Faculty of Science

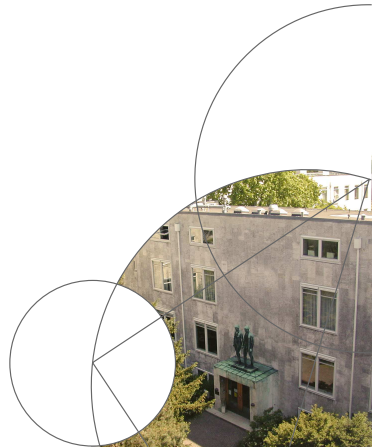# Linear Classification
## Machine Learning

Christian Igel
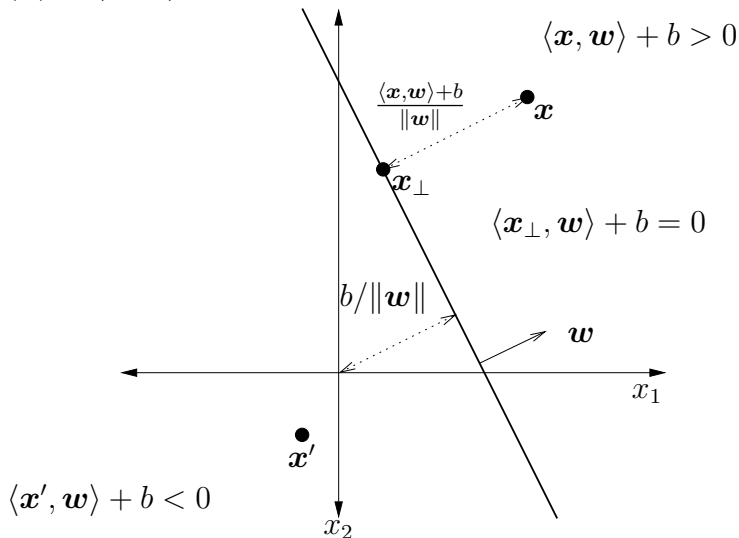Department of Computer Science

🐦 @christian_igel

# Outline

**①** Logistic Regression

**②** Linear Classification and Margins

**③** Perceptron Learning

**④** Convergence of Perceptron Learning

**⑤** Summary

# Outline

## Linear functions

$$f(\boldsymbol{x}) = \langle \boldsymbol{x}, \boldsymbol{w} \rangle + b$$



$\langle \boldsymbol{x}, \boldsymbol{w} \rangle + b > 0$

$\frac{\langle \boldsymbol{x}, \boldsymbol{w} \rangle + b}{\|\boldsymbol{w}\|}$

$\boldsymbol{x}$

$\boldsymbol{x}_\perp$

$\langle \boldsymbol{x}_\perp, \boldsymbol{w} \rangle + b = 0$

$b/\|\boldsymbol{w}\|$

$\boldsymbol{w}$

$x_1$

$\boldsymbol{x}'$

$\langle \boldsymbol{x}', \boldsymbol{w} \rangle + b < 0$

$x_2$

## Decision functions

- Classification assigns an input $x \in \mathcal{X}$ to one of a finite set of classes $\mathcal{Y} = \{\mathcal{C}_1, \ldots, \mathcal{C}_m\}$, $2 \leq m$.

- One approach is to learn *discriminant functions* $\delta_k : \mathcal{X} \to \mathbb{R}$, $1 \leq k \leq m$, and assign a pattern $x$ to class $\hat{y}$ using

$$\hat{y} = h(x) = \mathrm{argmax}_k \, \delta_k(x) \ .$$

## Linear classification

- We build affine linear decision functions

$$\delta(\boldsymbol{x}) = \sum_{i=1}^{d} w_i x_i + b = \boldsymbol{w}^\mathsf{T} \boldsymbol{x} + b$$

  with $\boldsymbol{w} \in \mathbb{R}^d$ and $b \in \mathbb{R}$.

- For convenience, we define $\tilde{\boldsymbol{x}}_i^\mathsf{T} = (x_1, \ldots, x_d, 1)$ for $i = 1, \ldots, N$ and $\tilde{\boldsymbol{w}}^\mathsf{T} = (w_1, \ldots, w_d, b)$ and consider the equivalent formulation

$$\delta(\tilde{\boldsymbol{x}}) = \sum_{i=1}^{d+1} \tilde{w}_i \tilde{x}_i = \tilde{\boldsymbol{w}}^\mathsf{T} \tilde{\boldsymbol{x}} \ .$$

  We omit the tilde in the following.

## Binary decision functions

- If we have only two classes, we can consider a single function

$$\delta(x) = \delta_1(x) - \delta_2(x)$$

  and the hypothesis

$$h(x) = \begin{cases} \mathcal{C}_1 & \text{if } \delta(x) > 0 \\ \mathcal{C}_2 & \text{otherwise} \end{cases}.$$

- For $\mathcal{Y} = \{-1, 1\}$ this is equal to

$$h(x) = \operatorname{sgn}(\delta(x)) = \begin{cases} 1 & \text{if } \delta(x) > 0 \\ -1 & \text{otherwise} \end{cases}.$$

# Decision functions and class posteriors

- If we know the class posteriors $P(Y \mid X)$ we can perform optimal classification: a pattern $x$ is assigned to class $\mathcal{C}_k$ with maximum $P(Y = \mathcal{C}_k \mid X = x)$, i.e.,

$$\hat{y} = h(x) = \operatorname{argmax}_k P(Y = \mathcal{C}_k \mid X = x)$$

  or in the binary case with $\mathcal{Y} = \{-1, 1\}$

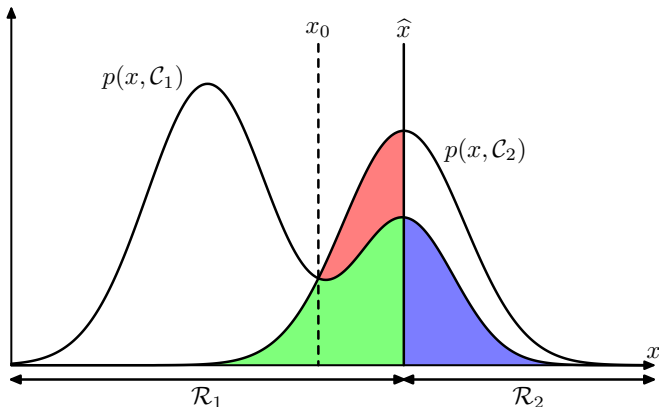$$\delta(x) = P(Y = \mathcal{C}_1 \mid X = x) - P(Y = \mathcal{C}_2 \mid X = x)$$

  and $\hat{y} = h(x) = \operatorname{sgn}(\delta(x))$.

- $P(Y = \mathcal{C}_k \mid X = x)$ is proportional to the class-conditional density $p(X = x \mid Y = \mathcal{C}_k)$ times the class prior $P(Y = \mathcal{C}_k)$:

$$P(Y = \mathcal{C}_k \mid X = x) = \frac{p(X = x \mid Y = \mathcal{C}_k) P(Y = \mathcal{C}_k)}{p(X = x)}$$
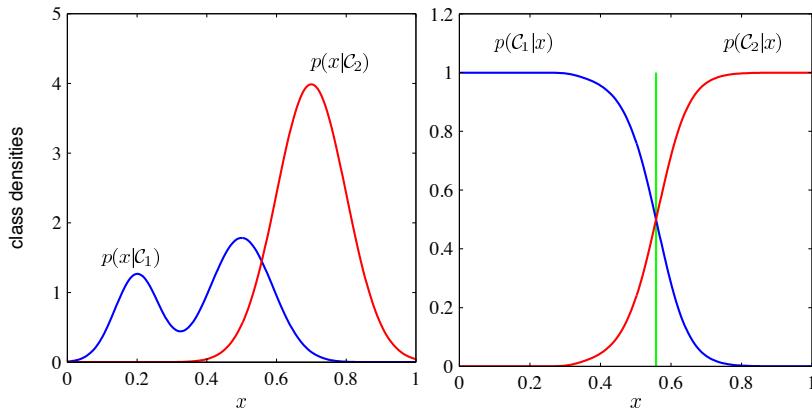
# Joint probabilities



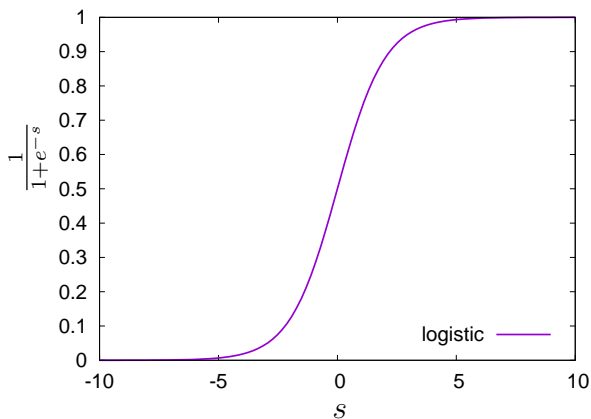C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer-Verlag, 2006

# Class-conditional densities



C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer-Verlag, 2006

# Logistic function



$$\theta(s) = \frac{1}{1 + e^{-s}}$$

## Predicting probabilities

- Instead of predicting the class label, we want to learn

$$f(\boldsymbol{x}) = P(Y = 1 \,|\, X = \boldsymbol{x})$$

assuming that the data is generated by

$$P(Y = y \,|\, X = \boldsymbol{x}) = \begin{cases} f(\boldsymbol{x}) & \text{for } y = 1 \\ 1 - f(\boldsymbol{x}) & \text{for } y = -1 \end{cases} .$$

- In the binary case, our model takes the form $h : \mathcal{X} \to [0, 1]$:

$$h(\boldsymbol{x}) = \theta(\boldsymbol{w}^{\mathsf{T}} \boldsymbol{x})$$

## Log-odds

Choice of logistic function is not arbitrary, it guarantees that
the (affine) linear part of the model encodes the log-odds:

$$\boldsymbol{w}^\mathsf{T}\boldsymbol{x} = \ln \frac{P(Y = 1 \mid X = \boldsymbol{x})}{P(Y = -1 \mid X = \boldsymbol{x})}$$

That is, $\boldsymbol{w}^\mathsf{T}\boldsymbol{x}$ encodes on log-scale how frequent class 1
occurs relative to class 0-1.

## Likelihood function

- Our hypothesis $h$ describes the probability distribution:

$$P(Y = y \mid X = \boldsymbol{x}; \boldsymbol{w}) = \theta(y\boldsymbol{w}^\mathsf{T}\boldsymbol{x}) = \begin{cases} \theta(\boldsymbol{w}^\mathsf{T}\boldsymbol{x}) & \text{for } y = 1 \\ 1 - \theta(\boldsymbol{w}^\mathsf{T}\boldsymbol{x}) & \text{for } y = -1 \end{cases}$$

- $S = \{(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_N, y_N)\} \subseteq (\mathbb{R}^n \times \{-1, 1\})^N$

- Likelihood (function) of the parameters $\boldsymbol{w}$ given training data $S$ is the probability of observing $S$ when the data is generated by $h$ with parameters $\boldsymbol{w}$.

- Likelihood for i.i.d. $S$:

$$\prod_{i=1}^N P(Y = y_i \mid X = \boldsymbol{x}_i; h) \text{ or short } \prod_{i=1}^N P(y_i \mid \boldsymbol{x}_i)$$

# Maximum likelihood

- Learning principle: Maximize the likelihood function!
- Equivalently, we can minimize the negative logarithmic likelihood.
- Negative log-likelihood (divided by $N$):

$$-\frac{1}{N} \ln \left( \prod_{i=1}^{N} P(y_i \,|\, \boldsymbol{x}_i) \right) = -\frac{1}{N} \sum_{i=1}^{N} \ln \left( P(y_i \,|\, \boldsymbol{x}_i) \right)$$

- Plugging in our linear hypothesis gives the error function:

$$-\frac{1}{N} \sum_{n=1}^{N} \ln \left( \theta(y_n \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n) \right) = \frac{1}{N} \sum_{n=1}^{N} \ln \left( 1 + e^{-y_n \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n} \right)$$

## Recall: Gradient

- The *gradient*

$$\nabla f(\boldsymbol{x}) = \left( \frac{\partial f(\boldsymbol{x})}{\partial x_1}, \frac{\partial f(\boldsymbol{x})}{\partial x_2}, \ldots, \frac{\partial f(\boldsymbol{x})}{\partial x_d} \right)^{\mathsf{T}}$$

points in the direction $\nabla f(\boldsymbol{x})/\|\nabla f(\boldsymbol{x})\|$ giving maximum rate of change $\|\nabla f(\boldsymbol{x})\|$.

## Gradient descent

- Consider learning by iteratively changing the parameters:

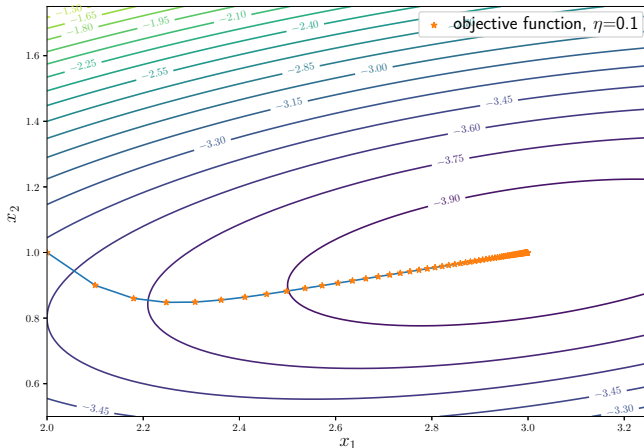$$\boldsymbol{w} \leftarrow \boldsymbol{w} + \Delta\boldsymbol{w}$$

- Simplest choice is (steepest) gradient descent

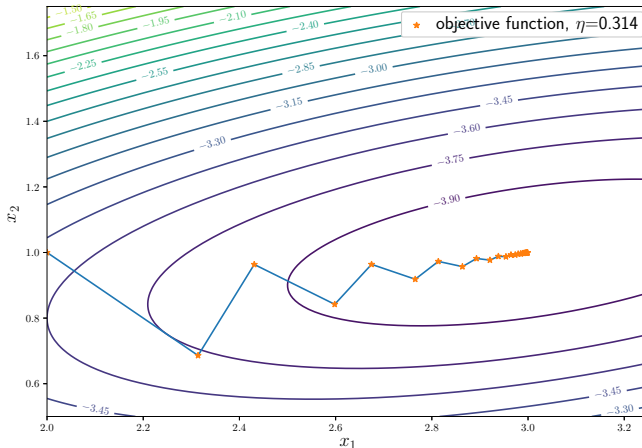$$\Delta\boldsymbol{w} = -\eta \nabla f|_{\boldsymbol{w}}$$

with learning rate $\eta > 0$.

# Gradient descent example, small learning rate

# Gradient descent, larger learning rate

# Gradient for training logistic regression

- For data $\{(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_N, y_N)\} \subseteq (\mathbb{R}^n \times \{-1, 1\})^N$, we have the following gradient of the negative log-likelihood:

$$-\frac{1}{N} \sum_{n=1}^{N} \frac{y_n \boldsymbol{x}_n}{1 + e^{y_n \boldsymbol{w}^\mathsf{T} \boldsymbol{x}_n}}$$

- This equals:

$$-\frac{1}{N} \sum_{n=1}^{N} \left[ \frac{y_n + 1}{2} - \theta(\boldsymbol{w}^\mathsf{T} \boldsymbol{x}_n) \right] \boldsymbol{x}_n$$

- Thus, for $\{(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_N, y_N)\} \subseteq (\mathbb{R}^n \times \{0, 1\})^N$, we have the gradient:

$$-\frac{1}{N} \sum_{n=1}^{N} \left[ y_n - \theta(\boldsymbol{w}^\mathsf{T} \boldsymbol{x}_n) \right] \boldsymbol{x}_n$$

# Logistic regression algorithm
# (steepest descent)

**Algorithm 1:** Logistic regression

**Input:** data $\{(\boldsymbol{x}_1, y_1), \ldots, \boldsymbol{x}_N, y_N)\} \subseteq (\mathbb{R}^n \times \{-1, 1\})^N$,
learning rate $\eta$

**Output:** weights of linear hypothesis $h(\boldsymbol{x}) = \langle \boldsymbol{w}, \boldsymbol{x} \rangle$

1   initialize $\boldsymbol{w}$

2   **repeat**

>    // gradient of negative log-likelihood over $N$

3    $\boldsymbol{g} \leftarrow -\frac{1}{N} \sum_{n=1}^{N} \frac{y_n \boldsymbol{x}_n}{1 + e^{y_n \boldsymbol{w}^\mathsf{T} \boldsymbol{x}_n}}$

>    // model parameter update

4    $\boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \boldsymbol{g}$

5   **until** *stopping criterion is met*

# Logistic regression algorithm
# (stochastic gradient descent, SGD)

**Algorithm 2:** Logistic regression

**Input:** data $\{(\boldsymbol{x}_1, y_1), \dots\} \subseteq (\mathbb{R}^n \times \{-1, 1\})^N$, learning
rate $\eta$

**Output:** weights of linear hypothesis $h(\boldsymbol{x}) = \langle \boldsymbol{w}, \boldsymbol{x} \rangle$

1   initialize $\boldsymbol{w}$
2   **repeat**
3      pick $(\boldsymbol{x}, y) \in S$
4      $\boldsymbol{w} \leftarrow \boldsymbol{w} + \eta \frac{y\boldsymbol{x}}{1 + e^{y\boldsymbol{w}^\mathsf{T}\boldsymbol{x}}}$
5   **until** *stopping criterion is met*

# Logistic regression algorithm
# (mini-batch gradient descent)

**Algorithm 3:** Logistic regression

**Input:** data $\{(\boldsymbol{x}_1, y_1), \dots\} \subseteq (\mathbb{R}^n \times \{-1, 1\})^N$, learning
rate $\eta$

**Output:** weights of linear hypothesis $h(\boldsymbol{x}) = \langle \boldsymbol{w}, \boldsymbol{x} \rangle$

1  initialize $\boldsymbol{w}$

2  **repeat**

3      pick $S' \subset S$

4      $\boldsymbol{g} \leftarrow -\frac{1}{|S'|} \sum_{(\boldsymbol{x}, y) \in S'} \frac{y\boldsymbol{x}}{1 + e^{y\boldsymbol{w}^\mathsf{T}\boldsymbol{x}}}$

5      $\boldsymbol{w} \leftarrow \boldsymbol{w} - \eta\boldsymbol{g}$

6  **until** *stopping criterion is met*

## Multiple classes

- Binary, single discriminant function, $y \in \{0, 1\}$:

$$P(Y = 1 \,|\, \boldsymbol{x}) = \frac{1}{1 + e^{-\delta(\boldsymbol{x})}} = \frac{e^{\delta(\boldsymbol{x})}}{1 + e^{\delta(\boldsymbol{x})}}$$

- Binary, two discriminant functions, $y \in \{1, 2\}$:

$$P(Y = y \,|\, \boldsymbol{x}) = \frac{e^{\delta_y(\boldsymbol{x})}}{e^{\delta_1(\boldsymbol{x})} + e^{\delta_2(\boldsymbol{x})}} = \frac{e^{\delta_y(\boldsymbol{x})+C}}{e^{\delta_1(\boldsymbol{x})+C} + e^{\delta_2(\boldsymbol{x})+C}}$$

for every constant $C$, thus logistic function is special case for $C = -\delta_1(\boldsymbol{x})$.

- Multiple classes, $y \in \{1, \dots, m\}$:

$$P(Y = y \,|\, \boldsymbol{x}) = \underbrace{\frac{e^{\delta_y(\boldsymbol{x})}}{\sum_{i=1}^{m} e^{\delta_i(\boldsymbol{x})}}}_{\text{softmax function}}$$

# Outline

# Margins I

The functional margin of an example $(\boldsymbol{x}_i, y_i)$ with respect to a hyperplane $(\boldsymbol{w}, b)$ is

$$\gamma_i := y_i(\langle \boldsymbol{w}, \boldsymbol{x}_i \rangle + b) \ .$$

The geometric margin of an example $(\boldsymbol{x}_i, y_i)$ with respect to a hyperplane $(\boldsymbol{w}, b)$ is

$$\rho_i := y_i(\langle \boldsymbol{w}, \boldsymbol{x}_i \rangle + b)/\|\boldsymbol{w}\| = \gamma_i/\|\boldsymbol{w}\| \ .$$

A positive margin implies correct classification.
The margin of a hyperplane $(\boldsymbol{w}, b)$ with respect to a training set $S$ is $\min_i \rho_i$. The margin of a training set $S$ is the maximum geometric margin over all hyperplanes. A hyperplane realizing this margin is called maximum margin hyperplane.

# Margins II

# Outline

**1** Logistic Regression

**2** Linear Classification and Margins

**3** Perceptron Learning

**4** Convergence of Perceptron Learning

**5** Summary

# Analyzing the Perceptron

Why should we look at the Perceptron?

- Linear classifiers such as perceptrons are the basis of technical neurocomputing

- Support Vector Machines are basically linear classifiers

- Basic concepts of learning theory can be explained easily:
  - Margins
  - Dual representation
  - Bounds involving margins and the radius of the ball containing the data

# Perceptron learning algorithm (primal form)

For simplicity, consider hyperplanes with no bias ($b = 0$), i.e.,
$\mathcal{H} = \{h(\boldsymbol{x}) = \mathrm{sgn}(\langle \boldsymbol{w}, \boldsymbol{x} \rangle) \,|\, \boldsymbol{w} \in \mathbb{R}^n\}$.

---

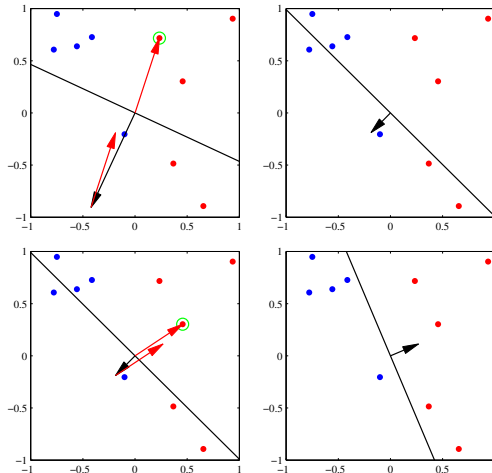**Algorithm 4:** Perceptron

---

**Input:** separable data $\{(\boldsymbol{x}_1, y_1), \dots\} \subseteq (\mathbb{R}^n \times \{-1, 1\})^N$
**Output:** hypothesis $h(\boldsymbol{x}) = \mathrm{sgn}(\langle \boldsymbol{w}_k, \boldsymbol{x} \rangle)$

1   $\boldsymbol{w}_0 \leftarrow \boldsymbol{0}; k \leftarrow 0$
2   **repeat**
3     **for** $i = 1, \dots, N$ **do**
4       **if** $y_i \langle \boldsymbol{w}_k, \boldsymbol{x}_i \rangle \leq 0$ **then**
5         $\boldsymbol{w}_{k+1} \leftarrow \boldsymbol{w}_k + y_i \boldsymbol{x}_i$
6         $k \leftarrow k + 1$

7   **until** *no mistake made within* **for** *loop*

---

# Perceptron learning in pictures



C. M. Bishop. *Pattern Recognition and Machine Learning.* Springer-Verlag, 2006

# Outline

# Novikoff

### Theorem (Novikoff)

Let $S = \{(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_N, y_N)\}$ be a non-trivial training set
(i.e., containing patterns of both classes), $\boldsymbol{w}_0 = \boldsymbol{0}$, and let

$$R \leftarrow \max_{1 \leq i \leq N} \|\boldsymbol{x}_i\| \ .$$

Suppose that there exists $\boldsymbol{w}_{opt}$ and $\rho > 0$ such that
$\|\boldsymbol{w}_{opt}\| = 1$ and

$$y_i \langle \boldsymbol{w}_{opt}, \boldsymbol{x}_i \rangle \geq \rho > 0$$

for $1 \leq i \leq N$. Then the number of updates $k$ made by the
online perceptron algorithm on $S$ is at most

$$\left( \frac{R}{\rho} \right)^2 \ .$$

# Novikoff, sketch of proof I

Let $i$ be the index of the example in update $k$

$$\begin{aligned}
\|\boldsymbol{w}_{k+1}\|^2 &= \langle \boldsymbol{w}_k + y_i \boldsymbol{x}_i, \boldsymbol{w}_k + y_i \boldsymbol{x}_i \rangle \\
&= \|\boldsymbol{w}_k\|^2 + 2y_i \langle \boldsymbol{w}_k, \boldsymbol{x}_i \rangle + \|\boldsymbol{x}_i\|^2 \\
&\leq \|\boldsymbol{w}_k\|^2 + R^2 \\
&\leq (k+1)R^2
\end{aligned}$$

# Novikoff, sketch of proof II

$$\langle \boldsymbol{w}_{\mathsf{opt}}, \boldsymbol{w}_{k+1} \rangle = \langle \boldsymbol{w}_{\mathsf{opt}}, \boldsymbol{w}_k \rangle + y_i \langle \boldsymbol{w}_{\mathsf{opt}}, \boldsymbol{x}_i \rangle$$
$$\geq \langle \boldsymbol{w}_{\mathsf{opt}}, \boldsymbol{w}_k \rangle + \rho$$
$$\geq (k+1)\rho$$

$$k^2 \rho^2 \leq \langle \boldsymbol{w}_{\mathsf{opt}}, \boldsymbol{w}_k \rangle^2 \leq \|\boldsymbol{w}_{\mathsf{opt}}\|^2 \|\boldsymbol{w}_k\|^2 \leq kR^2$$

$$k \leq \frac{R^2}{\rho^2}$$

## Dual representation

- Weight vector of hyperplane computed by online perceptron algorithm can be written as

$$\boldsymbol{w} = \sum_{i=1}^{N} \alpha_i y_i \boldsymbol{x}_i$$

- Function $h(\boldsymbol{x}) = \mathrm{sgn}(\delta(\boldsymbol{x}))$ can be written in dual coordinates

$$\begin{aligned} \delta(\boldsymbol{x}) &= \langle \boldsymbol{w}, \boldsymbol{x} \rangle \\ &= \left\langle \sum_{i=1}^{N} \alpha_i y_i \boldsymbol{x}_i, \boldsymbol{x} \right\rangle \\ &= \sum_{i=1}^{N} \alpha_i y_i \langle \boldsymbol{x}_i, \boldsymbol{x} \rangle \end{aligned}$$

# Perceptron learning algorithm (dual form)

**Algorithm 5:** Perceptron (dual form)

**Input:** separable data $\{(\boldsymbol{x}_1, y_1), \dots\} \subseteq (\mathbb{R}^n \times \{-1, 1\})^N$

**Output:** hypothesis $h(\boldsymbol{x}) = \mathrm{sgn}\left(\sum_{i=1}^{N} \alpha_i y_i \langle \boldsymbol{x}_i, \boldsymbol{x} \rangle\right)$

1 $\boldsymbol{\alpha} \leftarrow \boldsymbol{0}$

2 **repeat**

3     **for** $i = 1, \dots, N$ **do**

4        **if** $y_i \sum_{j=1}^{N} \alpha_j y_j \langle \boldsymbol{x}_j, \boldsymbol{x}_i \rangle \leq 0$ **then**

5           $\alpha_i \leftarrow \alpha_i + 1$

6 **until** *no mistake made within* **for** *loop*

# Outline

① Logistic Regression

② Linear Classification and Margins

③ Perceptron Learning

④ Convergence of Perceptron Learning

⑤ Summary

# Summary I

Logistic regression

- is easy to use, has in its simplest form no hyperparameters (not counting $\eta$),
- gives surprisingly good results, is highly recommended as baseline method,
- does typically not tend to overfit (assuming $d \ll N$), but does not capture non-linearities,
- can be used with non-linear transformations,
- can be parallelized and is applicable to "Big Data".

# Summary II

Hey, we also now know about

- perceptron learning,

- margins,

- dual representation,

- bounds involving margins and the radius of the ball containing the data.