# Machine learning: linear classification

# Linear classification framework

**training data**

|  | $x_1$ | $x_2$ | $y$ |
|---|---|---|---|
| **example** | 0 | 2 | 1 |
| **example** | -2 | 0 | 1 |
| **example** | 1 | -1 | -1 |

**learning algorithm** →

$[2, 0]$ **input**

↓

$f$ **classifier**

↓

-1 **label**



**decision boundary**

Design decisions:

Which classifiers are possible? **hypothesis class**

How good is a classifier? **loss function**
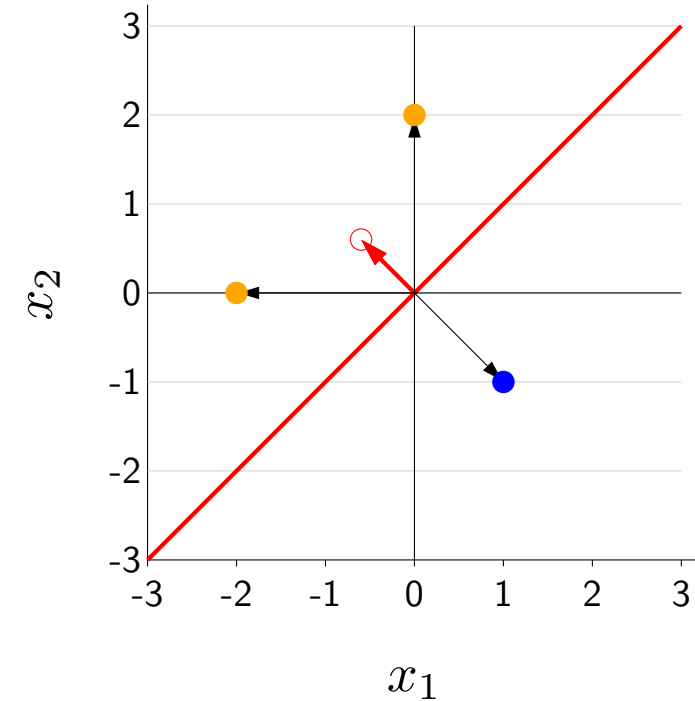
How do we compute the best classifier? **optimization algorithm**

# An example linear classifier

$$\overbrace{}^{\mathbf{w}} \quad \overbrace{}^{\phi(x)}$$

$$f(x) = \text{sign}([-0.6, 0.6] \cdot [x_1, x_2])$$

$$\text{sign}(z) = \begin{cases} +1 & \text{if } z > 0 \\ -1 & \text{if } z < 0 \\ 0 & \text{if } z = 0 \end{cases}$$

| $x_1$ | $x_2$ | $f(x)$ |
|-------|-------|--------|
| 0 | 2 | 1 |
| -2 | 0 | 1 |
| 1 | -1 | -1 |



$$f([0, 2]) = \text{sign}([-0.6, 0.6] \cdot [0, 2]) = \text{sign}(1.2) = 1$$

$$f([-2, 0]) = \text{sign}([-0.6, 0.6] \cdot [-2, 0]) = \text{sign}(1.2) = 1$$

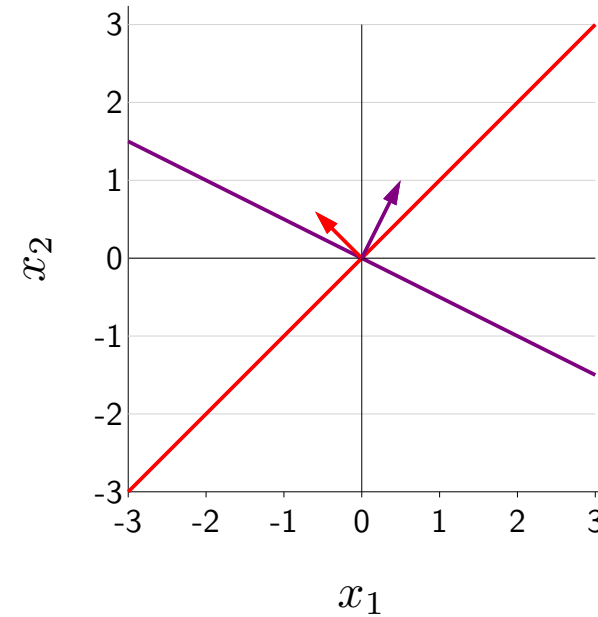$$f([1, -1]) = \text{sign}([-0.6, 0.6] \cdot [1, -1]) = \text{sign}(-1.2) = -1$$

Decision boundary: $x$ such that $\mathbf{w} \cdot \phi(x) = 0$

# Hypothesis class: which classifiers?



$$\phi(x) = [x_1, x_2]$$

$$f(x) = \text{sign}([-0.6, 0.6] \cdot \phi(x))$$

$$f(x) = \text{sign}([0.5, 1] \cdot \phi(x))$$

General binary classifier:

$$f_{\mathbf{w}}(x) = \text{sign}(\mathbf{w} \cdot \phi(x))$$

Hypothesis class:

$$\mathcal{F} = \{f_{\mathbf{w}} : \mathbf{w} \in \mathbb{R}^2\}$$

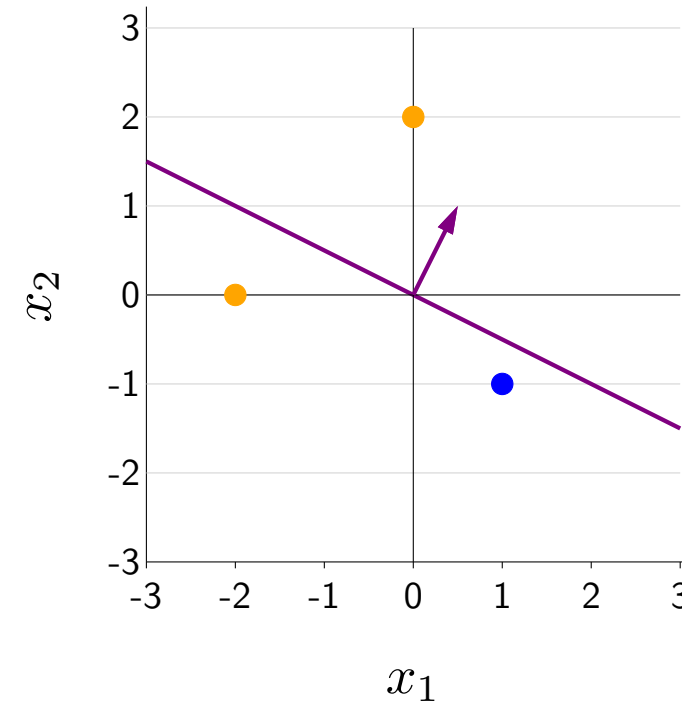# Loss function: how good is a classifier?

$$f_{\mathbf{w}}(x) = \mathbf{w} \cdot \phi(x)$$
$$\mathbf{w} = [0.5, 1]$$
$$\phi(x) = [x_1, x_2]$$

**training data** $\mathcal{D}_{\text{train}}$

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0 | 2 | 1 |
| -2 | 0 | 1 |
| 1 | -1 | -1 |



$$\text{Loss}_{\text{0-1}}(x, y, \mathbf{w}) = \mathbf{1}[f_{\mathbf{w}}(x) \neq y] \text{ **zero-one loss**}$$

$$\text{Loss}([0, 2], 1, [0.5, 1]) = \mathbf{1}[\text{sign}([0.5, 1] \cdot [0, 2]) \neq 1] = 0$$

$$\text{Loss}([-2, 0], 1, [0.5, 1]) = \mathbf{1}[\text{sign}([0.5, 1] \cdot [-2, 0]) \neq 1] = 1$$
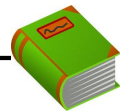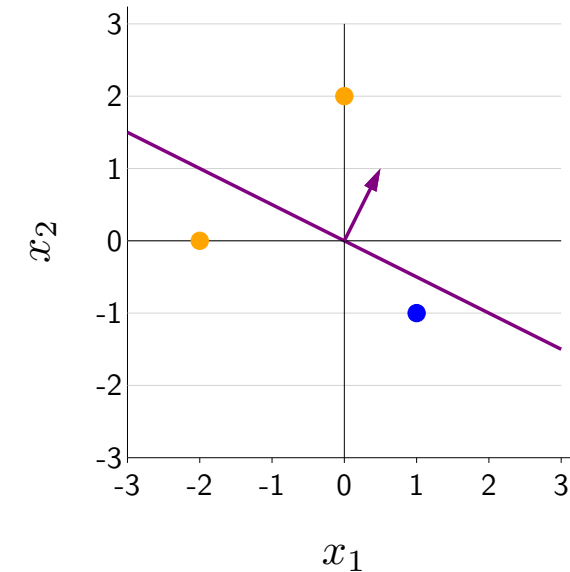
$$\text{Loss}([1, -1], -1, [0.5, 1]) = \mathbf{1}[\text{sign}([0.5, 1] \cdot [1, -1]) \neq -1] = 0$$

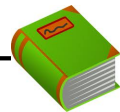$$\text{TrainLoss}([0.5, 1]) = 0.33$$

# Score and margin

Predicted label: $f_{\mathbf{w}}(x) = \mathsf{sign}(\mathbf{w} \cdot \phi(x))$

Target label: $y$
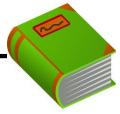


---

📗 **Definition: score**

The score on an example $(x, y)$ is $\mathbf{w} \cdot \phi(x)$, how **confident** we are in predicting $+1$.
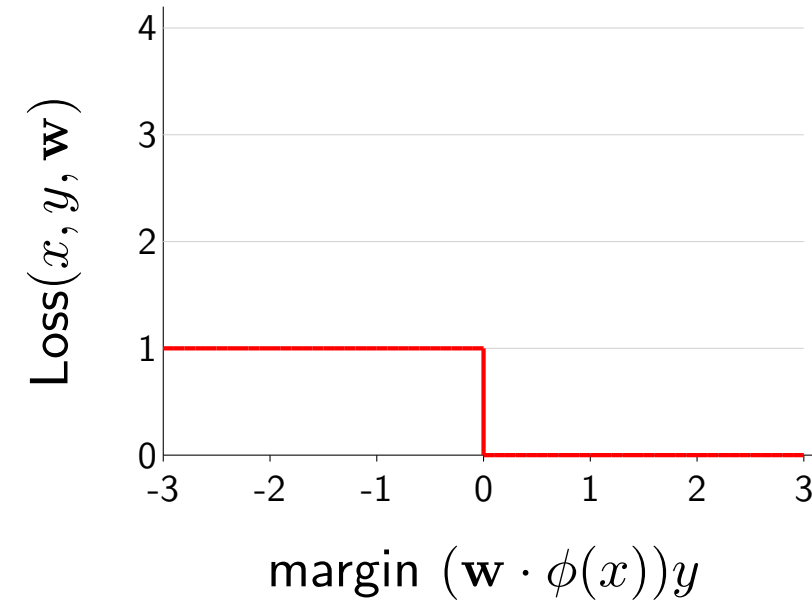
---

📗 **Definition: margin**

The margin on an example $(x, y)$ is $(\mathbf{w} \cdot \phi(x))y$, how **correct** we are.

# Zero-one loss rewritten

**Definition: zero-one loss**

$$\text{Loss}_{0\text{-}1}(x, y, \mathbf{w}) = \mathbf{1}[f_{\mathbf{w}}(x) \neq y]$$

$$= \mathbf{1}[\underbrace{(\mathbf{w} \cdot \phi(x))y}_{\text{margin}} \leq 0]$$
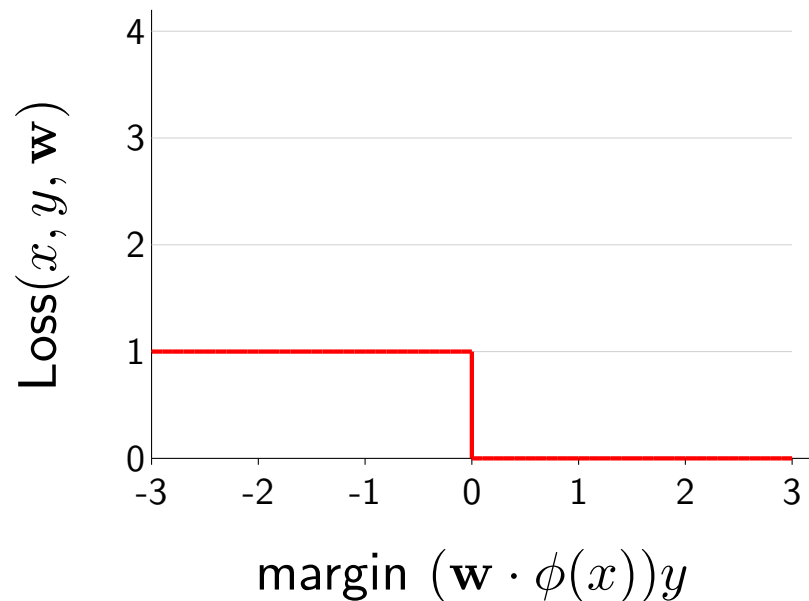
# Optimization algorithm: how to compute best?

Goal: $\min_{\mathbf{w}} \text{TrainLoss}(\mathbf{w})$

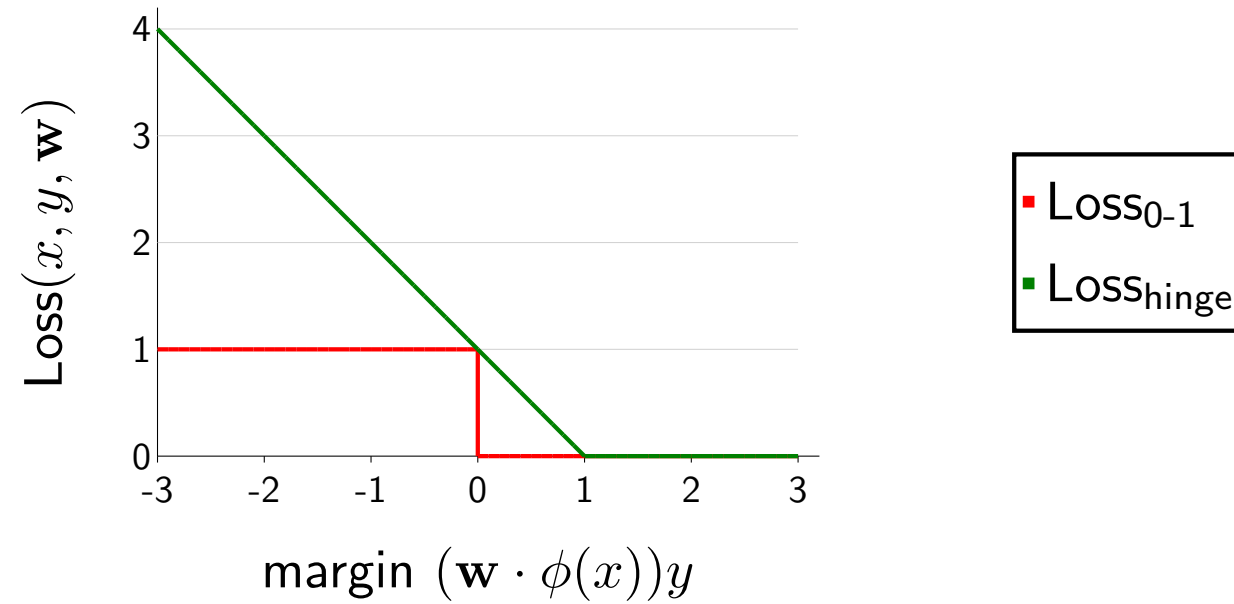To run gradient descent, compute the gradient:

$$\nabla_{\mathbf{w}} \text{TrainLoss}(\mathbf{w}) = \frac{1}{|\mathcal{D}_{\text{train}}|} \sum_{(x,y) \in \mathcal{D}_{\text{train}}} \nabla \text{Loss}_{\text{0-1}}(x, y, \mathbf{w})$$

$$\nabla_{\mathbf{w}} \text{Loss}_{\text{0-1}}(x, y, \mathbf{w}) = \nabla \mathbf{1}[(\mathbf{w} \cdot \phi(x))y \leq 0]$$
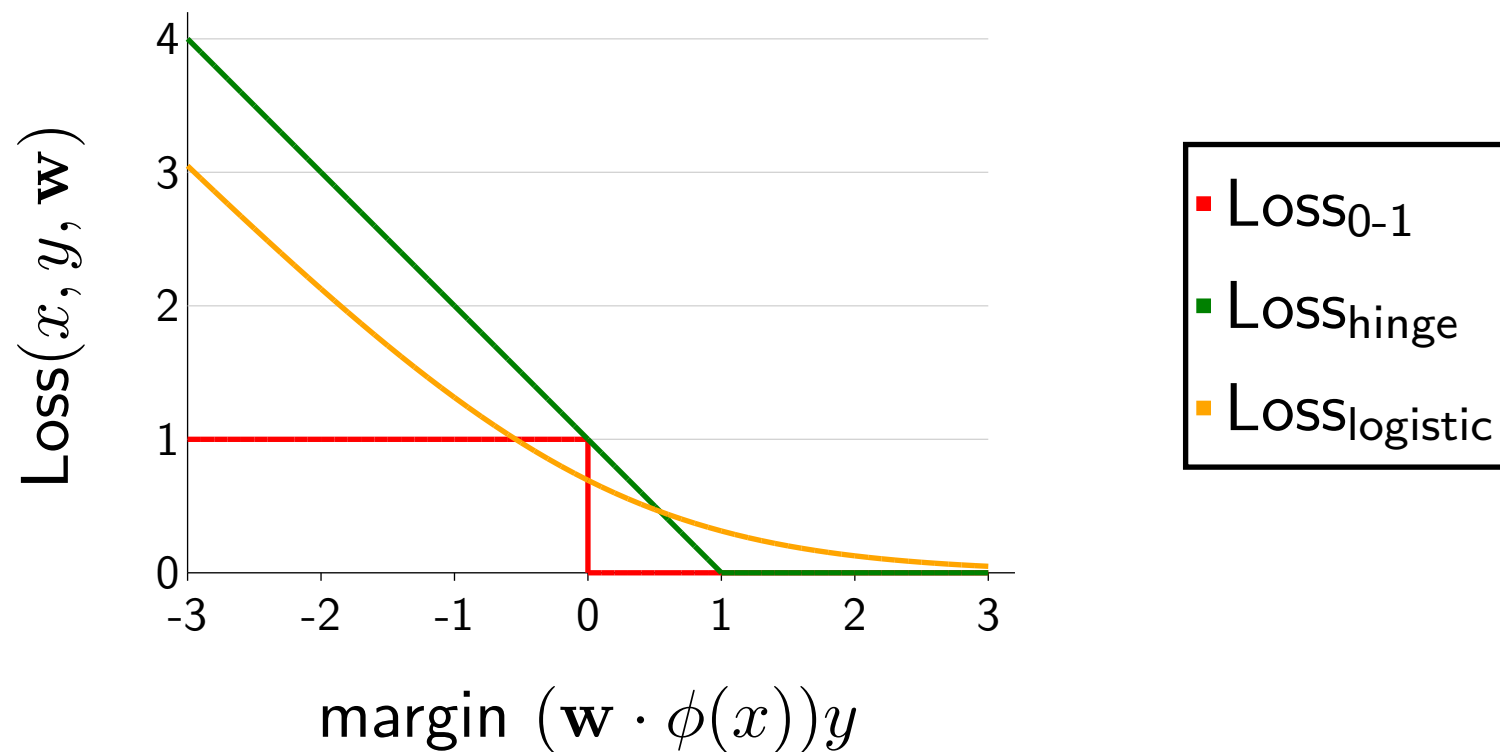


**Gradient is zero almost everywhere!**

# Hinge loss



$$\text{Loss}_{\text{hinge}}(x, y, \mathbf{w}) = \max\{1 - (\mathbf{w} \cdot \phi(x))y, 0\}$$
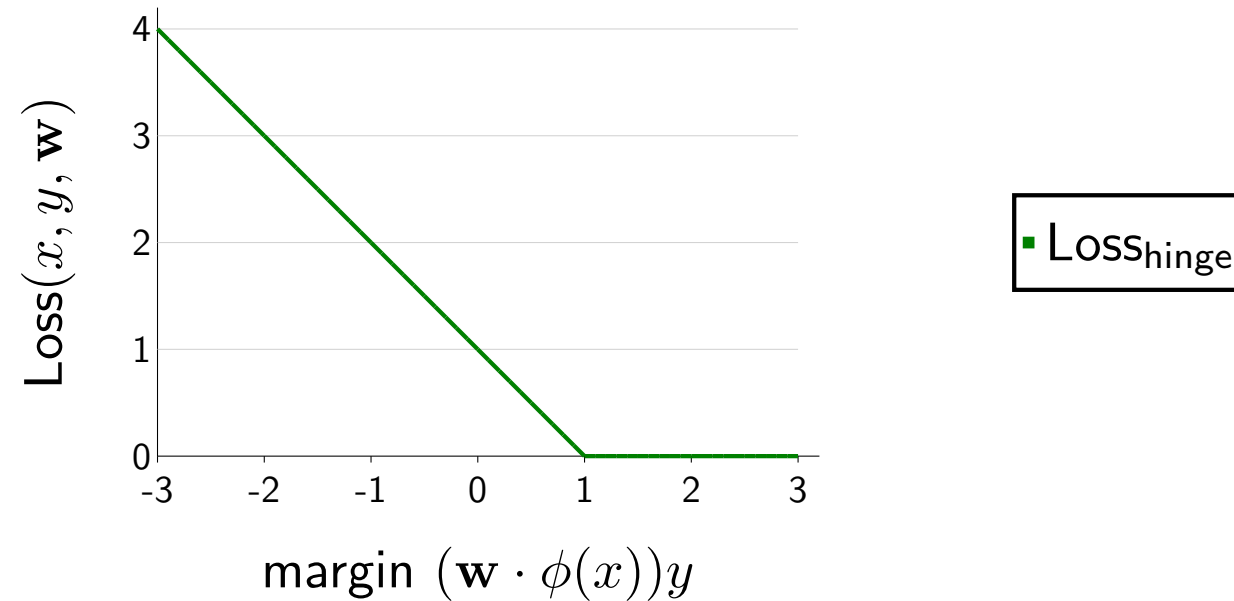
# Digression: logistic regression

$$\mathsf{Loss}_{\mathsf{logistic}}(x, y, \mathbf{w}) = \log(1 + e^{-(\mathbf{w} \cdot \phi(x))y})$$



Intuition: Try to increase margin even when it already exceeds 1

# Gradient of the hinge loss



$$\text{Loss}_{\text{hinge}}(x, y, \mathbf{w}) = \max\{1 - (\mathbf{w} \cdot \phi(x))y, 0\}$$

$$\nabla \text{Loss}_{\text{hinge}}(x, y, \mathbf{w}) = \begin{cases} -\phi(x)y & \text{if } \{1 - (\mathbf{w} \cdot \phi(x))y\} > \{0\} \\ 0 & \text{otherwise} \end{cases}$$
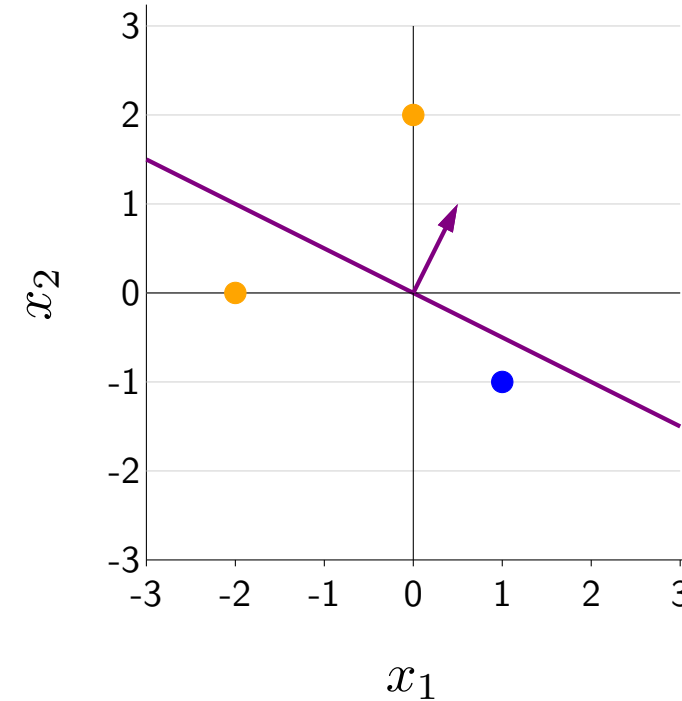
# Hinge loss on training data

$$f_{\mathbf{w}}(x) = \mathbf{w} \cdot \phi(x)$$

$$\mathbf{w} = [0.5, 1]$$

$$\phi(x) = [x_1, x_2]$$

**training data** $\mathcal{D}_{\text{train}}$

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0 | 2 | 1 |
| -2 | 0 | 1 |
| 1 | -1 | -1 |



$$\text{Loss}_{\text{hinge}}(x, y, \mathbf{w}) = \max\{1 - (\mathbf{w} \cdot \phi(x))y, 0\}$$

$\text{Loss}([0, 2], 1, [0.5, 1]) = \max\{1 - [0.5, 1] \cdot [0, 2](1), 0\} = 0$ $\qquad$ $\nabla\text{Loss}([0, 2], 1, [0.5, 1]) = [0, 0]$

$\text{Loss}([-2, 0], 1, [0.5, 1]) = \max\{1 - [0.5, 1] \cdot [-2, 0](1), 0\} = 2$ $\qquad$ $\nabla\text{Loss}([-2, 0], 1, [0.5, 1]) = [2, 0]$

$\text{Loss}([1, -1], -1, [0.5, 1]) = \max\{1 - [0.5, 1] \cdot [1, -1](-1), 0\} = 0.5$ $\qquad$ $\nabla\text{Loss}([1, -1], -1, [0.5, 1]) = [1, -1]$

$\text{TrainLoss}([0.5, 1]) = 0.83$ $\qquad$ $\nabla\text{TrainLoss}([0.5, 1]) = [1, -0.33]$

# Gradient descent (hinge loss) in Python

[code]

# Summary so far

$$\underbrace{\mathbf{w} \cdot \phi(x)}_{\text{score}}$$

|  | Regression | Classification |
|---|---|---|
| Prediction $f_{\mathbf{w}}(x)$ | score | sign(score) |
| Relate to target $y$ | residual (score $- y$) | margin (score $y$) |
| Loss functions | squared<br>absolute deviation | zero-one<br>hinge<br>logistic |
| Algorithm | gradient descent | gradient descent |