

Statistical Machine Learning (BE4M33SSU)

Lecture 11: Structured Output Support Vector Machines

Czech Technical University in Prague
V. Franc

Linear classifier

Two-class linear classifier:

- ◆ \mathcal{X} is a set of observations and $\mathcal{Y} = \{+1, -1\}$ is a set of hidden labels
- ◆ $\phi: \mathcal{X} \rightarrow \mathbb{R}^n$ - feature map embedding observations from \mathcal{X} to \mathbb{R}^n
- ◆ Two-class linear classifier $h: \mathcal{X} \rightarrow \mathcal{Y}$

$$h(x; \mathbf{w}, b) = \text{sign}(\langle \mathbf{w}, \phi(x) \rangle + b) = \begin{cases} +1 & \text{if } \langle \mathbf{w}, \phi(x) \rangle + b \geq 0 \\ -1 & \text{if } \langle \mathbf{w}, \phi(x) \rangle + b < 0 \end{cases}$$

A generic linear classifier:

- ◆ \mathcal{X} is set of observations and \mathcal{Y} is a finite set
- ◆ $\phi: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^n$ is a joint feature map embedding $\mathcal{X} \times \mathcal{Y}$ to \mathbb{R}^n
- ◆ Generic linear classifier $h: \mathcal{X} \rightarrow \mathcal{Y}$

$$h(x; \mathbf{w}) \in \underset{y \in \mathcal{Y}}{\text{Argmax}} \langle \mathbf{w}, \phi(x, y) \rangle$$

Example: multi-class linear classifier

- ◆ \mathcal{X} is a set of observations and $\mathcal{Y} = \{1, \dots, Y\}$ is a set of class labels
- ◆ Multi-class linear classifier $h: \mathcal{X} \rightarrow \mathcal{Y}$

$$h(x; \mathbf{w}) \in \underset{y \in \mathcal{Y}}{\text{Argmax}} \langle \mathbf{w}_y, \phi(x) \rangle$$

where $\phi: \mathcal{X} \rightarrow \mathbb{R}^d$ is a feature map $\mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_Y) \in \mathbb{R}^{d \cdot Y}$ are parameters.

- ◆ We can write the score function as

$$\langle \mathbf{w}_y, \phi(x) \rangle = \langle \mathbf{w}, \phi(x, y) \rangle$$

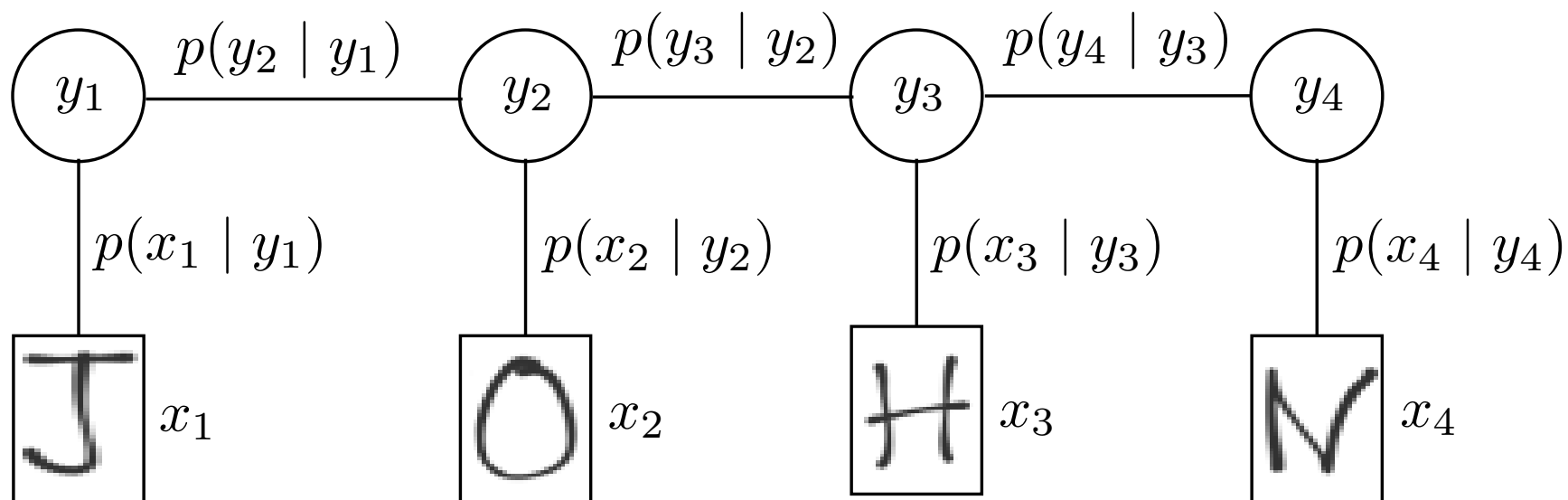
where $\phi: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^{d \cdot Y}$ is

$$\phi(x, y) = (0, \dots, \underbrace{\phi(x)}_{y\text{-th slot}}, \dots, 0)$$

Example: sequence classifier for OCR

- ◆ $\mathbf{x} = (x_1, \dots, x_L) \in \mathcal{I}^L$ - sequence of images with characters
- ◆ $\mathbf{y} = (y_1, \dots, y_L) \in \mathcal{A}^L$ - seq. of chars. from $\mathcal{A} = \{A, \dots, Z\}$
- ◆ $p(x_i | y_i)$ - appearance model for characters
- ◆ $p(y_i | y_{i-1})$ - language model
- ◆ Finding the most probable sequence of characters:

$$\hat{\mathbf{y}} \in \underset{\mathbf{y} \in \mathcal{A}^L}{\text{Argmax}} \left(\underbrace{p(y_1) \prod_{i=2}^L p(y_i | y_{i-1}) \prod_{i=1}^L p(x_i | y_i)}_{p(x_1, \dots, x_L, y_1, \dots, y_L)} \right)$$



Example: sequence classifier for OCR

- ◆ The MAP estimate from HMM:

$$\hat{\mathbf{y}} \in \underset{\mathbf{y} \in \mathcal{A}^L}{\text{Argmax}} \left(\log p(y_1) + \sum_{i=2}^L \log p(y_i | y_{i-1}) + \sum_{i=1}^L \log p(x_i | y_i) \right)$$

- ◆ Let us assume the following parametrization:

$$\begin{aligned} \log p(y_1) &= \langle \mathbf{w}, \phi(y_1) \rangle \\ \log p(y_i | y_{i-1}) &= \langle \mathbf{w}, \phi(y_{i-1}, y_i) \rangle \\ \log p(x_i | y_i) &= \langle \mathbf{w}, \phi(x_i, y_i) \rangle \end{aligned}$$

- ◆ The MAP estimate becomes a linear classifier:

$$\hat{\mathbf{y}} = \underset{(y_1, \dots, y_k) \in \mathcal{A}^L}{\text{Argmax}} \left\langle \mathbf{w}, \underbrace{\phi(y_1) + \sum_{i=2}^L \phi(y_{i-1}, y_i) + \sum_{i=1}^L \phi(x_i, y_i)}_{\phi(\mathbf{x}, \mathbf{y})} \right\rangle$$

Learning by Empirical Risk Minimization

- ◆ $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow [0, \infty)$ - loss such that $\ell(y, y') = 0$ iff $y = y'$.
- ◆ Find parameters \mathbf{w} of $h(x; \mathbf{w})$ which minimize the expected risk

$$R(\mathbf{w}) = \mathbb{E}_{(x,y) \sim p} \left(\ell(y, h(x; \mathbf{w})) \right)$$

- ◆ The Empirical Risk Minimization principle leads to solving

$$\mathbf{w}^* \in \underset{\mathbf{w} \in \mathbb{R}^n}{\text{Argmin}} R_{\mathcal{T}^m}(\mathbf{w})$$

where the empirical risk is

$$R_{\mathcal{T}^m}(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \ell(y^i, h(x^i; \mathbf{w}))$$

and $\mathcal{T}^m = \{(x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}) \mid i = 1, \dots, m\}$ are training examples drawn from i.i.d. with distribution $p(x, y)$.

Learning linear classifier from separable examples

- ◆ A correctly classified example (x^i, y^i) , that is,

$$y^i = h(x^i; \mathbf{w}) = \underset{y \in \mathcal{Y}}{\text{Argmax}} \langle \mathbf{w}, \phi(x^i, y) \rangle$$

implies

$$\langle \phi(x^i, y^i), \mathbf{w} \rangle > \langle \phi(x^i, y), \mathbf{w} \rangle, \quad \forall y \in \mathcal{Y} \setminus \{y^i\}$$

Definition 1. *The examples $\mathcal{T}^m = \{(x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}) \mid i = 1, \dots, m\}$ are linearly separable w.r.t. joint feature map $\phi: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^n$ if there exists $\mathbf{w} \in \mathbb{R}^n$ such that*

$$\langle \phi(x^i, y^i), \mathbf{w} \rangle > \langle \phi(x^i, y), \mathbf{w} \rangle, \quad \forall i \in \{1, \dots, m\}, y \in \mathcal{Y} \setminus \{y^i\}$$

(Generic) Perceptron algorithm

- ◆ **Task:** given a set of points $\{\mathbf{a}^i \in \mathbb{R}^n \mid i = 1, 2, \dots, K\}$ we want to find $\mathbf{w} \in \mathbb{R}^n$ such that

$$\langle \mathbf{w}, \mathbf{a}^i \rangle > 0, \quad \forall i \in \{1, 2, \dots, K\} \quad (1)$$

- ◆ **Perceptron:**

1. $\mathbf{w} \leftarrow \mathbf{0}$
2. Find a violating $\langle \mathbf{w}, \mathbf{a}^i \rangle \leq 0, i \in \{1, 2, \dots, K\}$
3. If there is no violating inequality return \mathbf{w} otherwise update

$$\mathbf{w} \leftarrow \mathbf{w} + \mathbf{a}^i$$

and go to step 2.

- ◆ If the set of inequalities (1) is solvable then the Perceptron algorithm exits in a finite number of steps which does not depend on m .

Structured Output Perceptron

- Learning $h(x; \mathbf{w}) \in \text{Argmax}_{y \in \mathcal{Y}} \langle \mathbf{w}, \phi(x, y) \rangle$ from examples $\mathcal{T}^m = \{(x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}) \mid i = 1, \dots, m\}$ leads to solving

$$\langle \phi(x^i, y^i), \mathbf{w} \rangle - \langle \phi(x^i, y), \mathbf{w} \rangle > 0, \quad \forall i \in \{1, \dots, m\}, y \in \mathcal{Y} \setminus \{y^i\}$$

- Structured Output Perceptron:**

1. $\mathbf{w} \leftarrow \mathbf{0}$

2. Find a misclassified example $(x^i, y^i) \in \mathcal{T}^m$ such that

$$y^i \neq \hat{y}^i \in \text{Argmax}_{y \in \mathcal{Y}} \langle \mathbf{w}, \phi(x^i, y) \rangle \quad \text{prediction problem}$$

3. If there is no misclassified example return \mathbf{w} otherwise update

$$\mathbf{w} \leftarrow \mathbf{w} + \phi(x^i, y^i) - \phi(x^i, \hat{y}^i) \quad \text{parameter update}$$

and go to step 2.

Structured Output SVM

- ◆ Learning $h(x; \mathbf{w}) \in \text{Argmax}_{y \in \mathcal{Y}} \langle \mathbf{w}, \phi(x, y) \rangle$ from examples $\mathcal{T}^m = \{(x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}) \mid i = 1, \dots, m\}$ by ERM leads to

$$\mathbf{w}^* \in \underset{\mathbf{w} \in \mathbb{R}^n}{\text{Argmin}} R_{\mathcal{T}^m}(\mathbf{w}) \quad \text{where} \quad R_{\mathcal{T}^m}(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \ell(y^i, h(x^i; \mathbf{w}))$$

- ◆ The SO-SVM approximates the ERM by a convex problem

$$\mathbf{w}^* \in \underset{\mathbf{w} \in \mathcal{W}_r}{\text{Argmin}} R^\psi(\mathbf{w}) \quad \text{where} \quad R^\psi(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \psi(x^i, y^i, \mathbf{w})$$

where

- $\mathcal{W}_r \subseteq \mathbb{R}^n$ - convex feasible set; e.g. $\mathcal{W}_r = \{\mathbf{w} \in \mathbb{R}^n \mid \|\mathbf{w}\| \leq r\}$
- $\psi: \mathcal{X} \times \mathcal{Y} \times \mathbb{R}^n \rightarrow \mathbb{R}$ - convex proxy approximating the true loss ℓ

Margin rescaling loss

- ◆ We require that the score of the correct label y^i is higher than the score of the incorrect label y by margin proportional to the loss $\ell(y^i, y)$:

$$\langle \mathbf{w}, \phi(x^i, y^i) \rangle \geq \langle \mathbf{w}, \phi(x^i, y) \rangle + \ell(y^i, y), \quad \forall y \in \mathcal{Y} \setminus \{y^i\}$$

- ◆ The margin rescaling loss

$$\psi(x^i, y^i, \mathbf{w}) = \max \left\{ 0, \max_{y \in \mathcal{Y} \setminus \{y^i\}} \{ \ell(y^i, y) + \langle \mathbf{w}, \phi(x^i, y) \rangle - \langle \mathbf{w}, \phi(x^i, y^i) \rangle \} \right\}$$

- ◆ Upper bounds of the true loss:

$$y^i \neq \hat{y} = h(x^i; \mathbf{w}) \in \underset{y \in \mathcal{Y}}{\text{Argmax}} \langle \mathbf{w}, \phi(x^i, y) \rangle$$

implies $\langle \mathbf{w}, \phi(x^i, \hat{y}) \rangle - \langle \mathbf{w}, \phi(x^i, y^i) \rangle \geq 0$ and hence

$$\psi(x^i, y^i, \mathbf{w}) \geq \ell(y^i, h(x^i, \mathbf{w})), \quad \forall \mathbf{w} \in \mathbb{R}^n$$

SO-SVM with margin-rescaling loss

- Using shortcuts $\ell_i(y) = \ell(y^i, y)$ and $\phi_i(y) = \phi(x^i, y) - \phi(x^i, y^i)$ we can simplify the margin rescaling loss:

$$\begin{aligned}
 \psi(x^i, y^i, \mathbf{w}) &= \max\{0, \max_{y \in \mathcal{Y} \setminus \{y^i\}} \{\ell(y^i, y) + \langle \mathbf{w}, \phi(x^i, y) \rangle - \langle \mathbf{w}, \phi(x^i, y^i) \rangle\}\} \\
 &= \max_{y \in \mathcal{Y}} \{\ell(y^i, y) + \langle \mathbf{w}, \phi(x^i, y) \rangle - \langle \mathbf{w}, \phi(x^i, y^i) \rangle\} \\
 &= \max_{y \in \mathcal{Y}} \{\ell_i(y) + \langle \mathbf{w}, \phi_i(y) \rangle\}
 \end{aligned}$$

- The SO-SVM leads to a convex constrained optimization problem:

$$\mathbf{w}^* \in \underset{\mathbf{w} \in \mathcal{W}_r}{\text{Argmin}} R^\psi(\mathbf{w}) \quad \text{where} \quad R^\psi(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \max_{y \in \mathcal{Y}} \{\ell_i(y) + \langle \mathbf{w}, \phi_i(y) \rangle\}$$

and $\mathcal{W}_r = \{\mathbf{w} \in \mathbb{R}^n \mid \|\mathbf{w}\| \leq r\}$ is a convex set.

SO-SVM as a convex quadratic program

- ◆ The SO-SVM problem can be written as unconstrained problem:

$$\mathbf{w}^* \in \underset{\mathbf{w} \in \mathbb{R}^n}{\text{Argmin}} \left(\frac{\lambda}{2} \|\mathbf{w}\|^2 + R^\psi(\mathbf{w}) \right)$$

- ◆ After introducing slack variables it can be further rewritten as a constrained quadratic program:

$$\mathbf{w}^* = \underset{\mathbf{w} \in \mathbb{R}^n, \boldsymbol{\xi} \in \mathbb{R}^m}{\text{argmin}} \left(\frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \right)$$

subject to

$$\xi_i \geq \ell_i(y) + \langle \mathbf{w}, \phi_i(y) \rangle, \quad \forall i \in \{1, \dots, m\}, \forall y \in \mathcal{Y}$$

- ◆ Note that the QP has $m|\mathcal{Y}|$ linear constraints !

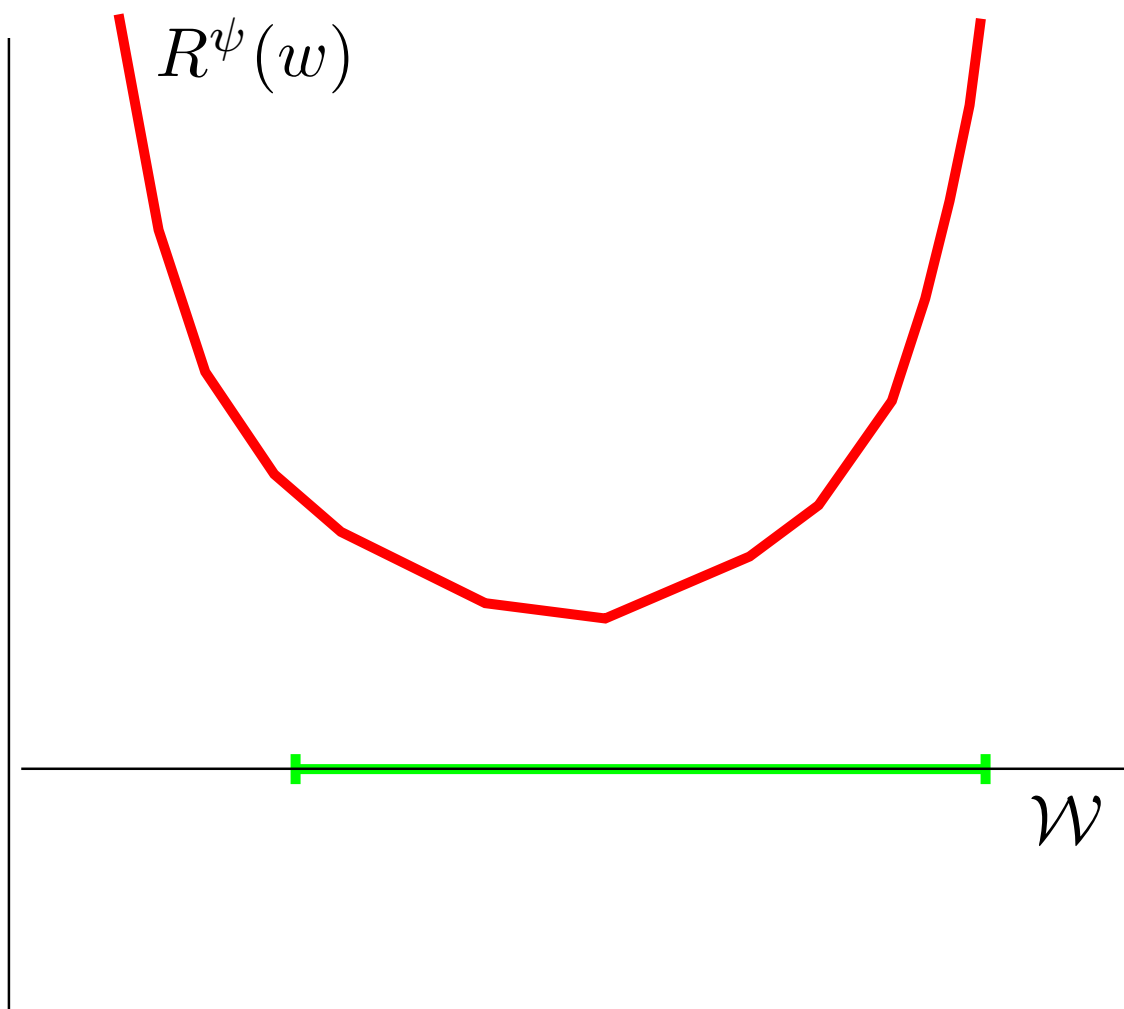
Cutting plane algorithm



$$R^\psi(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \max_{\hat{y}^i \in \mathcal{Y}} (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), \mathbf{w} \rangle) = \max_{\substack{\hat{y}^1 \in \mathcal{Y} \\ \vdots \\ \hat{y}^m \in \mathcal{Y}}} \frac{1}{m} \sum_{i=1}^m (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), \mathbf{w} \rangle)$$

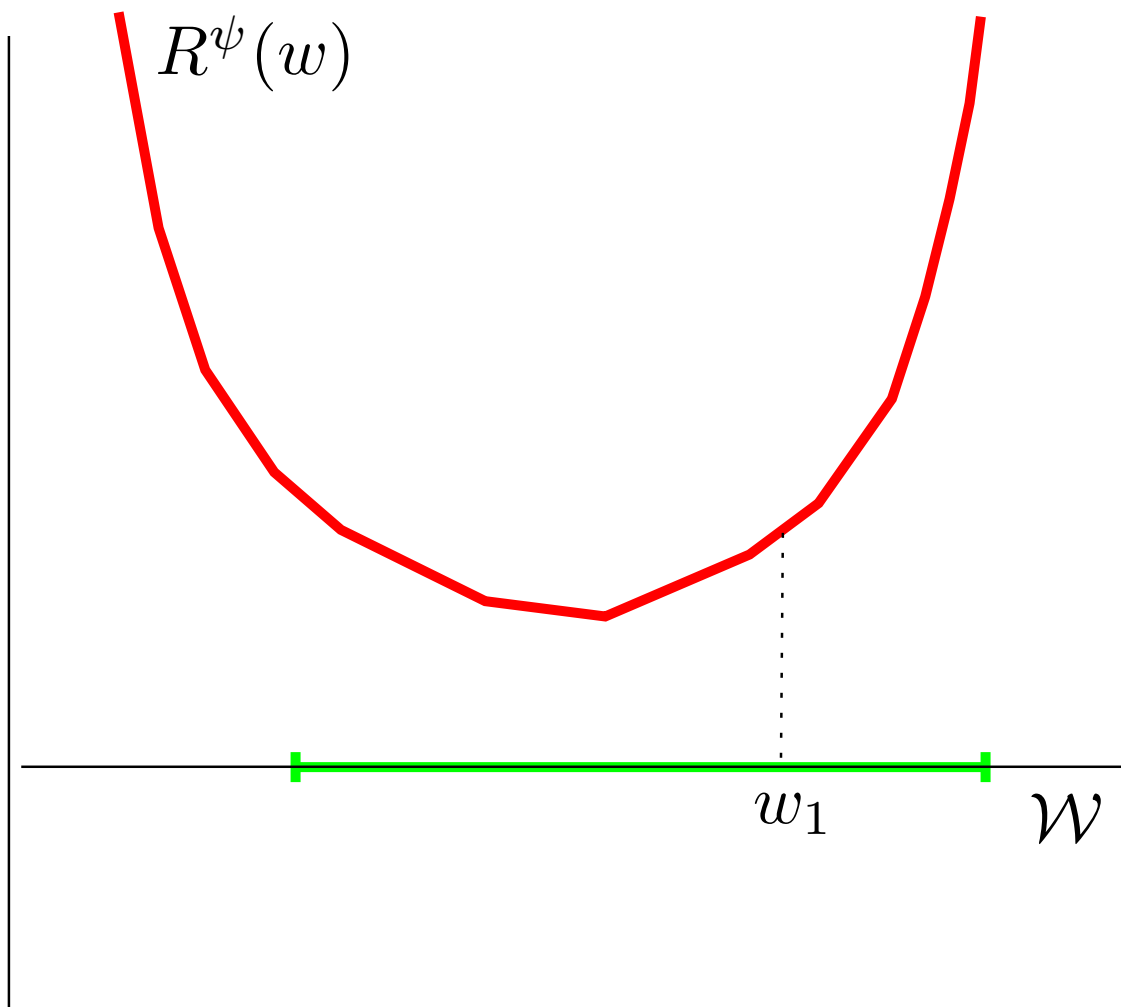
Cutting plane algorithm

$$R^\psi(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \max_{\hat{y}^i \in \mathcal{Y}} (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), \mathbf{w} \rangle) = \max_{\substack{\hat{y}^1 \in \mathcal{Y} \\ \vdots \\ \hat{y}^m \in \mathcal{Y}}} \frac{1}{m} \sum_{i=1}^m (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), \mathbf{w} \rangle)$$



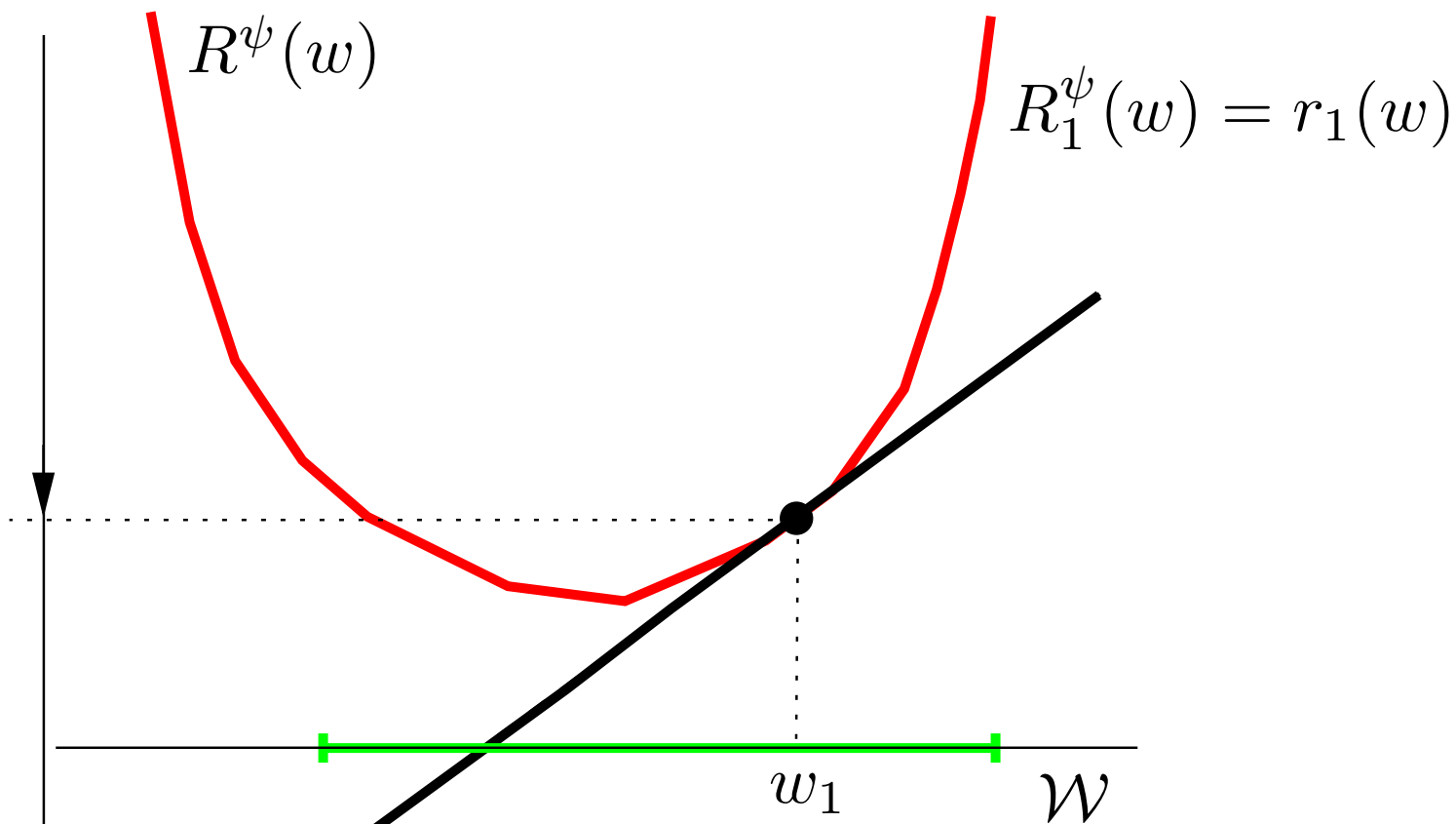
Cutting plane algorithm

$$R^\psi(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \max_{\hat{y}^i \in \mathcal{Y}} (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), \mathbf{w} \rangle) = \max_{\substack{\hat{y}^1 \in \mathcal{Y} \\ \vdots \\ \hat{y}^m \in \mathcal{Y}}} \frac{1}{m} \sum_{i=1}^m (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), \mathbf{w} \rangle)$$



Cutting plane algorithm

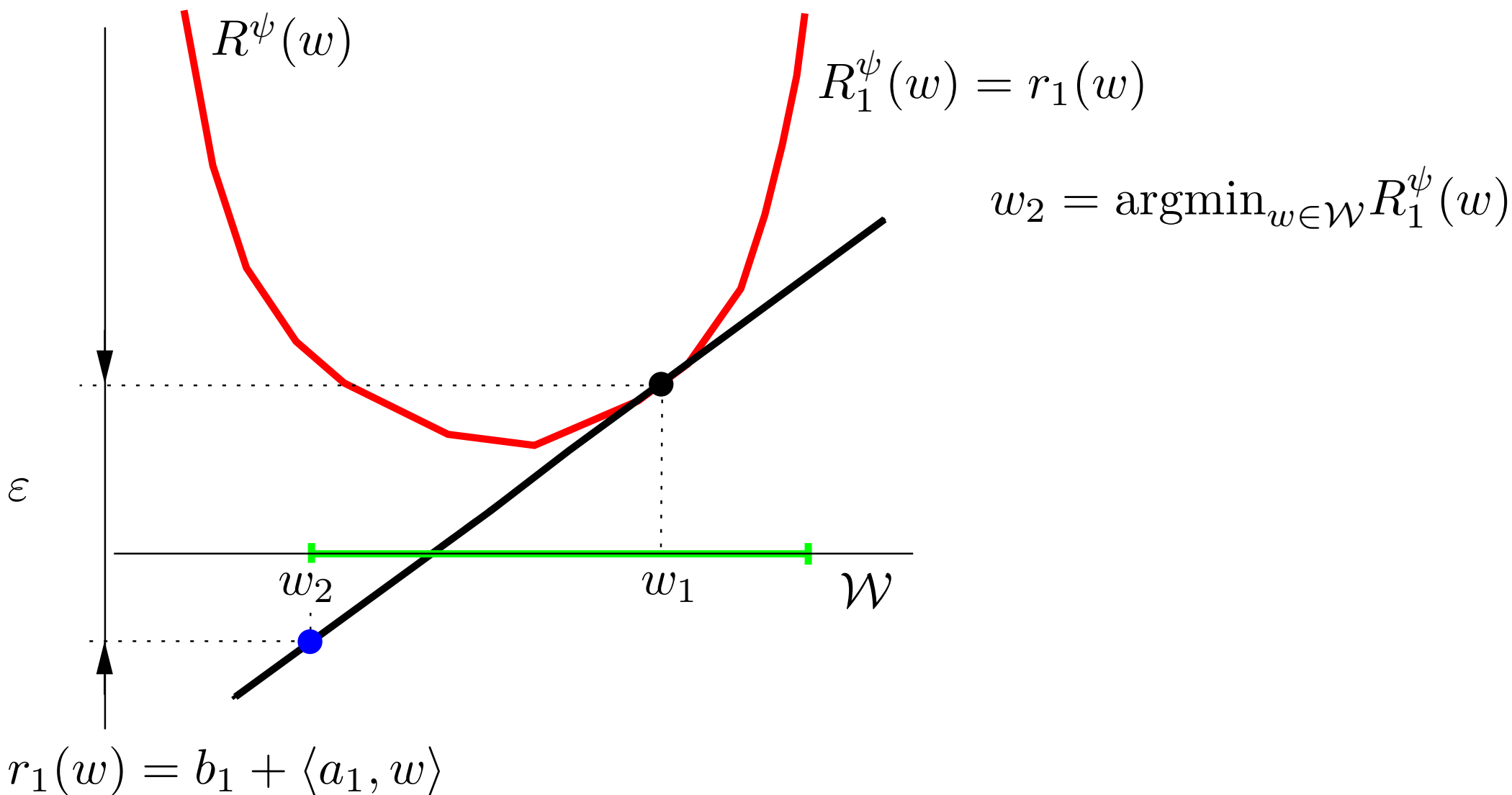
$$R^\psi(w) = \frac{1}{m} \sum_{i=1}^m \max_{\hat{y}^i \in \mathcal{Y}} (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), w \rangle) = \max_{\substack{\hat{y}^1 \in \mathcal{Y} \\ \vdots \\ \hat{y}^m \in \mathcal{Y}}} \frac{1}{m} \sum_{i=1}^m (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), w \rangle)$$



$$r_1(w) = b_1 + \langle a_1, w \rangle$$

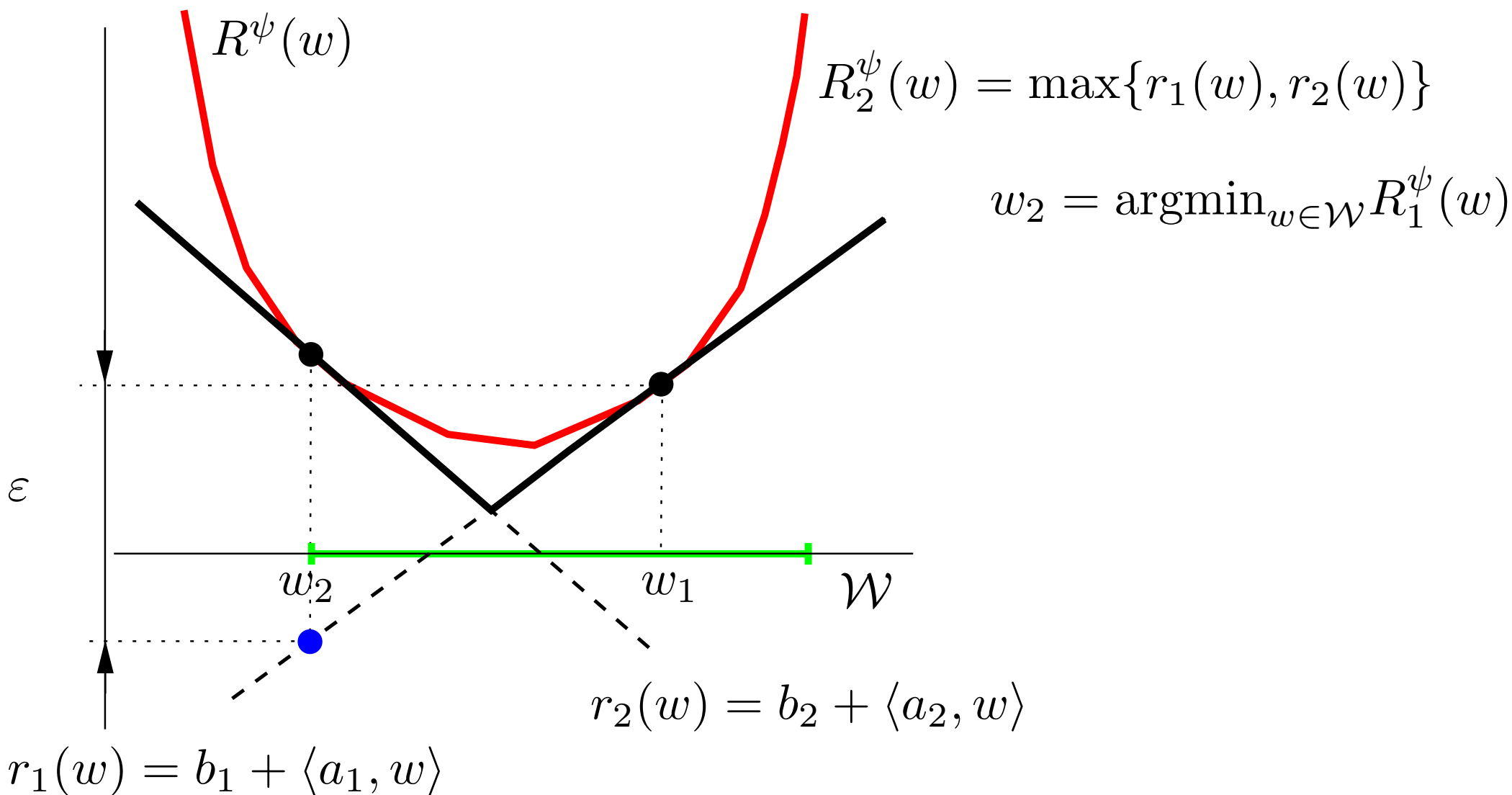
Cutting plane algorithm

$$R^\psi(w) = \frac{1}{m} \sum_{i=1}^m \max_{\hat{y}^i \in \mathcal{Y}} (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), w \rangle) = \max_{\substack{\hat{y}^1 \in \mathcal{Y} \\ \vdots \\ \hat{y}^m \in \mathcal{Y}}} \frac{1}{m} \sum_{i=1}^m (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), w \rangle)$$



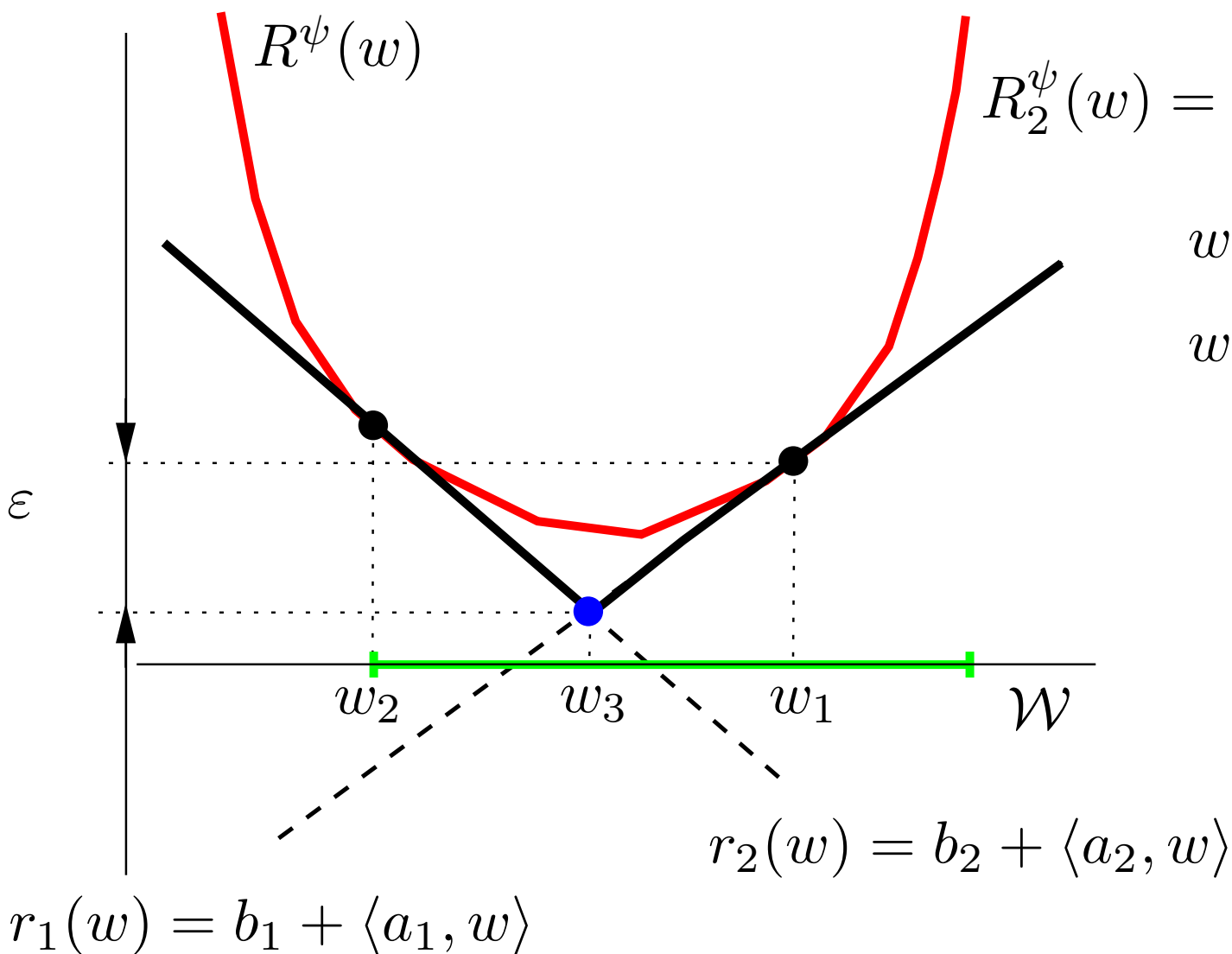
Cutting plane algorithm

$$R^\psi(w) = \frac{1}{m} \sum_{i=1}^m \max_{\hat{y}^i \in \mathcal{Y}} (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), w \rangle) = \max_{\substack{\hat{y}^1 \in \mathcal{Y} \\ \vdots \\ \hat{y}^m \in \mathcal{Y}}} \frac{1}{m} \sum_{i=1}^m (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), w \rangle)$$



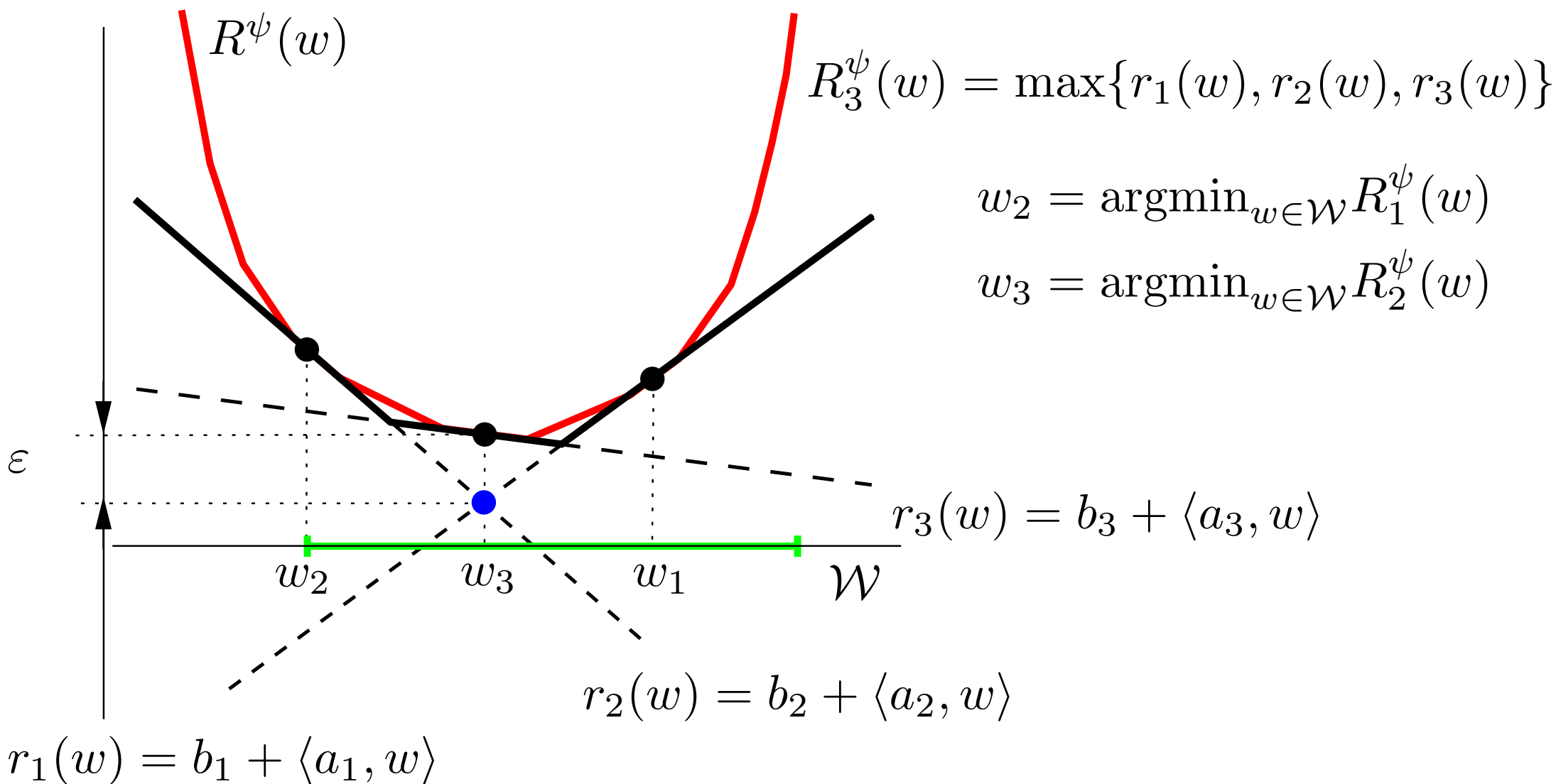
Cutting plane algorithm

$$R^\psi(w) = \frac{1}{m} \sum_{i=1}^m \max_{\hat{y}^i \in \mathcal{Y}} (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), w \rangle) = \max_{\substack{\hat{y}^1 \in \mathcal{Y} \\ \vdots \\ \hat{y}^m \in \mathcal{Y}}} \frac{1}{m} \sum_{i=1}^m (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), w \rangle)$$



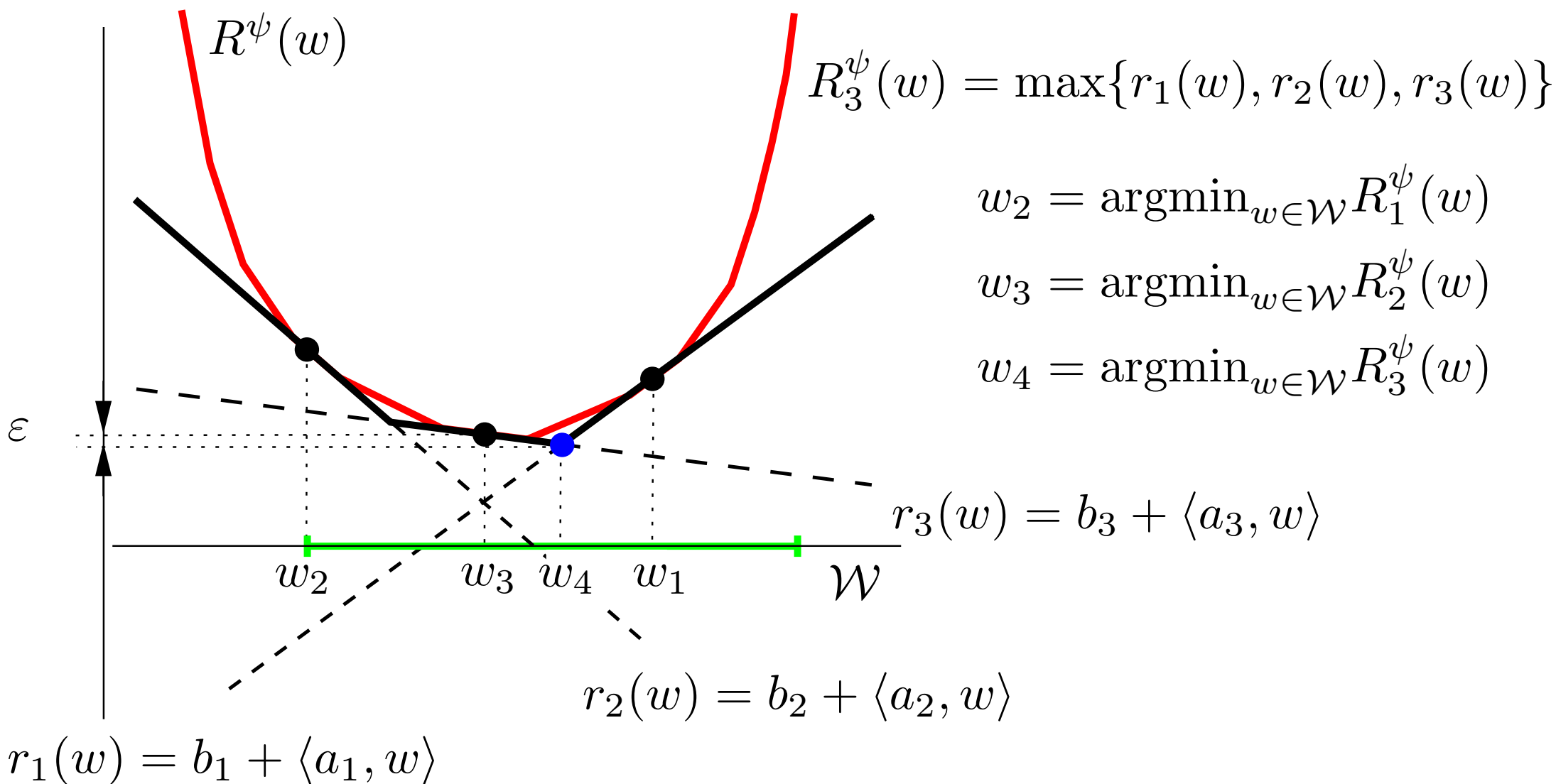
Cutting plane algorithm

$$R^\psi(w) = \frac{1}{m} \sum_{i=1}^m \max_{\hat{y}^i \in \mathcal{Y}} (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), w \rangle) = \max_{\substack{\hat{y}^1 \in \mathcal{Y} \\ \vdots \\ \hat{y}^m \in \mathcal{Y}}} \frac{1}{m} \sum_{i=1}^m (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), w \rangle)$$



Cutting plane algorithm

$$R^\psi(w) = \frac{1}{m} \sum_{i=1}^m \max_{\hat{y}^i \in \mathcal{Y}} (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), w \rangle) = \max_{\substack{\hat{y}^1 \in \mathcal{Y} \\ \vdots \\ \hat{y}^m \in \mathcal{Y}}} \frac{1}{m} \sum_{i=1}^m (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), w \rangle)$$



Cutting plane algorithm

1. $\mathbf{w}_1 \in \mathcal{W}$, $t \leftarrow 1$

2. Compute a new cutting plane and the objective value:

$$\mathbf{a}_t = \frac{1}{m} \sum_{i=1}^m \phi_i(\hat{y}^i), \quad b_t = \frac{1}{m} \sum_{i=1}^m \ell_i(\hat{y}^i), \quad R^\psi(\mathbf{w}_t) = b_t + \langle \mathbf{w}_t, \mathbf{a}_t \rangle$$

where \hat{y}^i is a solutions of **loss augmented prediction** problem:

$$\hat{y}^i = \operatorname{argmax}_{y \in \mathcal{Y}} (\ell_i(y) + \langle \mathbf{w}, \phi_i(y) \rangle) = \operatorname{argmax}_{y \in \mathcal{Y}} (\ell(y^i, y) + \langle \mathbf{w}, \phi(x^i, y) \rangle)$$

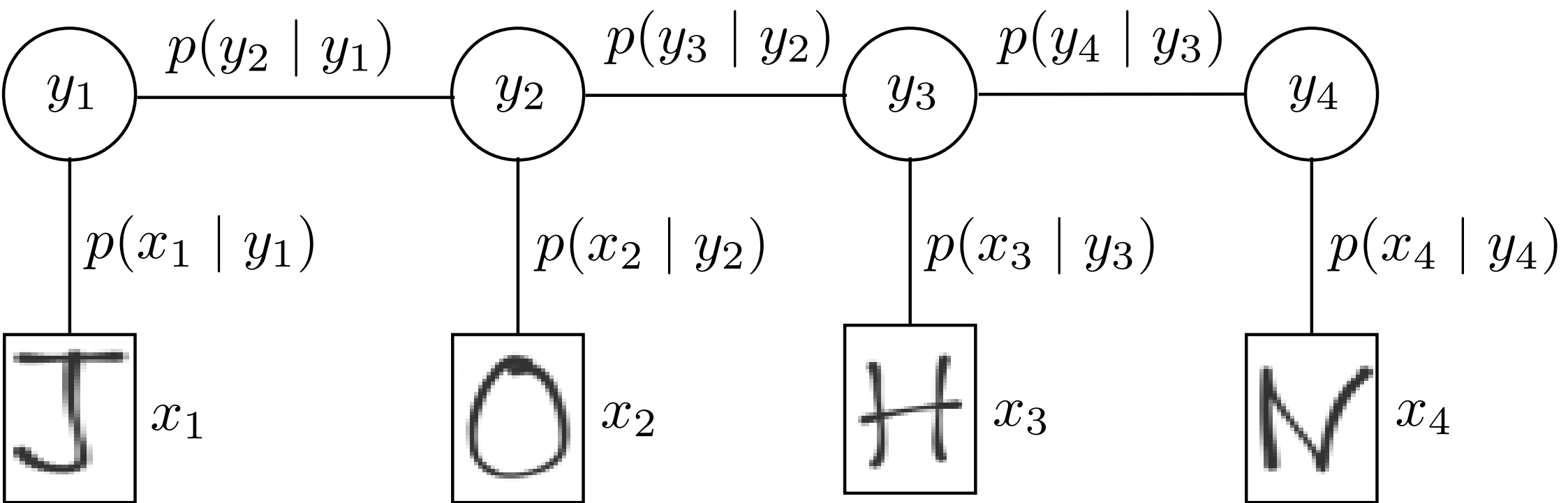
3. Solve a reduced problem

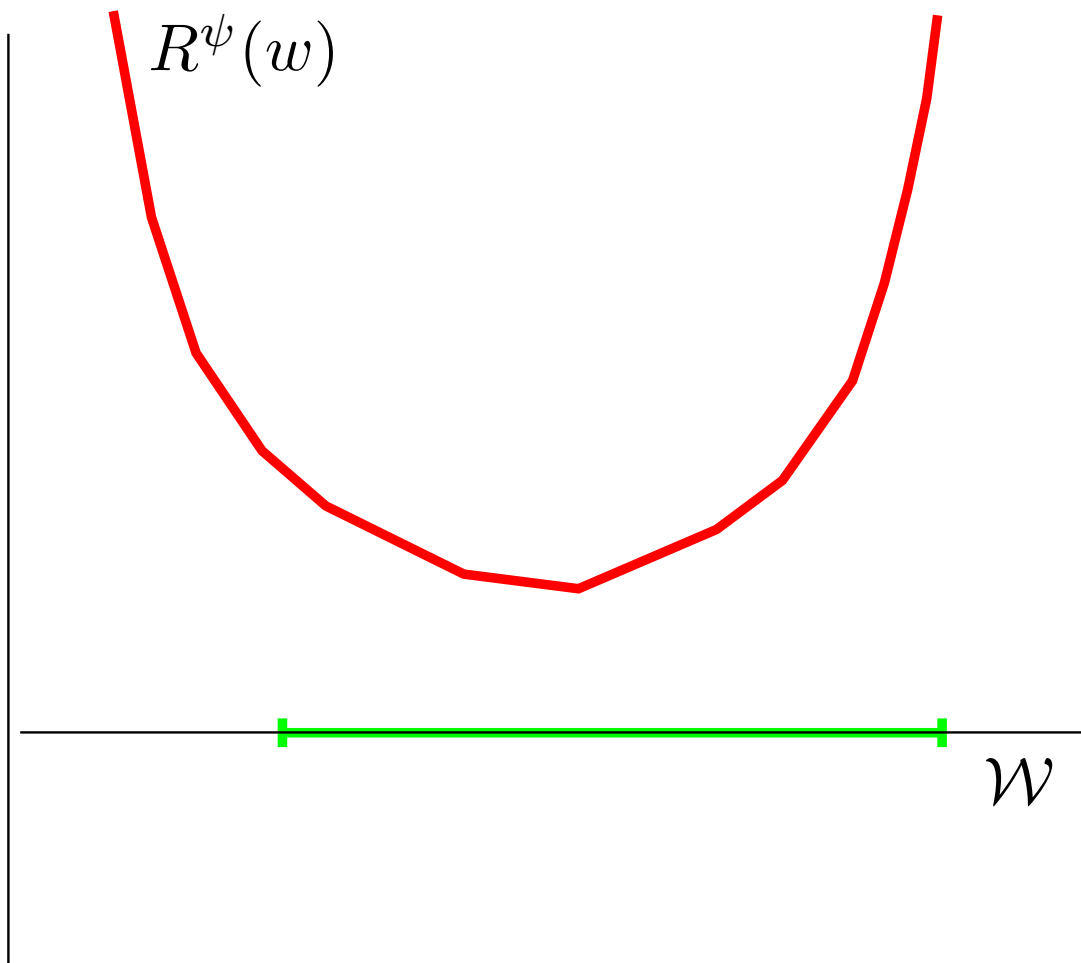
$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} R_t^\psi(\mathbf{w}), \quad \text{where} \quad R_t^\psi(\mathbf{w}) = \max_{i=1, \dots, t} (b_i + \langle \mathbf{w}, \mathbf{a}_i \rangle)$$

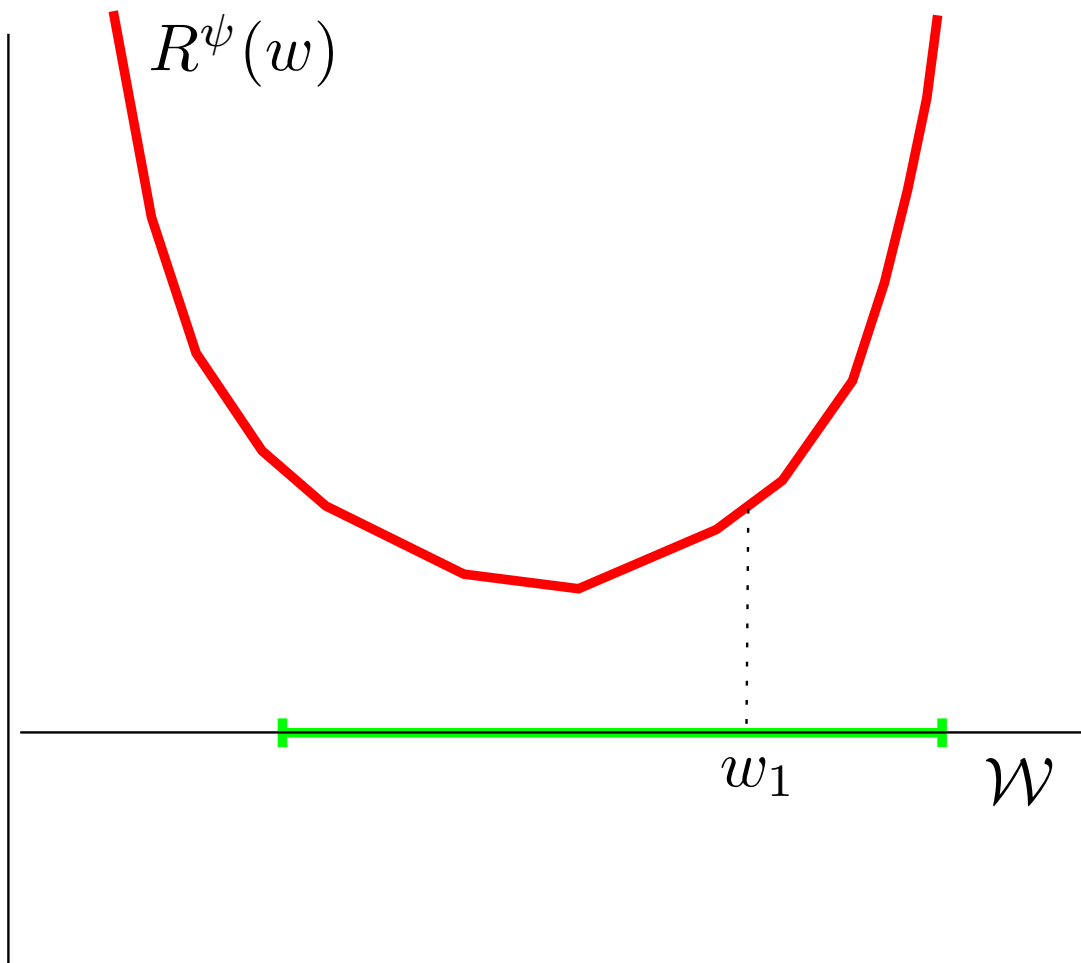
4. If $\min_{i=1, \dots, t+1} R(\mathbf{w}_t) - R^\psi(\mathbf{w}_{t+1}) \leq \varepsilon$ exit else $t \leftarrow t + 1$ and go to 2.

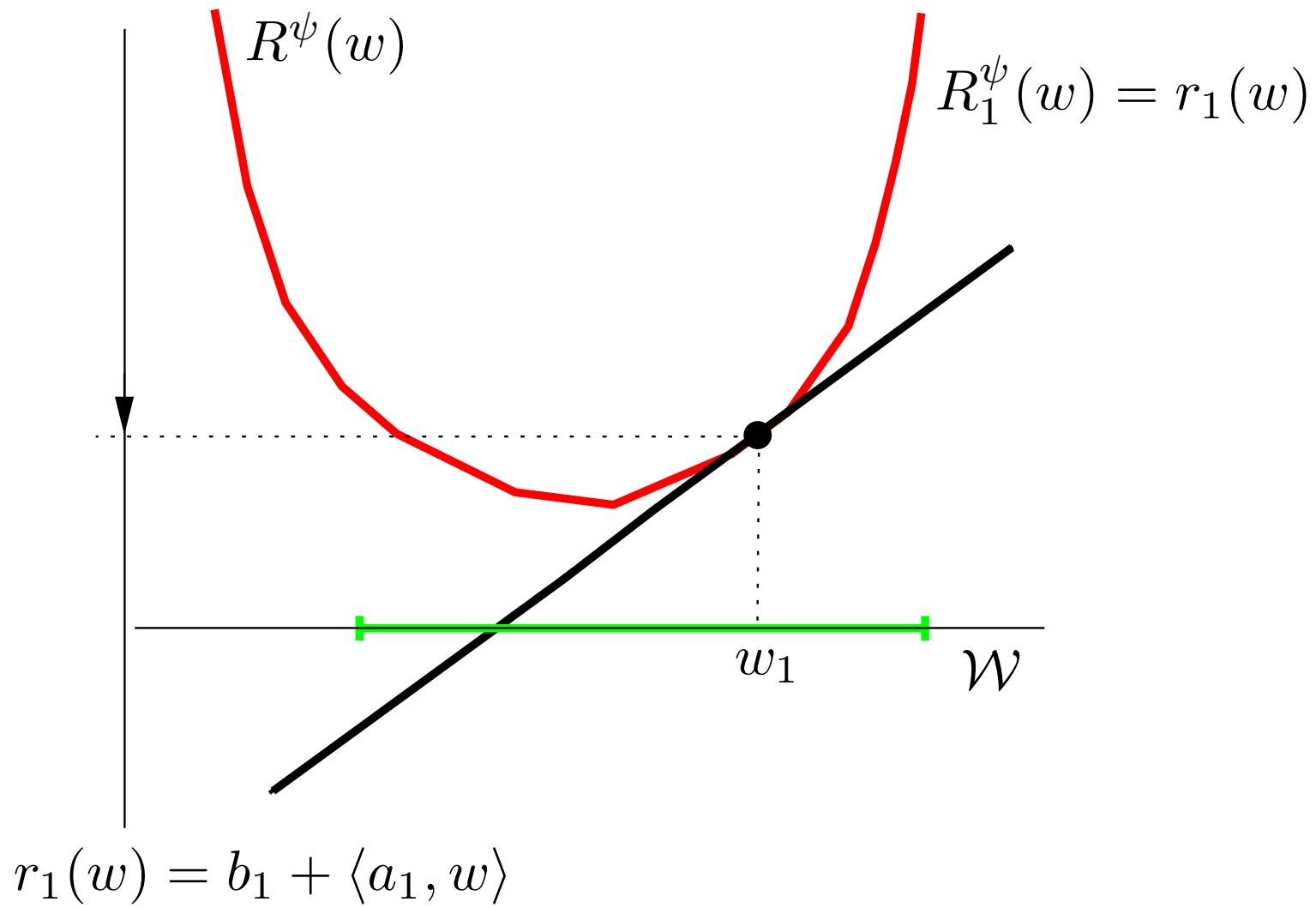
Summary

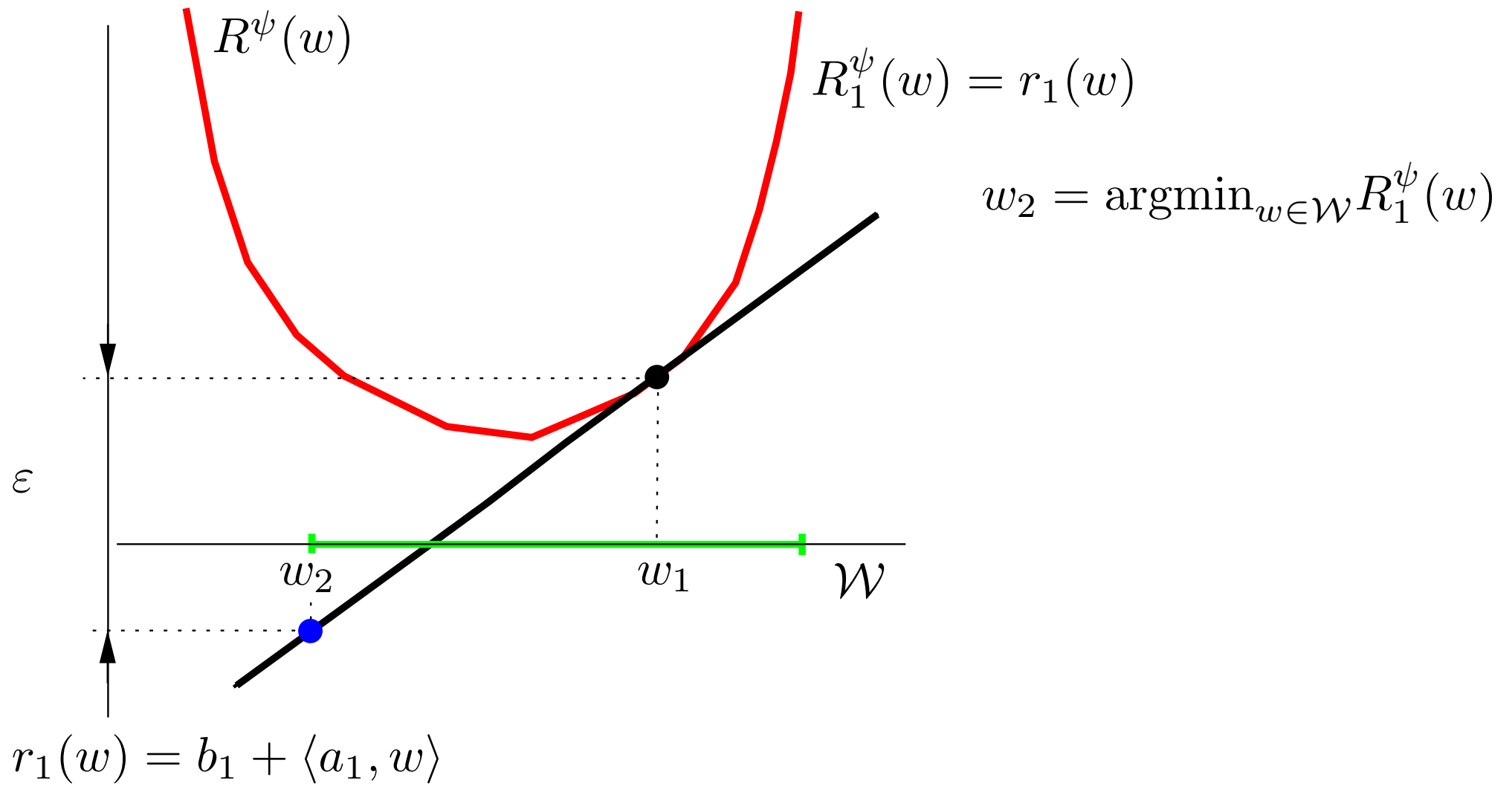
- ◆ Generalized linear classifier
- ◆ Structured Output Perceptron
- ◆ Structured Output Support Vector Machines
- ◆ Cutting Plane Algorithm

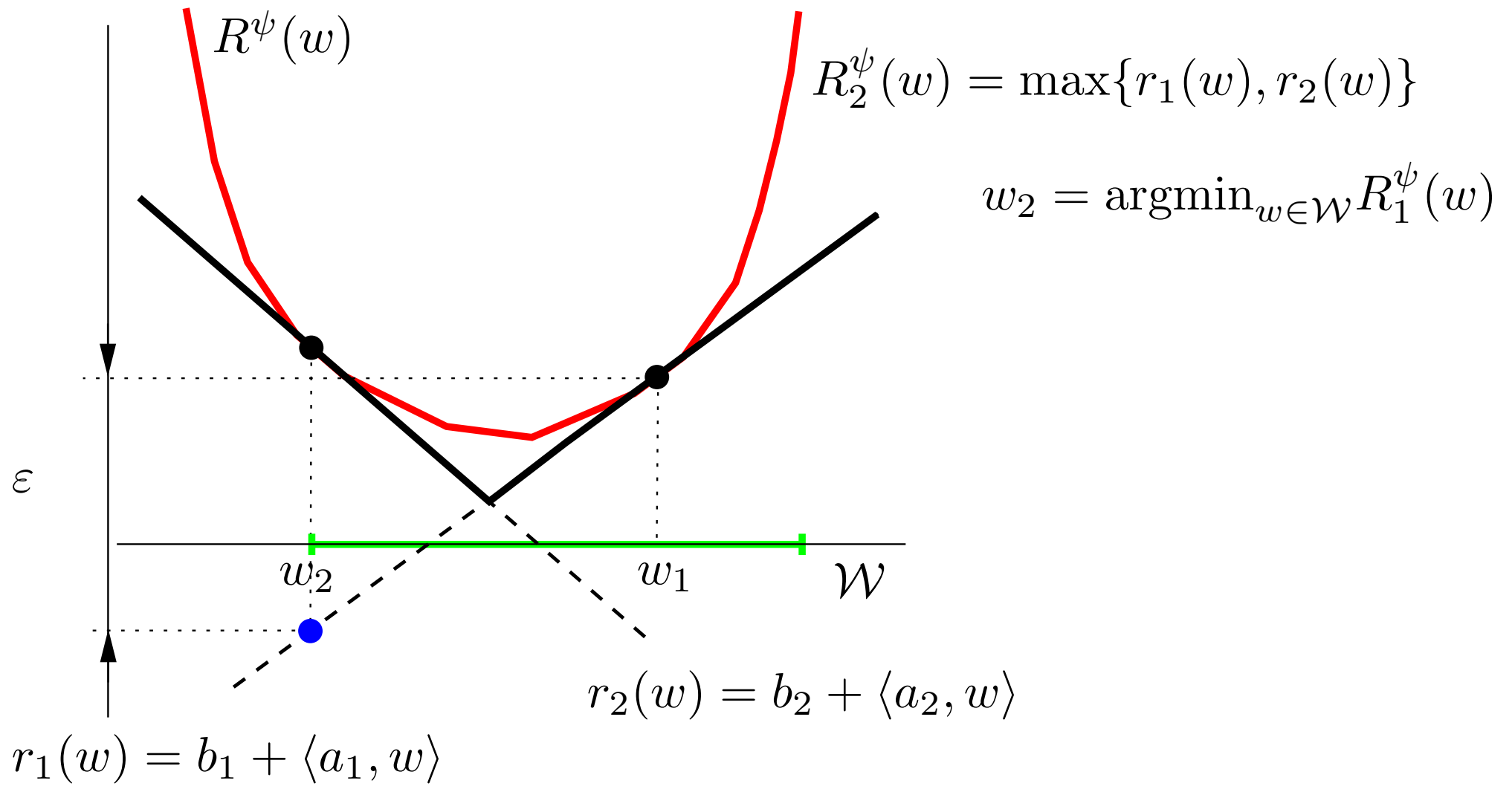


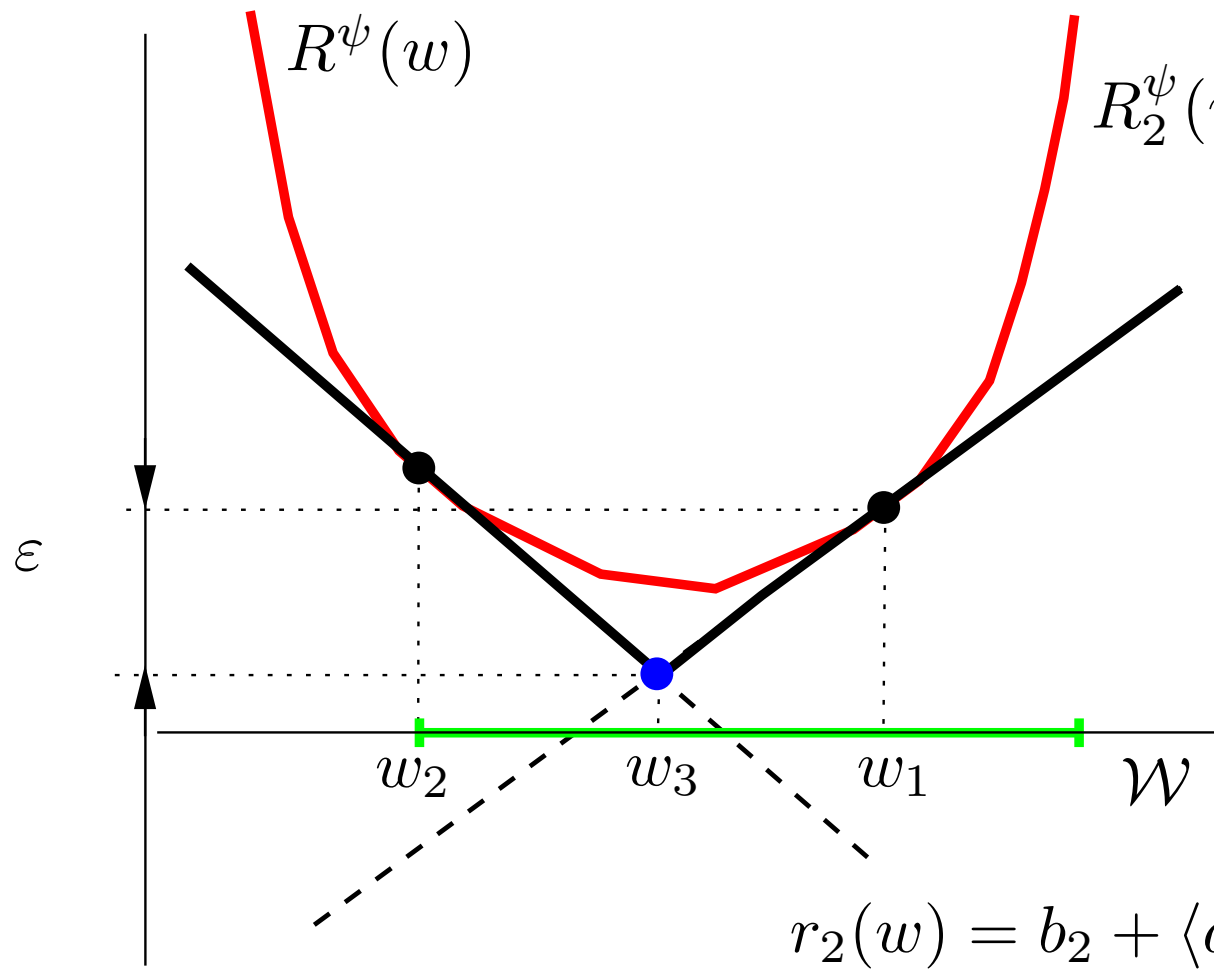












$$r_1(w) = b_1 + \langle a_1, w \rangle$$

$$r_2(w) = b_2 + \langle a_2, w \rangle$$

$$R_2^\psi(w) = \max\{r_1(w), r_2(w)\}$$

$$w_2 = \operatorname{argmin}_{w \in \mathcal{W}} R_1^\psi(w)$$

$$w_3 = \operatorname{argmin}_{w \in \mathcal{W}} R_2^\psi(w)$$

