{\mathbf{v} \cdot \mathbf{w}_{T+1}}{\|\mathbf{v}\|}$$

$$\leq \|\mathbf{w}_{T+1}\| \quad \text{(Cauchy-Schwarz ineq.)}$$

$$= \|\mathbf{w}_{t_m} + y_{t_m} \mathbf{x}_{t_m}\| \quad (t_m \text{ largest } t \text{ in } I)$$

$$= \left[ \|\mathbf{w}_{t_m}\|^2 + \|\mathbf{x}_{t_m}\|^2 + \underbrace{2 y_{t_m} \mathbf{w}_{t_m} \cdot \mathbf{x}_{t_m}}_{\leq 0} \right]^{1/2}$$

$$\leq \left[ \|\mathbf{w}_{t_m}\|^2 + R^2 \right]^{1/2}$$

$$\leq \left[ M R^2 \right]^{1/2} = \sqrt{M} R. \quad \text{(applying the same to previous } t\text{s in } I)$$

- Notes:
  - bound independent of dimension and tight.
  - convergence can be slow for small margin, it can be in $\Omega(2^N)$.
  - among the many variants: voted perceptron algorithm. Predict according to

    $$\mathrm{sgn}\left(\left(\sum_{t \in I} c_t \mathbf{w}_t\right) \cdot \mathbf{x}\right),$$

    where $c_t$ is the number of iterations $\mathbf{w}_t$ survives.
  - $\{x_t : t \in I\}$ are the support vectors for the perceptron algorithm.
  - non-separable case: does not converge.

# Leave-One-Out Error

- Definition: let $h_S$ be the hypothesis output by learning algorithm $L$ after receiving sample $S$ of size $m$. Then, the leave-one-out error of $L$ over $S$ is:

$$\widehat{R}_{\mathrm{loo}}(L) = \frac{1}{m} \sum_{i=1}^{m} 1_{h_{S-\{x_i\}}(x_i) \neq f(x_i)}.$$

- Property: unbiased estimate of expected error of hypothesis trained on sample of size $m-1$,

$$\boxed{\operatorname*{E}_{S \sim D^m}[\widehat{R}_{\mathrm{loo}}(L)]} = \frac{1}{m} \sum_{i=1}^{m} \operatorname*{E}_{S}[1_{h_{S-\{x_i\}}(x_i) \neq f(x_i)}] = \operatorname*{E}_{S}[1_{h_{S-\{x\}}(x) \neq f(x)}]$$

$$= \operatorname*{E}_{S' \sim D^{m-1}}[\operatorname*{E}_{x \sim D}[1_{h_{S'}(x) \neq f(x)}]] \boxed{= \operatorname*{E}_{S' \sim D^{m-1}}[R(h_{S'})].}$$

# Perceptron - Leave-One-Out Analysis

- Theorem: Assume that the data is separable. Let $h_S$ be the hypothesis returned by the Perceptron algorithm after training on sample $S \sim D^{m+1}$ (repeated passes) and let $M(S)$ be the number of updates made and let $R(h_S)$ be the error of $h_S$. Then,

$$\underset{S \sim D^m}{\mathrm{E}}[R(h_S)] \leq \underset{S \sim D^{m+1}}{\mathrm{E}}\left[\frac{\min(M(S), R_{m+1}^2/\rho_{m+1}^2)}{m+1}\right].$$

- Proof: Let $\mathbf{x}$ be a point in sample $S$. Then, If $h_{S-\{\mathbf{x}\}}$ misclassifies $\mathbf{x}$, there must have been an update at $\mathbf{x}$ during training to obtain $h_S$. Thus,

$$\widehat{R}_{\mathrm{loo}}(\mathrm{perceptron}) \leq \frac{M(S)}{m+1}.$$

# Dual Perceptron Algorithm

$\mathrm{DUAL\text{-}PERCEPTRON}(\boldsymbol{\alpha}_0)$

1    $\boldsymbol{\alpha}_1 \leftarrow \boldsymbol{\alpha}_0$      $\triangleright$ typically $\boldsymbol{\alpha}_0 = \mathbf{0}$

2    **for** $t \leftarrow 1$ **to** $T$ **do**

3          $\mathrm{RECEIVE}(\mathbf{x}_t)$

4          $\widehat{y}_t \leftarrow \mathrm{sgn}\big(\sum_{s=1}^{T} \alpha_s y_s (\mathbf{x}_s \cdot \mathbf{x}_t)\big)$

5          $\mathrm{RECEIVE}(y_t)$

6          **if** $(\widehat{y}_t \neq y_t)$ **then**

7              $\alpha_{t+1} \leftarrow \alpha_t + 1$

8          **else** $\alpha_{t+1} \leftarrow \alpha_t$

9    **return** $\boldsymbol{\alpha}$

# Kernel Perceptron Algorithm

$K$ PDS kernel.

$\text{KERNEL-PERCEPTRON}(\boldsymbol{\alpha}_0)$

1  $\boldsymbol{\alpha}_1 \leftarrow \boldsymbol{\alpha}_0$ $\qquad \triangleright$ typically $\boldsymbol{\alpha}_0 = \mathbf{0}$

2  **for** $t \leftarrow 1$ **to** $T$ **do**

3  $\qquad \text{RECEIVE}(x_t)$

4  $\qquad \widehat{y}_t \leftarrow \text{sgn}(\sum_{s=1}^{T} \alpha_s y_s K(x_s, x_t))$

5  $\qquad \text{RECEIVE}(y_t)$

6  $\qquad$ **if** $(\widehat{y}_t \neq y_t)$ **then**

7  $\qquad\qquad \alpha_{t+1} \leftarrow \alpha_t + 1$

8  $\qquad$ **else** $\alpha_{t+1} \leftarrow \alpha_t$

9  **return** $\boldsymbol{\alpha}$

# XOR Problem

- Use second-degree polynomial kernel with $c = 1$:



Linearly non-separable

Linearly separable by
$x_1 x_2 = 0$.

# Non-Separable Case

■ Theorem: Let $\mathbf{v}$ be any vector with $\|\mathbf{v}\| = 1$ and let $\rho > 0$. Define the deviation of $\mathbf{x}_t$ by:

$$d_t = \max\{0, \rho - y_t(\mathbf{v} \cdot \mathbf{x}_t)\},$$

and let $D = \sqrt{\sum_{t=1}^{T} d_t^2}$. Then, the number of perceptron updates after processing $\mathbf{x}_1, \ldots, \mathbf{x}_T$ is bounded by

$$\left[\frac{R + D}{\rho}\right]^2.$$

- **Proof**: Reduce problem to separable case in higher dimension.

  - Mapping (similar to trivial mapping):

$(N+t)$th component

$$\mathbf{x}_t = \begin{bmatrix} x_{t,1} \\ \vdots \\ x_{t,N} \end{bmatrix} \rightarrow \mathbf{x}'_t = \begin{bmatrix} x_{t,1} \\ \vdots \\ x_{t,N} \\ 0 \\ \vdots \\ 0 \\ \Delta \\ 0 \\ \vdots \\ 0 \end{bmatrix} \qquad \mathbf{v} \rightarrow \mathbf{v}' = \begin{bmatrix} v_1/Z \\ \vdots \\ v_N/Z \\ y_1 d_1/(\Delta Z) \\ \vdots \\ y_T d_T/(\Delta Z) \end{bmatrix}$$

$$\|\mathbf{v}'\| = 1 \Rightarrow Z = \sqrt{1 + \frac{D^2}{\Delta^2}}.$$

- **Now,** $y_t(\mathbf{v}' \cdot \mathbf{x}'_t) = y_t \left( \dfrac{\mathbf{v} \cdot \mathbf{x}_t}{Z} + \Delta \dfrac{y_t d_t}{Z\Delta} \right)$

$$= \frac{y_t \mathbf{v} \cdot \mathbf{x}_t}{Z} + \frac{d_t}{Z}$$

$$\geq \frac{y_t \mathbf{v} \cdot \mathbf{x}_t}{Z} + \frac{\rho - y_t(\mathbf{v} \cdot \mathbf{x}_t)}{Z} = \frac{\rho}{Z}.$$

- Since $\|\mathbf{x}'_t\|^2 \leq R^2 + \Delta^2$, the bound of the separable case applies: $\dfrac{(R^2 + \Delta^2)(1 + D^2/\Delta^2)}{\rho^2}$.

- With $\Delta = \sqrt{RD}$, this bound is minimized and equal to: $\dfrac{(R+D)^2}{\rho^2}$.

- Predictions made by the perceptron in the higher-dimension coincide with those of the perceptron in the original space.

# Winnow Algorithm

$\text{WINNOW}(\eta)$

1   $w_1 \leftarrow \mathbf{1}/N$

2   **for** $t \leftarrow 1$ **to** $T$ **do**

3        $\text{RECEIVE}(\mathbf{x}_t)$

4        $\widehat{y}_t \leftarrow \text{sgn}(\mathbf{w}_t \cdot \mathbf{x}_t)$         $\triangleright$    $y_t \in \{-1, +1\}$

5        $\text{RECEIVE}(y_t)$

6        **if** $(\widehat{y}_t \neq y_t)$ **then**

7            $Z_t \leftarrow \sum_{i=1}^{N} w_{t,i} \exp(\eta y_t x_{t,i})$

8            **for** $i \leftarrow 1$ **to** $N$ **do**

9                $w_{t+1,i} \leftarrow \dfrac{w_{t,i} \exp(\eta y_t x_{t,i})}{Z_t}$

10       **else** $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t$

11  **return** $\mathbf{w}_{T+1}$

# Winnow - Notes

- Winnow$=$weighted majority:
  - for $y_{t,i} = x_{t,i} \in \{-1, +1\}, \mathrm{sgn}(\mathbf{w}_t \cdot \mathbf{x}_t)$ coincides with the majority vote.
  - multiplying by $e^{\eta}$ or $e^{-\eta}$ the weight of correct or incorrect experts, is equivalent to multiplying by $\beta = e^{-2\eta}$ the weight of incorrect ones.

- Relationships with other algorithms: e.g., boosting and Perceptron (Winnow and Perceptron can be viewed as special instances of a general family).

- Motivation: large number of irrelevant features.

# Winnow Algorithm - Bound

■ **Theorem:** Assume that $\|x_t\|_\infty \leq R_\infty$ for all $t \in [1, T]$ and that for some $\rho_\infty > 0$ and $\mathbf{v} \in \mathbb{R}^N, \mathbf{v} \geq 0$ for all $t \in [1, T]$,

$$\rho_\infty \leq \frac{y_t(\mathbf{v} \cdot \mathbf{x}_t)}{\|\mathbf{v}\|_1}.$$

Then, the number of mistakes made by the Winnow algorithm is bounded by $\boxed{2\left(R_\infty^2 / \rho_\infty^2\right) \log N}$.

■ **Proof:** Let $I$ be the set of $t$s at which there is an update and let $M$ be the total number of updates.

# Winnow Algorithm - Bound

- Potential: $\Phi_t = \sum_{i=1}^{N} \frac{v_i}{\|\mathbf{v}\|} \log \frac{v_i/\|\mathbf{v}\|}{w_{t,i}}$.      (relative entropy)

- Upper bound: for each $t$ in $I$,

$$\Phi_{t+1} - \Phi_t = \sum_{i=1}^{N} \frac{v_i}{\|\mathbf{v}\|_1} \log \frac{w_{t,i}}{w_{t+1,i}}$$

$$= \sum_{i=1}^{N} \frac{v_i}{\|\mathbf{v}\|_1} \log \frac{Z_t}{\exp(\eta y_t x_{t,i})}$$

$$= \log Z_t - \eta \sum_{i=1}^{N} \frac{v_i}{\|\mathbf{v}\|_1} y_t x_{t,i}$$

$$\leq \log \Big[ \sum_{i=1}^{N} w_{t,i} \exp(\eta y_t x_{t,i}) \Big] - \eta \rho_{\infty}$$

$$= \log \mathop{\mathrm{E}}_{\mathbf{w}_t} \big[ \exp(\eta y_t x_t) \big] - \eta \rho_{\infty}$$

$$(\text{Hoeffding}) \leq \log \big[ \exp(\eta^2 (2R_{\infty})^2/8) \big] + \eta y_t \mathbf{w}_t \cdot \mathbf{x}_t - \eta \rho_{\infty}$$

$$\leq \eta^2 R_{\infty}^2/2 - \eta \rho_{\infty}.$$

# Winnow Algorithm - Bound

- **Upper bound:** summing up the inequalities yields

$$\Phi_{T+1} - \Phi_1 \leq M(\eta^2 R_\infty^2 / 2 - \eta \rho_\infty).$$

- **Lower bound:** note that

$$\Phi_1 = \sum_{i=1}^{N} \frac{v_i}{\|\mathbf{v}\|_1} \log \frac{v_i / \|\mathbf{v}\|_1}{1/N} = \log N + \sum_{i=1}^{N} \frac{v_i}{\|\mathbf{v}\|_1} \log \frac{v_i}{\|\mathbf{v}\|_1} \leq \log N$$

and for all $t, \Phi_t \geq 0$ (property or relative entropy).

Thus, $\Phi_{T+1} - \Phi_1 \geq 0 - \log N = -\log N.$

- **Comparison:** $-\log N \leq M(\eta^2 R_\infty^2 / 2 - \eta \rho_\infty).$ **For** $\eta = \frac{\rho_\infty}{R_\infty^2}$ we obtain

$$M \leq 2 \log N \frac{R_\infty^2}{\rho_\infty^2}.$$

# Notes

- Comparison with perceptron bound:
  - dual norms: norms for $x_t$ and $v$.
  - similar bounds with different norms.
  - each advantageous in different cases:
    - Winnow bound favorable when a sparse set of experts can predict well. For example, if $v = e_1$ and $x_t \in \{\pm 1\}^N$, $\log N$ **vs** $N$.
    - Perceptron favorable in opposite situation.