

# Cutting Plane Algorithm

Czech Technical University in Prague  
V. Franc

February 24, 2020

- ◆ Cutting Plane Algorithm
- ◆ Bundle Method for Risk Minimization
- ◆ Subgradients

**XEP33SML – Structured Model Learning, Summer 2020**

## Structured Output SVM

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

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

## Structured Output SVM

- ◆ Learning  $h(x; \mathbf{w}) = \text{Argmax}_{y \in \mathcal{Y}} \langle \mathbf{w}, \phi(x, y) \rangle$  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

$$\mathbf{w}^* \in \underset{\mathbf{w} \in \mathbb{R}^n}{\text{Argmin}} \left( \frac{\lambda}{2} \|\mathbf{w}\|^2 + R^\psi(\mathbf{w}) \right) \quad \text{where} \quad R^\psi(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \psi(x^i, y^i, \mathbf{w})$$

## Structured Output SVM

- ◆ Learning  $h(x; \mathbf{w}) = \text{Argmax}_{y \in \mathcal{Y}} \langle \mathbf{w}, \phi(x, y) \rangle$  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

$$\mathbf{w}^* \in \underset{\mathbf{w} \in \mathbb{R}^n}{\text{Argmin}} \left( \frac{\lambda}{2} \|\mathbf{w}\|^2 + R^\psi(\mathbf{w}) \right) \quad \text{where} \quad R^\psi(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \psi(x^i, y^i, \mathbf{w})$$

- ◆ The surrogate loss  $\psi: \mathcal{X} \times \mathcal{Y} \times \mathbb{R}^n \rightarrow \mathbb{R}$  is an upper bound:

$$\ell(y, h(x; \mathbf{w})) \leq \psi(x, y, \mathbf{w}), \quad \forall (x, y, \mathbf{w}) \in (\mathcal{X} \times \mathcal{Y} \times \mathbb{R}^n)$$

which is convex in  $\mathbf{w}$  for any  $(x, y)$ .

## SO-SVM leads to a convex QP

- ◆ The SO-SVM with margin-rescaling loss:

$$\mathbf{w}^* \in \underset{\mathbf{w} \in \mathbb{R}^n}{\text{Argmin}} \left( \frac{\lambda}{2} \|\mathbf{w}\|^2 + \underbrace{\frac{1}{m} \sum_{i=1}^m \max_{y \in \mathcal{Y}} \{ \ell_i(y) + \langle \mathbf{w}, \phi_i(y) \rangle \}}_{R^\psi(\mathbf{w})} \right)$$

## SO-SVM leads to a convex QP

- ◆ The SO-SVM with margin-rescaling loss:

$$\mathbf{w}^* \in \underset{\mathbf{w} \in \mathbb{R}^n}{\text{Argmin}} \left( \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \max_{y \in \mathcal{Y}} \{ \ell_i(y) + \langle \mathbf{w}, \phi_i(y) \rangle \}}_{R\psi(\mathbf{w})} \right)$$

- ◆ By using slack variables it can be rewritten as a 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

- ◆ The SO-SVM with margin-rescaling loss:

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

- ◆ Equivalent formulation: for any  $\lambda > 0$  there exists  $r > 0$  such that

$$\mathbf{w}^* \in \underset{\mathbf{w} \in \mathcal{W}}{\text{Argmin}} R^\psi(\mathbf{w}) \tag{1}$$

where  $\mathcal{W} = \{\mathbf{w} \in \mathbb{R}^n \mid \|\mathbf{w}\| \leq r\}$  is a ball of radius  $r$ .

## Cutting Plane Algorithm

- ◆ The SO-SVM with margin-rescaling loss:

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

- ◆ Equivalent formulation: for any  $\lambda > 0$  there exists  $r > 0$  such that

$$\mathbf{w}^* \in \underset{\mathbf{w} \in \mathcal{W}}{\text{Argmin}} R^\psi(\mathbf{w}) \quad (1)$$

where  $\mathcal{W} = \{\mathbf{w} \in \mathbb{R}^n \mid \|\mathbf{w}\| \leq r\}$  is a ball of radius  $r$ .

- ◆ CP algorithm: approximate (1) by a series of simpler problems

$$\mathbf{w}_{t+1} \in \underset{\mathbf{w} \in \mathcal{W}}{\text{Argmin}} R_t^\psi(\mathbf{w}), \quad t = 1, 2, \dots$$

where  $R_t^\psi(\mathbf{w})$  is a successively tighter lower bound of  $R^\psi(\mathbf{w})$ .

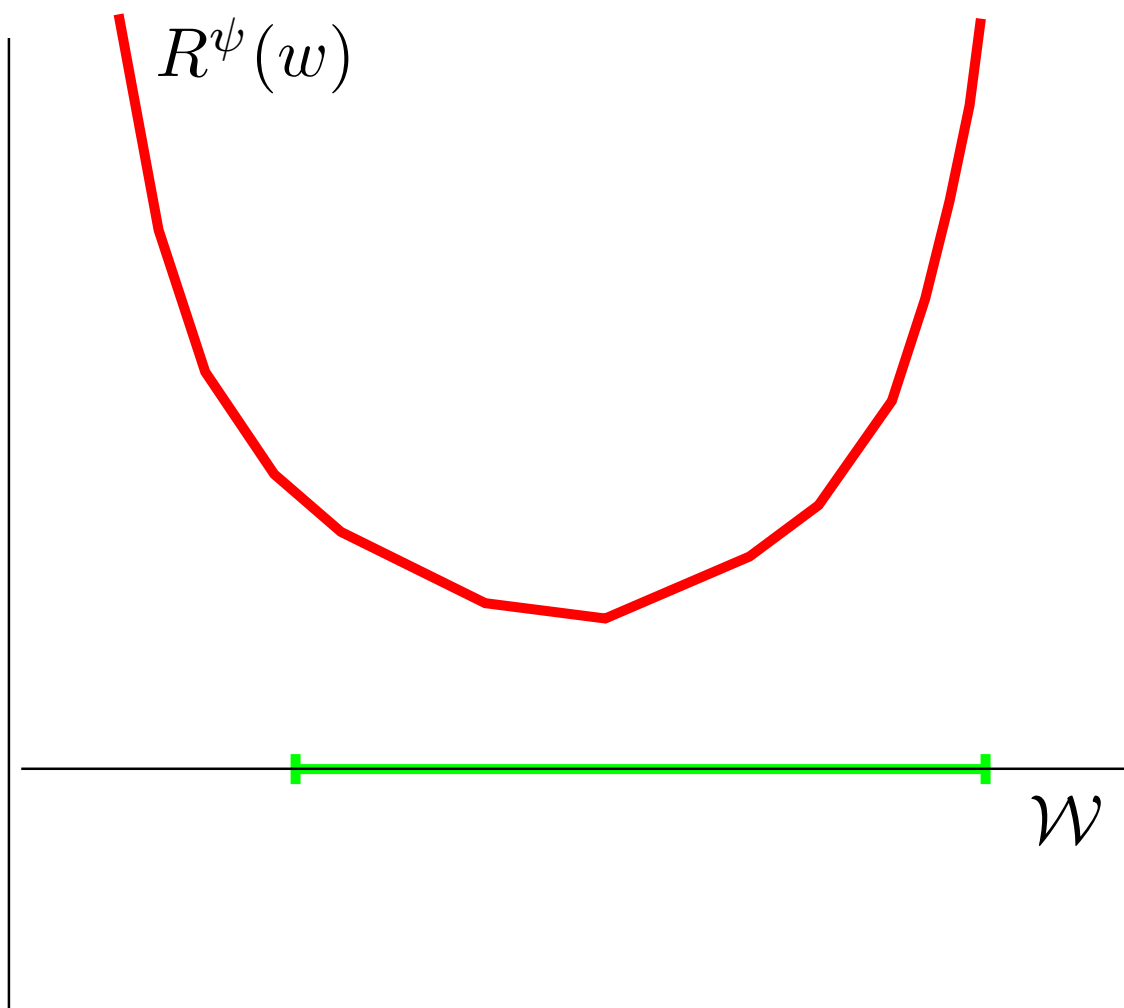


# Cutting Plane Algorithm

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

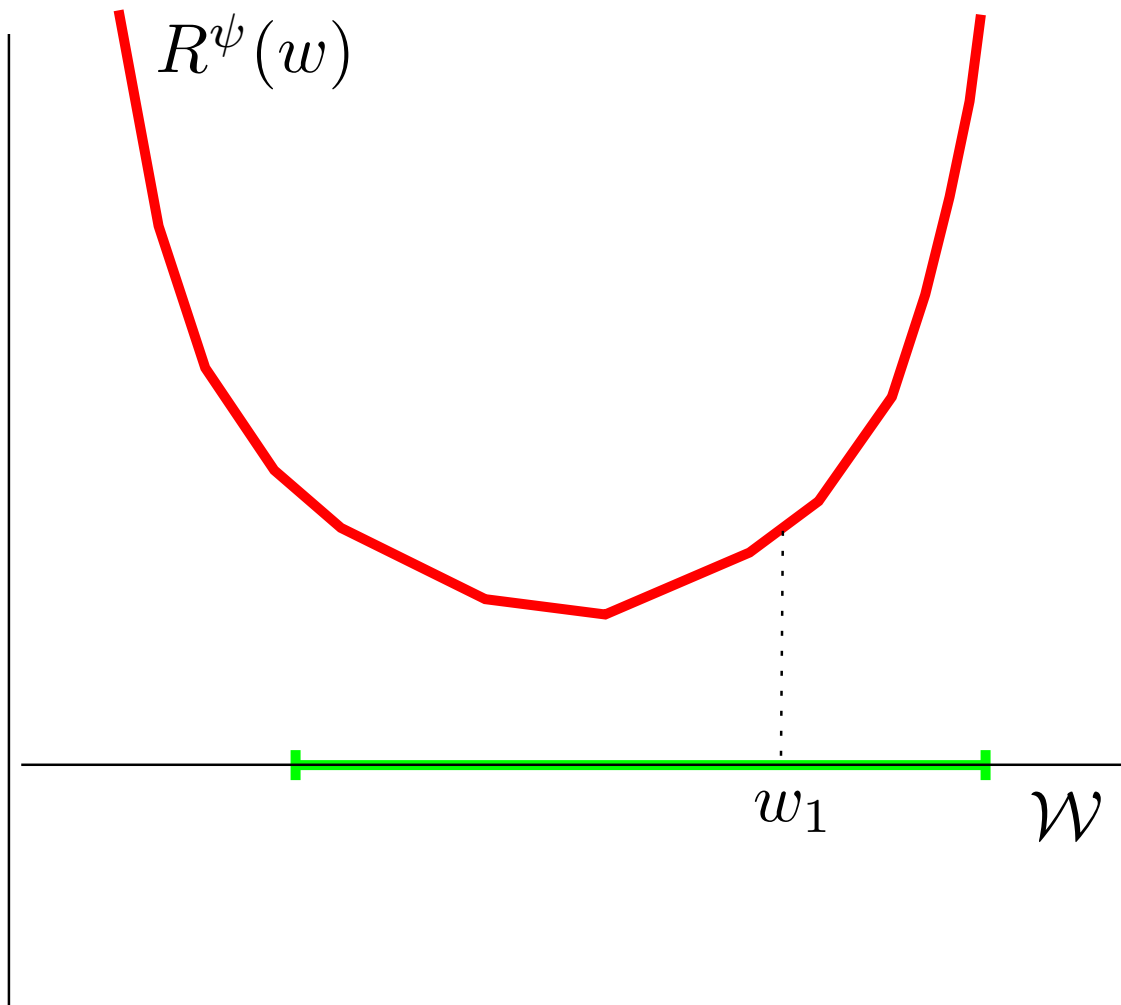
# Cutting Plane Algorithm

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



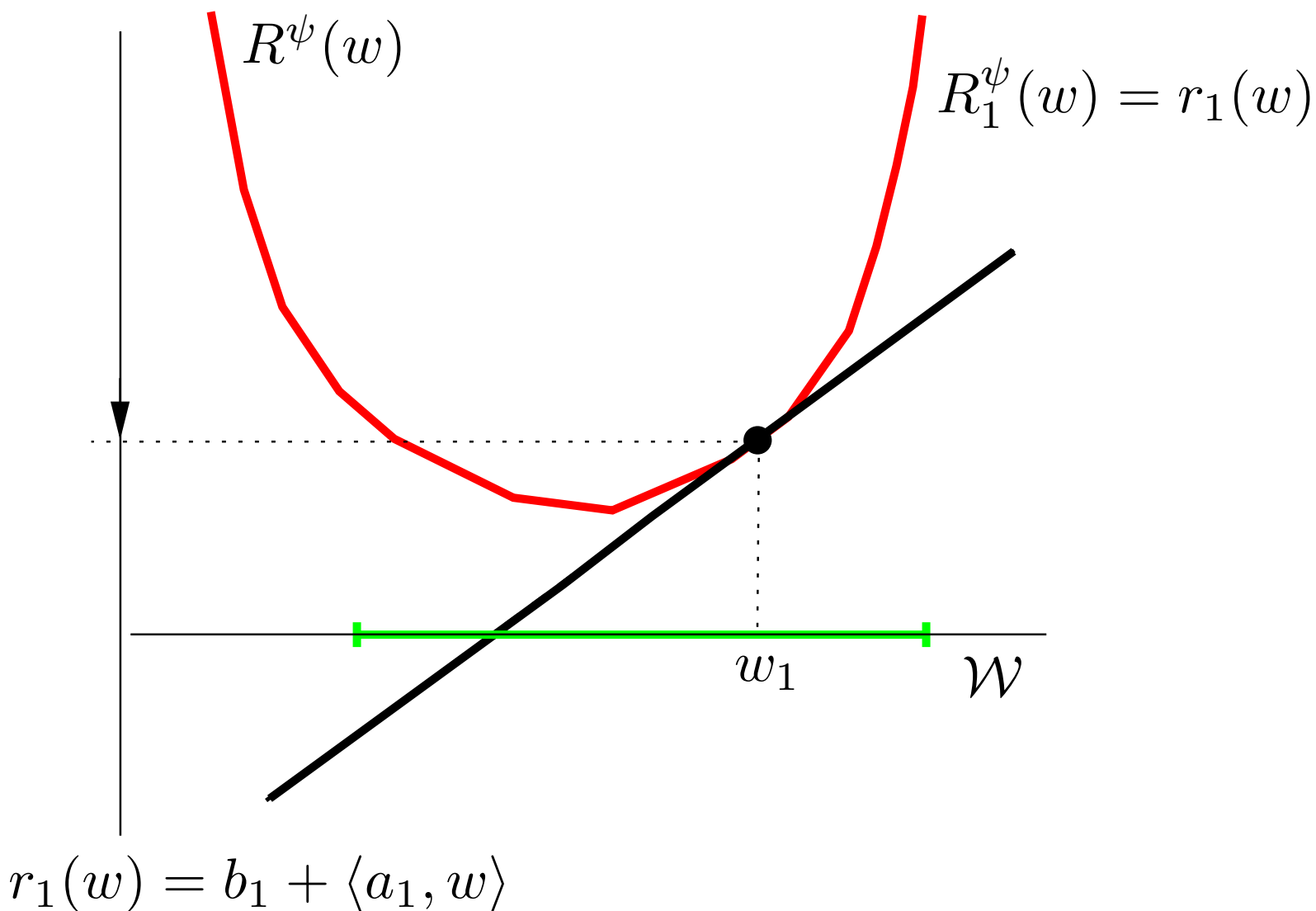
# Cutting Plane Algorithm

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



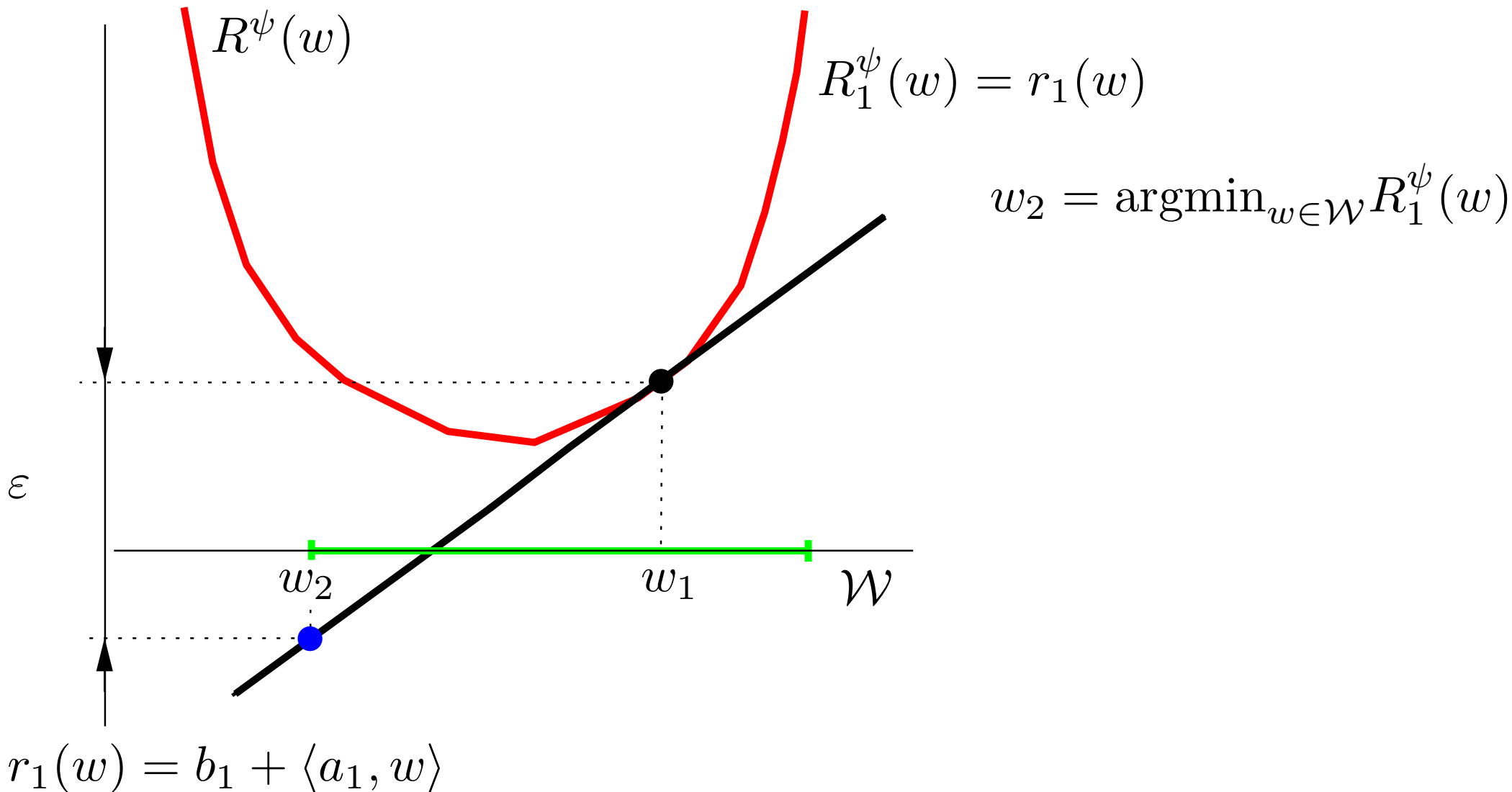
# Cutting Plane Algorithm

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



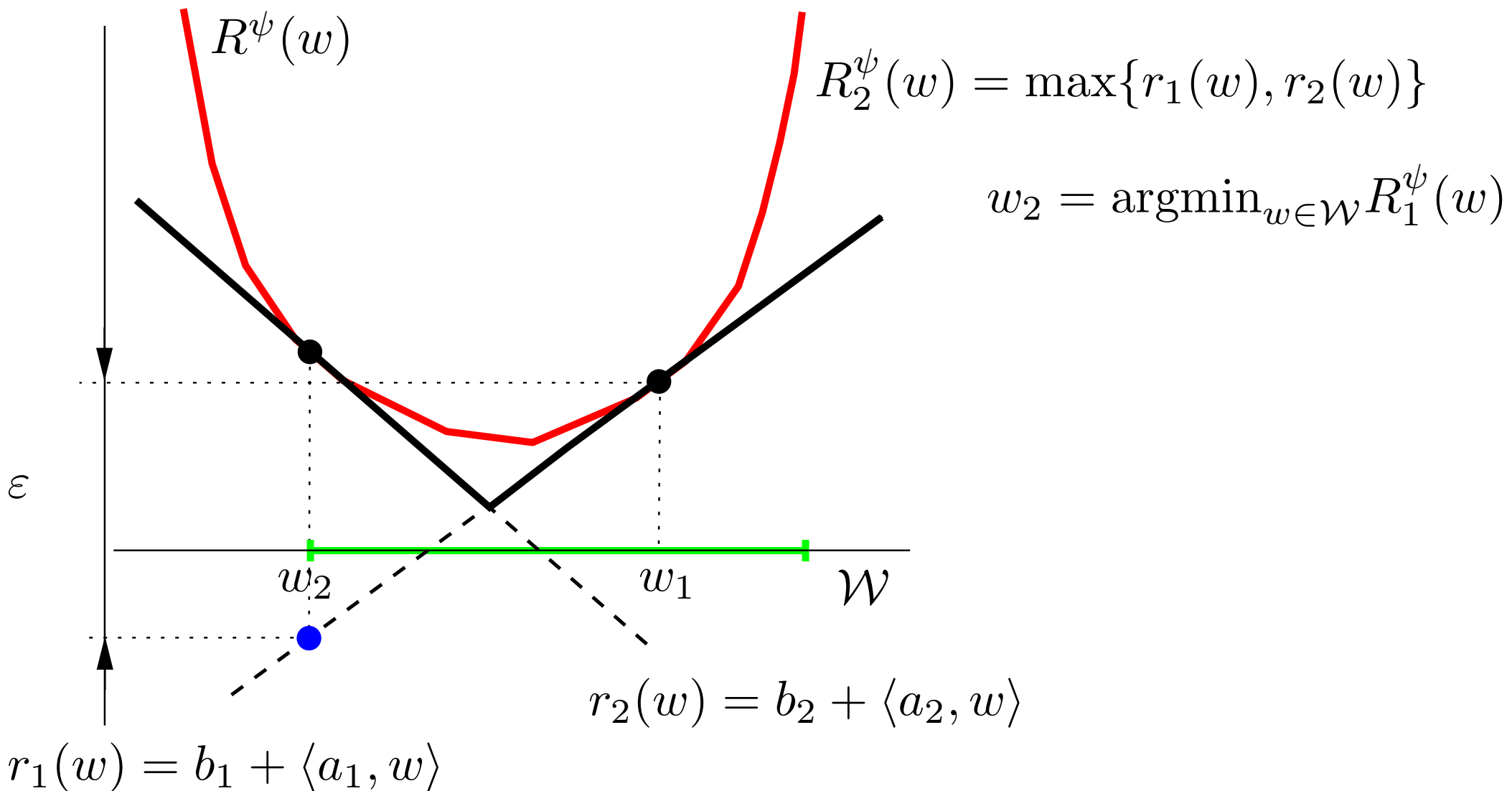
# Cutting Plane Algorithm

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



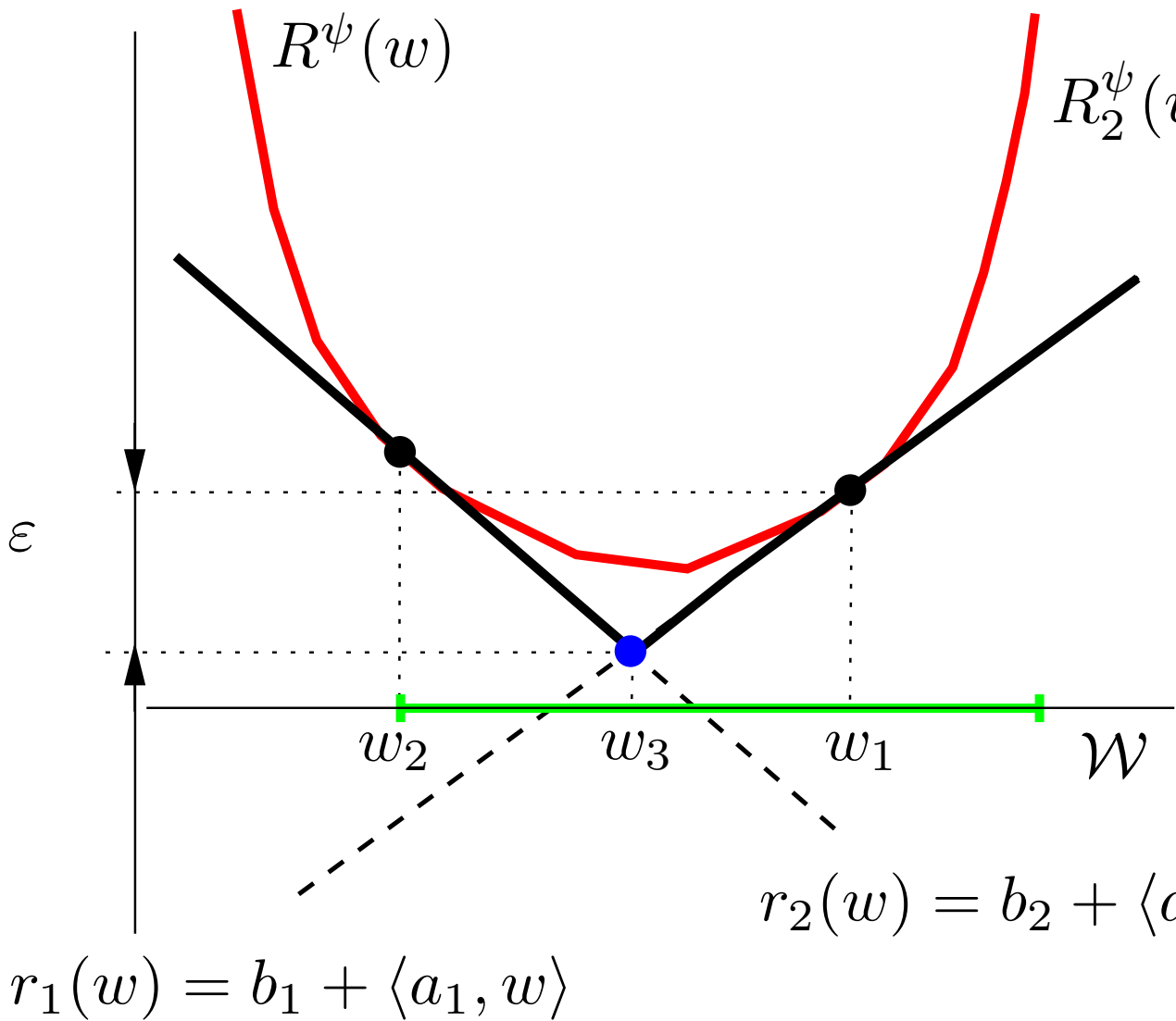
# Cutting Plane Algorithm

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



# Cutting Plane Algorithm

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



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

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

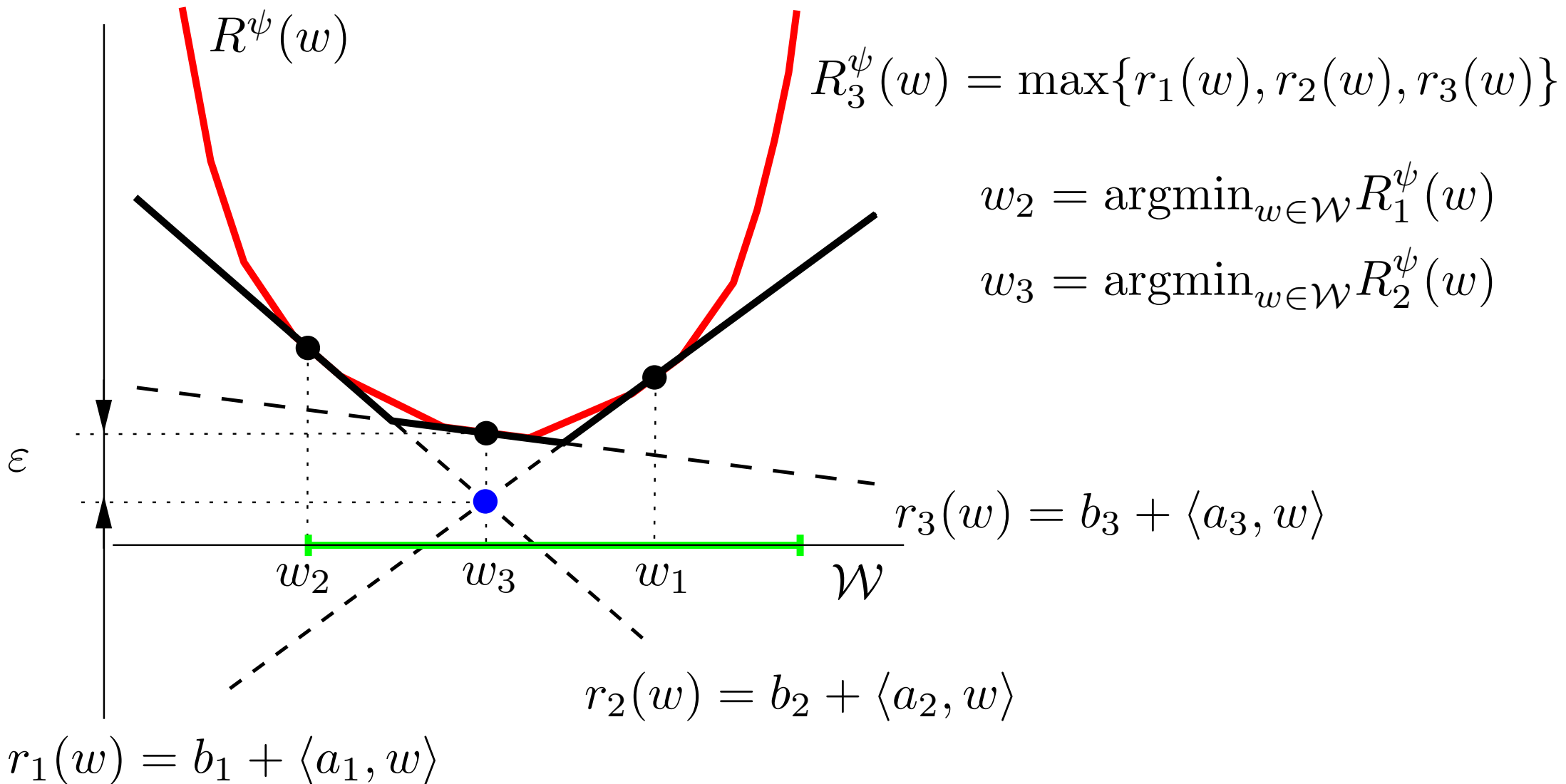
$$w_3 = \operatorname{argmin}_{w \in \mathcal{W}} R_2^\psi(w)$$

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

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

# Cutting Plane Algorithm

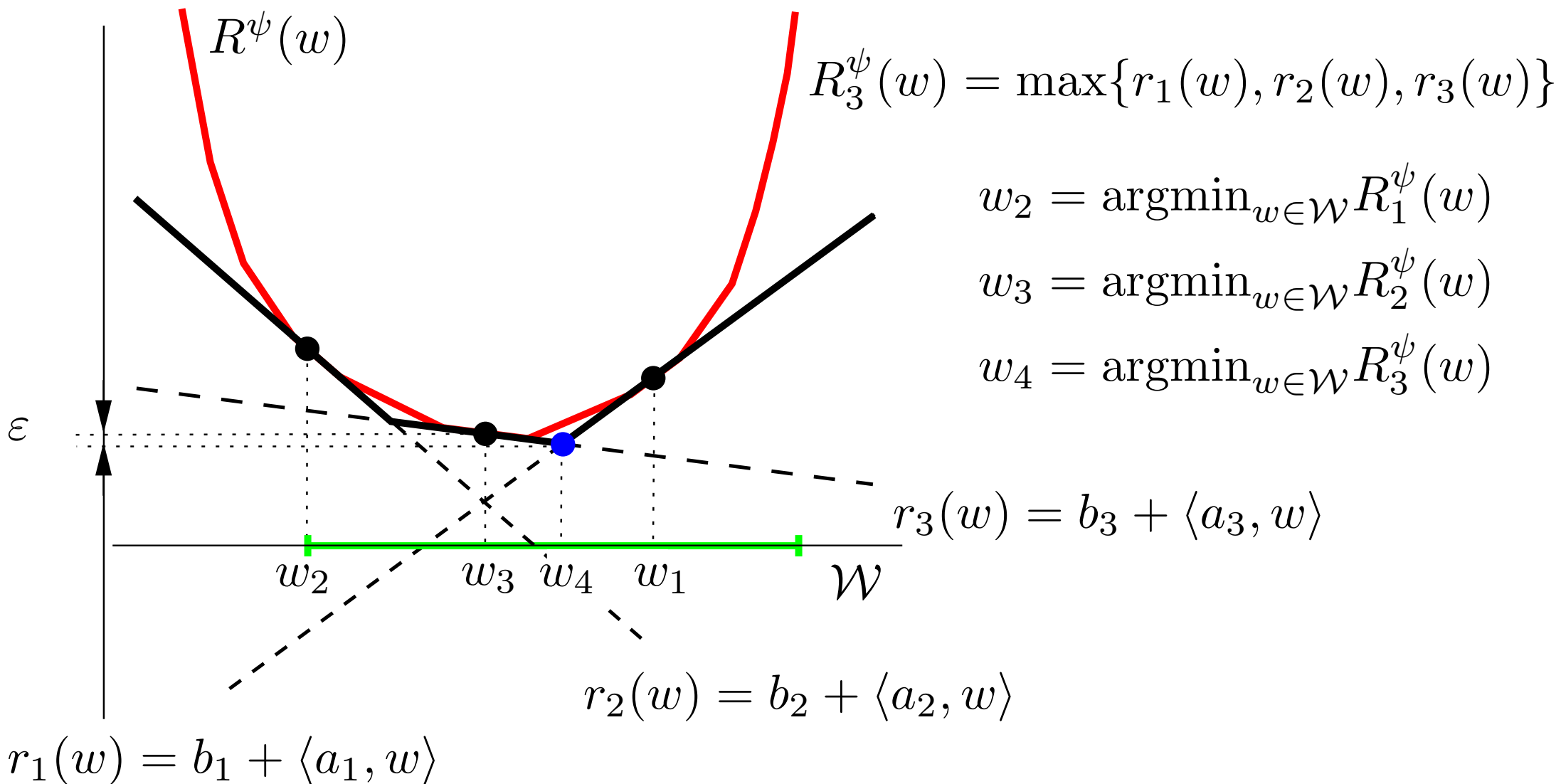
$$R^\psi(w) = \frac{1}{m} \sum_{i=1}^m \max_{\hat{y}^i \in \mathcal{Y}} (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), w \rangle) = \max_{\substack{\hat{y}^1 \in \mathcal{Y} \\ \vdots \\ \hat{y}^m \in \mathcal{Y}}} \frac{1}{m} \sum_{i=1}^m (\ell_i(\hat{y}^i) + \langle \phi_i(\hat{y}^i), w \rangle)$$





# 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

1.  $\mathbf{w}_1 \in \mathcal{W} = \{\mathbf{w} \in \mathbb{R}^n \mid \|\mathbf{w}\| \leq r\}$ ,  $t \leftarrow 1$

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

## Bundle Method for Risk Minimization

1.  $\mathbf{w}_1 \in \mathbb{R}^n$ ,  $t \leftarrow 1$
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

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \frac{\lambda}{2} \|\mathbf{w}\|^2 + R_t^\psi(\mathbf{w}) \right) \quad \text{where} \quad R_t^\psi(\mathbf{w}) = \max_{i=1, \dots, t} (b_i + \langle \mathbf{w}, \mathbf{a}_i \rangle)$$

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

## Bundle Method for Risk Minimization

- ◆ The original convex problem

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} F(\mathbf{w}) := \left( \frac{\lambda}{2} \|\mathbf{w}\|^2 + R(\mathbf{w}) \right)$$

## Bundle Method for Risk Minimization

- ◆ The original convex problem

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} F(\mathbf{w}) := \left( \frac{\lambda}{2} \|\mathbf{w}\|^2 + R(\mathbf{w}) \right)$$

is reduced to a sequence of reduced convex problems

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} F_t(\mathbf{w}) := \left( \frac{\lambda}{2} \|\mathbf{w}\|^2 + R_t(\mathbf{w}) \right)$$

## Bundle Method for Risk Minimization

- ◆ The original convex problem

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} F(\mathbf{w}) := \left( \frac{\lambda}{2} \|\mathbf{w}\|^2 + R(\mathbf{w}) \right)$$

is reduced to a sequence of reduced convex problems

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} F_t(\mathbf{w}) := \left( \frac{\lambda}{2} \|\mathbf{w}\|^2 + R_t(\mathbf{w}) \right)$$

where  $R_t(\mathbf{w})$  is the “cutting plane model”

$$R_t(\mathbf{w}) = \max_{i=0, \dots, t} \left[ R(\mathbf{w}_i) + \langle \mathbf{g}_i, \mathbf{w} - \mathbf{w}_i \rangle \right]$$

and  $\mathbf{g}_i = \partial R(\mathbf{w}_i)$  is a subgradient of  $R(\mathbf{w})$  at  $\mathbf{w}_i$ .

## Bundle Method for Risk Minimization

- ◆ The original convex problem

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} F(\mathbf{w}) := \left( \frac{\lambda}{2} \|\mathbf{w}\|^2 + R(\mathbf{w}) \right)$$

is reduced to a sequence of reduced convex problems

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} F_t(\mathbf{w}) := \left( \frac{\lambda}{2} \|\mathbf{w}\|^2 + R_t(\mathbf{w}) \right)$$

where  $R_t(\mathbf{w})$  is the “cutting plane model”

$$R_t(\mathbf{w}) = \max_{i=0, \dots, t} \left[ R(\mathbf{w}_i) + \langle \mathbf{g}_i, \mathbf{w} - \mathbf{w}_i \rangle \right]$$

and  $\mathbf{g}_i = \partial R(\mathbf{w}_i)$  is a subgradient of  $R(\mathbf{w})$  at  $\mathbf{w}_i$ .

- ◆ By construction it holds that  $R_t(\mathbf{w}) \leq R(\mathbf{w}), \forall \mathbf{w} \in \mathbb{R}^n$ .

## Subgradient

- ◆ Let  $f: \mathcal{X} \rightarrow \mathbb{R}$  be a convex function where  $\mathcal{X} \subseteq \mathbb{R}^n$  is a convex set.



## Subgradient

- ◆ Let  $f: \mathcal{X} \rightarrow \mathbb{R}$  be a convex function where  $\mathcal{X} \subseteq \mathbb{R}^n$  is a convex set.
- ◆ For differentiable  $f$  the gradient  $\nabla f(\mathbf{x}') \in \mathbb{R}^n$  at point  $\mathbf{x}' \in \mathcal{X}$  determines a global under-estimator of  $f$

$$f(\mathbf{x}) \geq f(\mathbf{x}') + \nabla f(\mathbf{x}')^T (\mathbf{x} - \mathbf{x}'), \quad \forall \mathbf{x} \in \mathcal{X}.$$

## Subgradient

- ◆ Let  $f: \mathcal{X} \rightarrow \mathbb{R}$  be a convex function where  $\mathcal{X} \subseteq \mathbb{R}^n$  is a convex set.
- ◆ For differentiable  $f$  the gradient  $\nabla f(\mathbf{x}') \in \mathbb{R}^n$  at point  $\mathbf{x}' \in \mathcal{X}$  determines a global under-estimator of  $f$

$$f(\mathbf{x}) \geq f(\mathbf{x}') + \nabla f(\mathbf{x}')^T (\mathbf{x} - \mathbf{x}'), \quad \forall \mathbf{x} \in \mathcal{X}.$$

- ◆ For non-differentiable  $f$  we can still construct a global under-estimator: the vector  $\mathbf{g} \in \mathbb{R}^n$  is a **subgradient** of  $f$  at point  $\mathbf{x}' \in \mathcal{X}$  if

$$f(\mathbf{x}) \geq f(\mathbf{x}') + \mathbf{g}^T (\mathbf{x} - \mathbf{x}'), \quad \forall \mathbf{x} \in \mathcal{X}.$$

## Subgradient

- ◆ Let  $f: \mathcal{X} \rightarrow \mathbb{R}$  be a convex function where  $\mathcal{X} \subseteq \mathbb{R}^n$  is a convex set.
- ◆ For differentiable  $f$  the gradient  $\nabla f(\mathbf{x}') \in \mathbb{R}^n$  at point  $\mathbf{x}' \in \mathcal{X}$  determines a global under-estimator of  $f$

$$f(\mathbf{x}) \geq f(\mathbf{x}') + \nabla f(\mathbf{x}')^T (\mathbf{x} - \mathbf{x}'), \quad \forall \mathbf{x} \in \mathcal{X}.$$

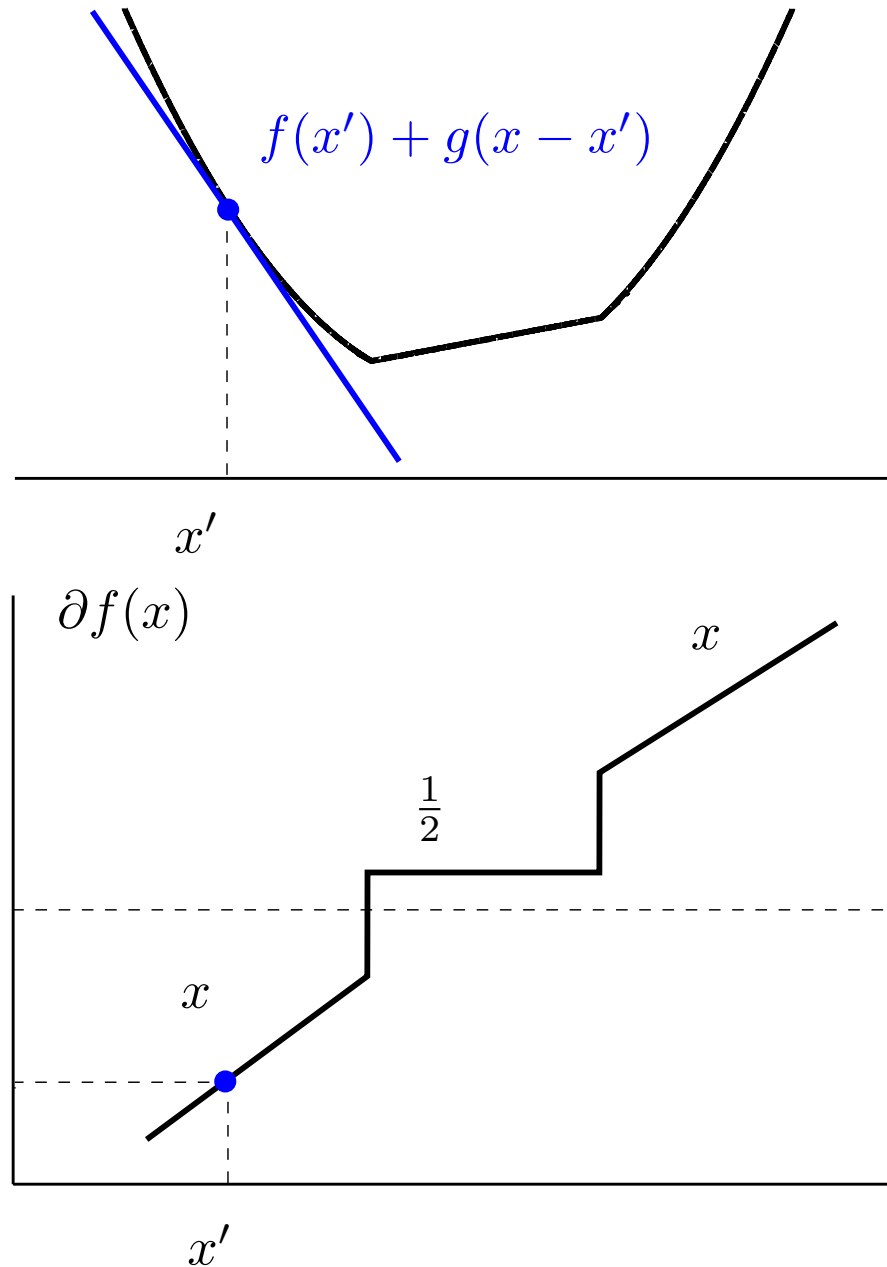
- ◆ For non-differentiable  $f$  we can still construct a global under-estimator: the vector  $\mathbf{g} \in \mathbb{R}^n$  is a subgradient of  $f$  at point  $\mathbf{x}' \in \mathcal{X}$  if

$$f(\mathbf{x}) \geq f(\mathbf{x}') + \mathbf{g}^T (\mathbf{x} - \mathbf{x}'), \quad \forall \mathbf{x} \in \mathcal{X}.$$

- ◆ There can be more than one subgradient at given point: the collection of subgradients of  $f$  at point  $\mathbf{x} \in \mathcal{X}$  is the subdifferential  $\partial f(\mathbf{x})$ .

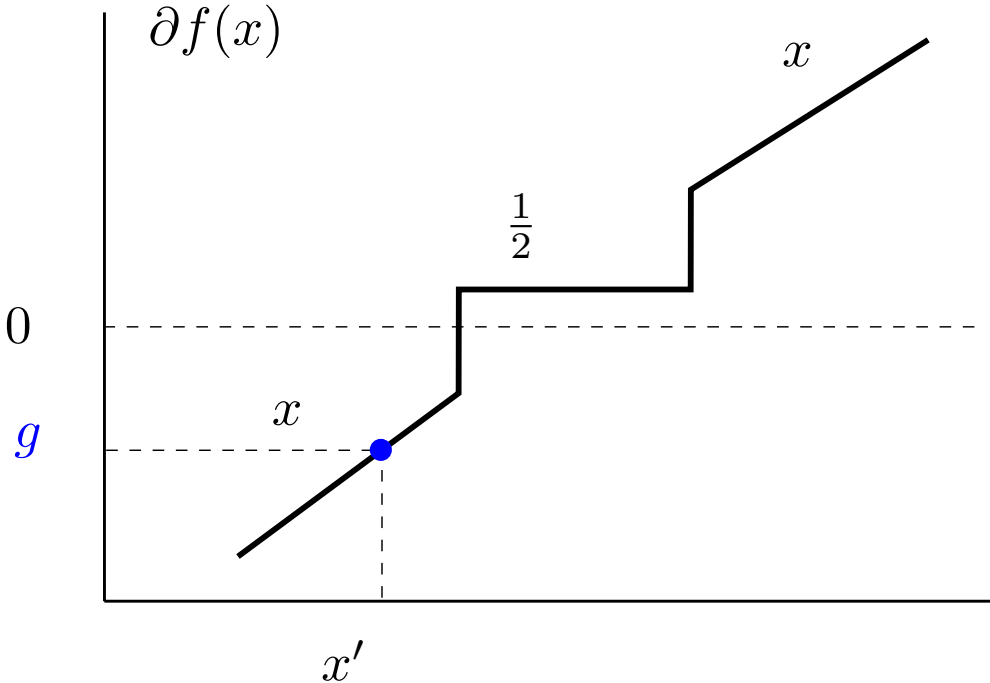
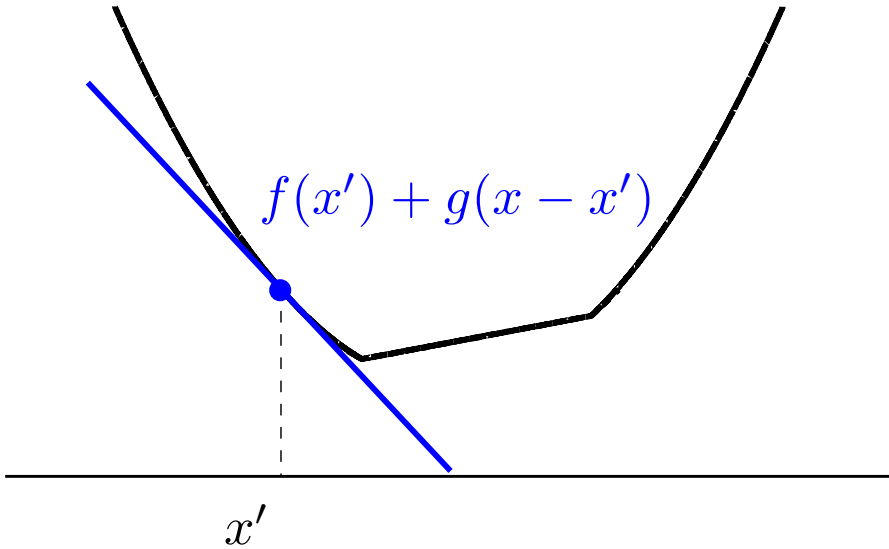
# An example: non-differentiable function and its subgradients

$$f(x) = \max\left\{\frac{1}{2}x^2, \frac{1}{2}x + 2\right\}$$



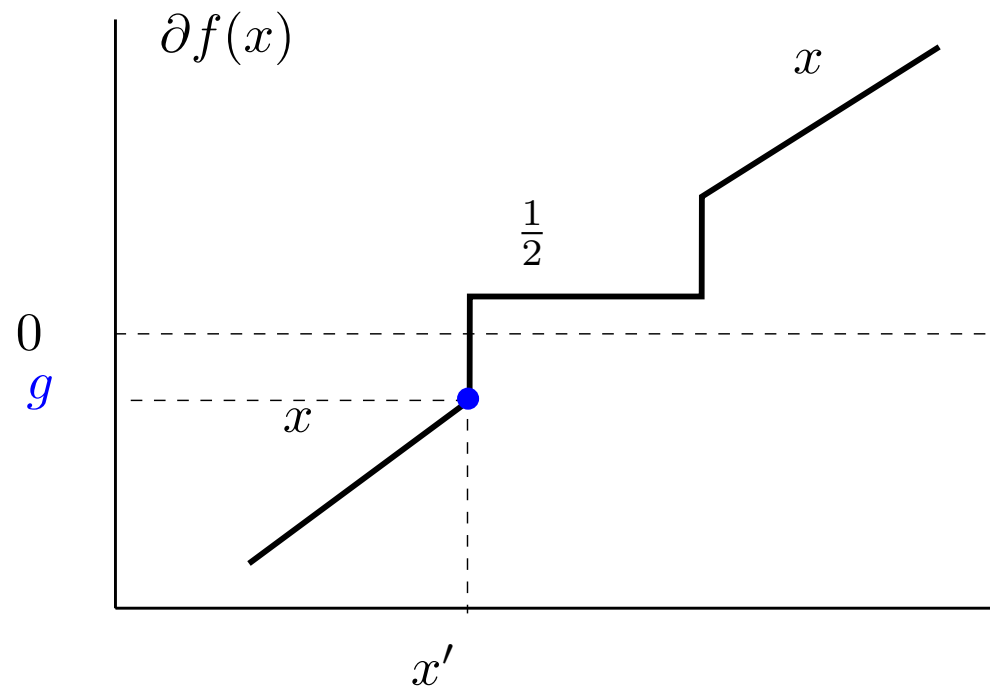
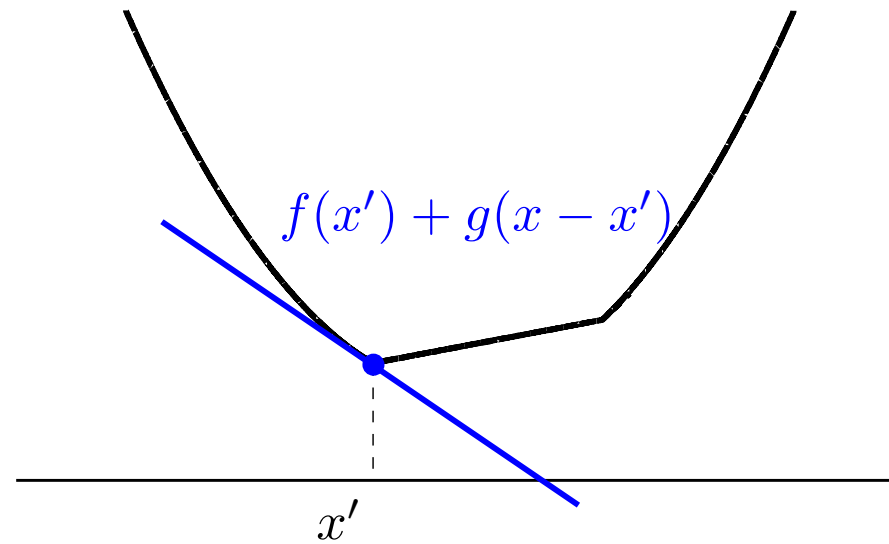
# An example: non-differentiable function and its subgradients

$$f(x) = \max\left\{\frac{1}{2}x^2, \frac{1}{2}x + 2\right\}$$



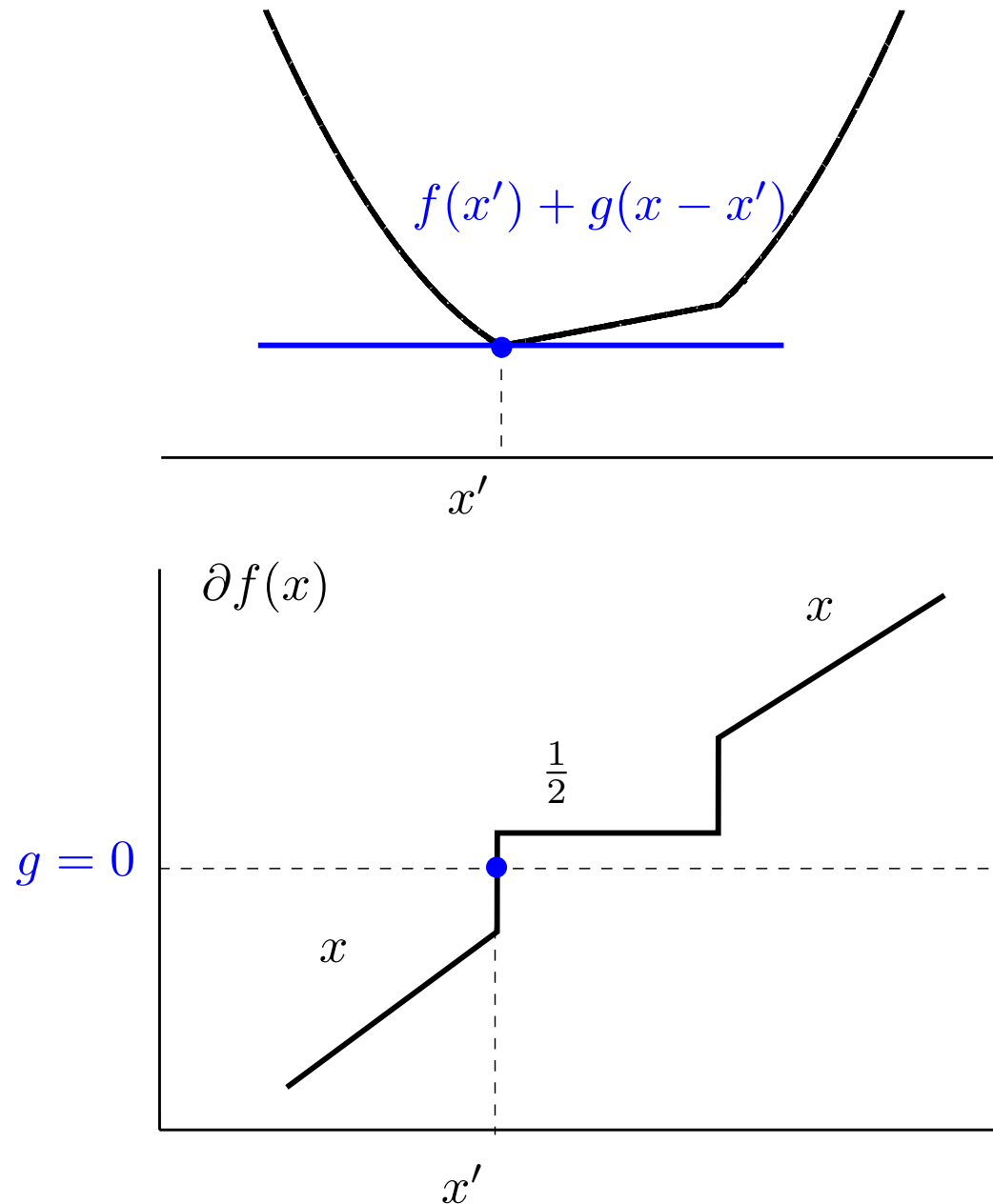
# An example: non-differentiable function and its subgradients

$$f(x) = \max\left\{\frac{1}{2}x^2, \frac{1}{2}x + 2\right\}$$



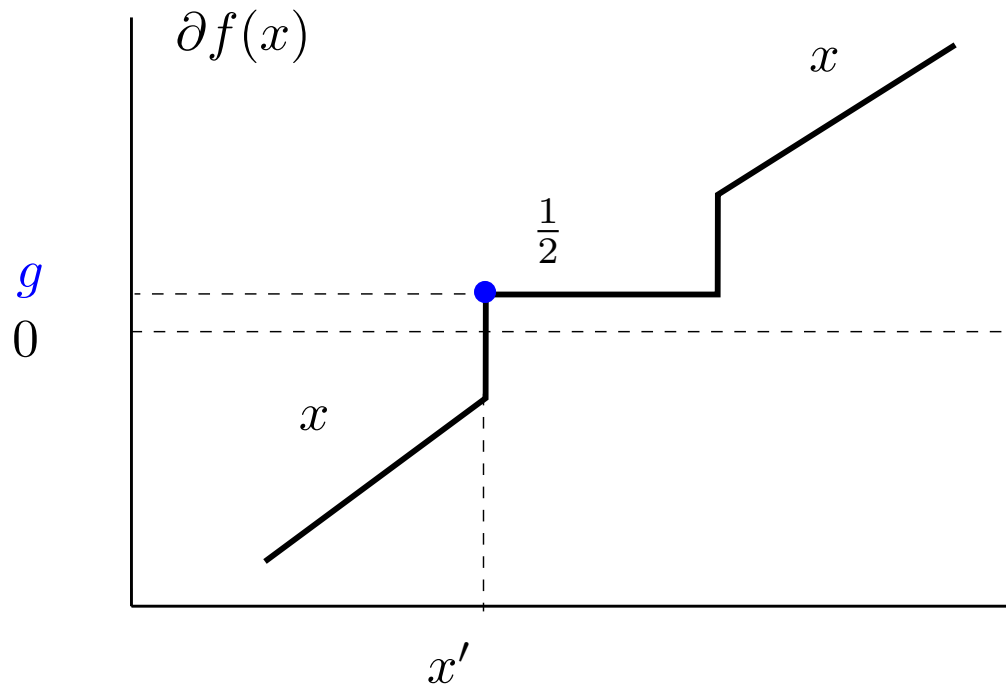
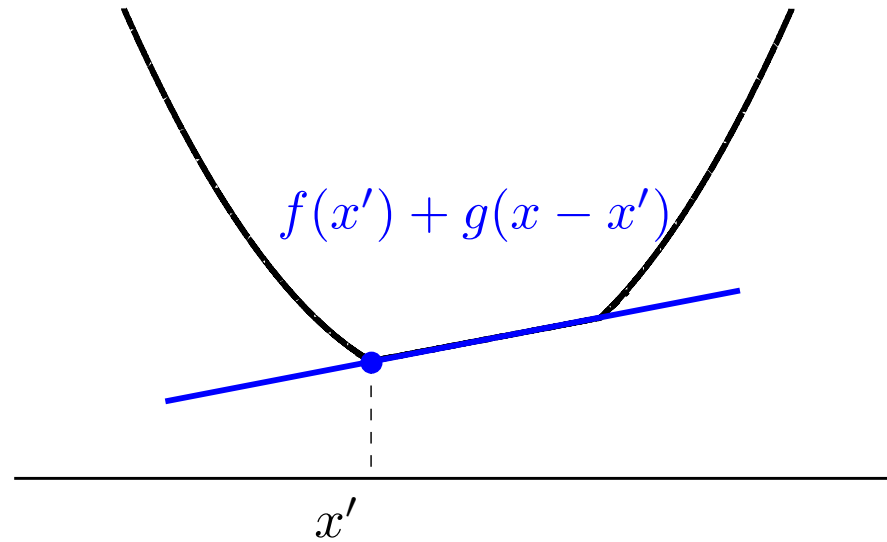
# An example: non-differentiable function and its subgradients

$$f(x) = \max\left\{\frac{1}{2}x^2, \frac{1}{2}x + 2\right\}$$



# An example: non-differentiable function and its subgradients

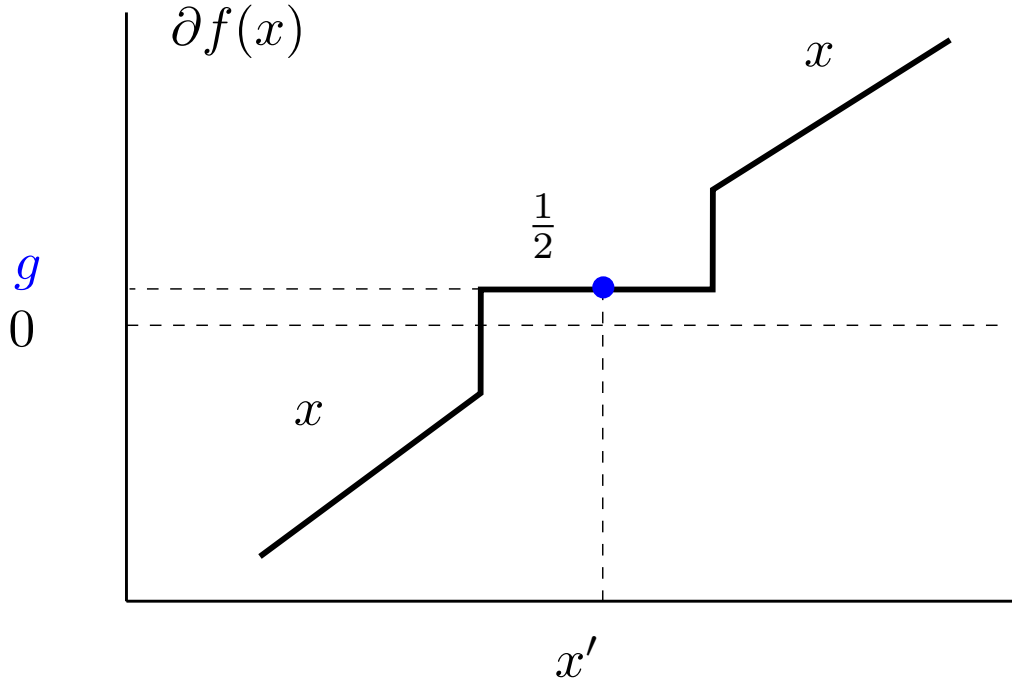
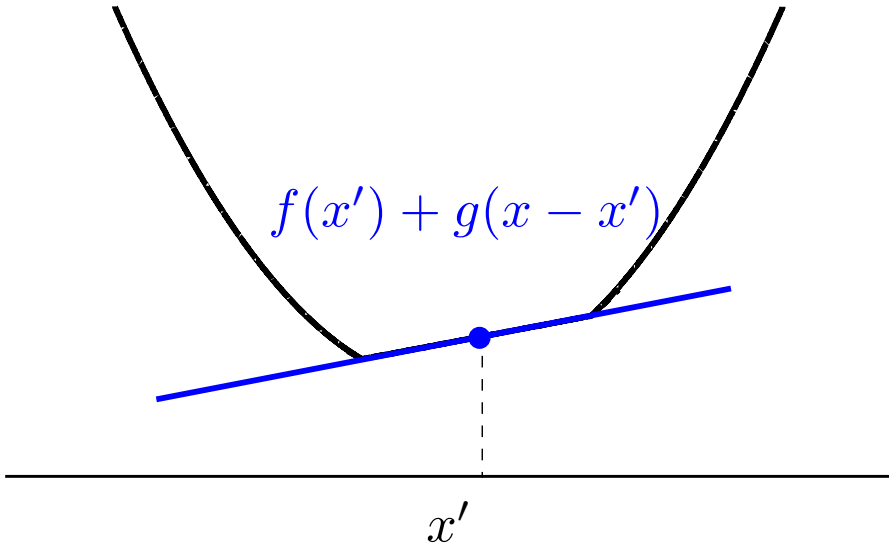
$$f(x) = \max\left\{\frac{1}{2}x^2, \frac{1}{2}x + 2\right\}$$





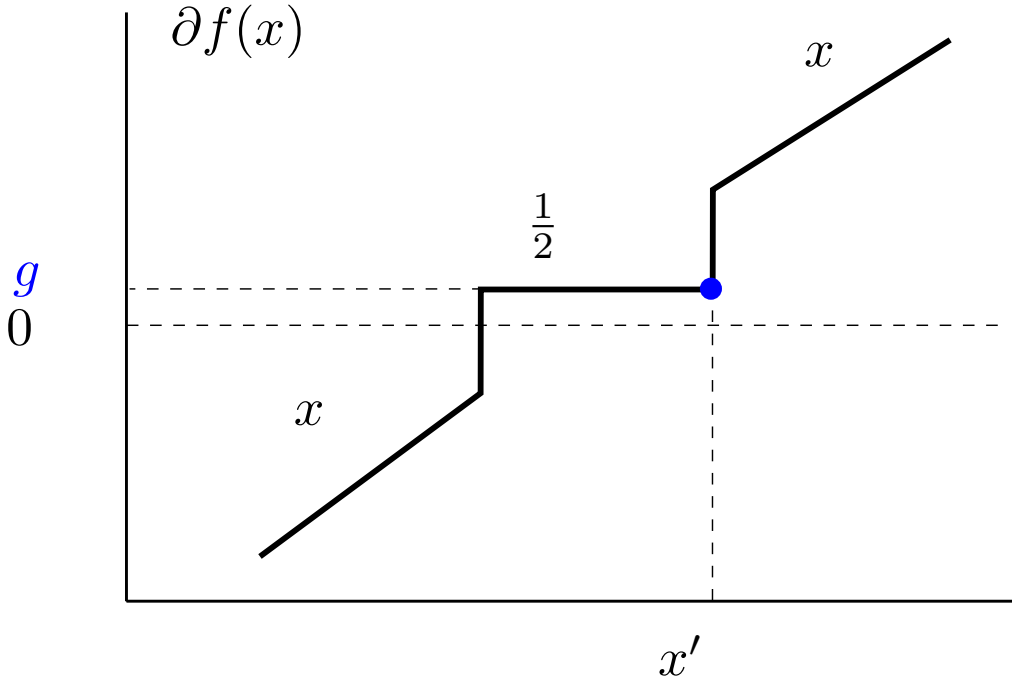
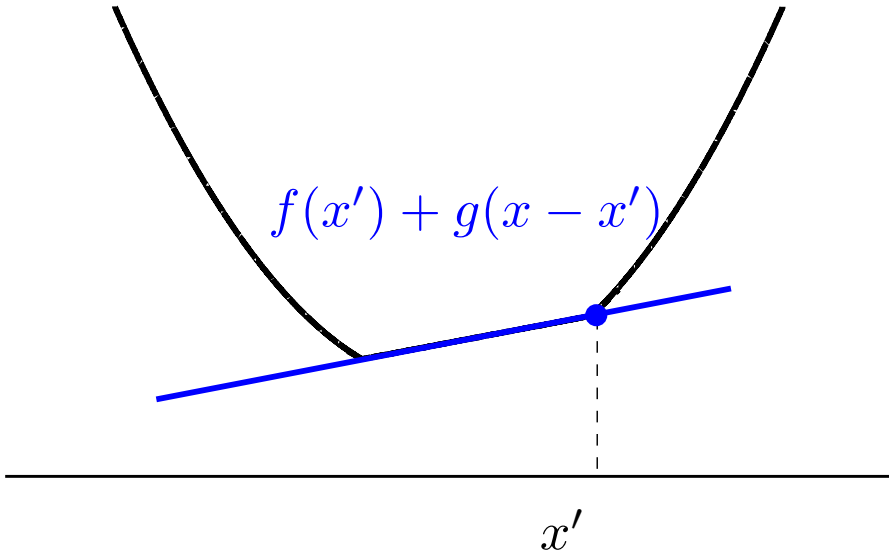
# An example: non-differentiable function and its subgradients

$$f(x) = \max\{\frac{1}{2}x^2, \frac{1}{2}x + 2\}$$



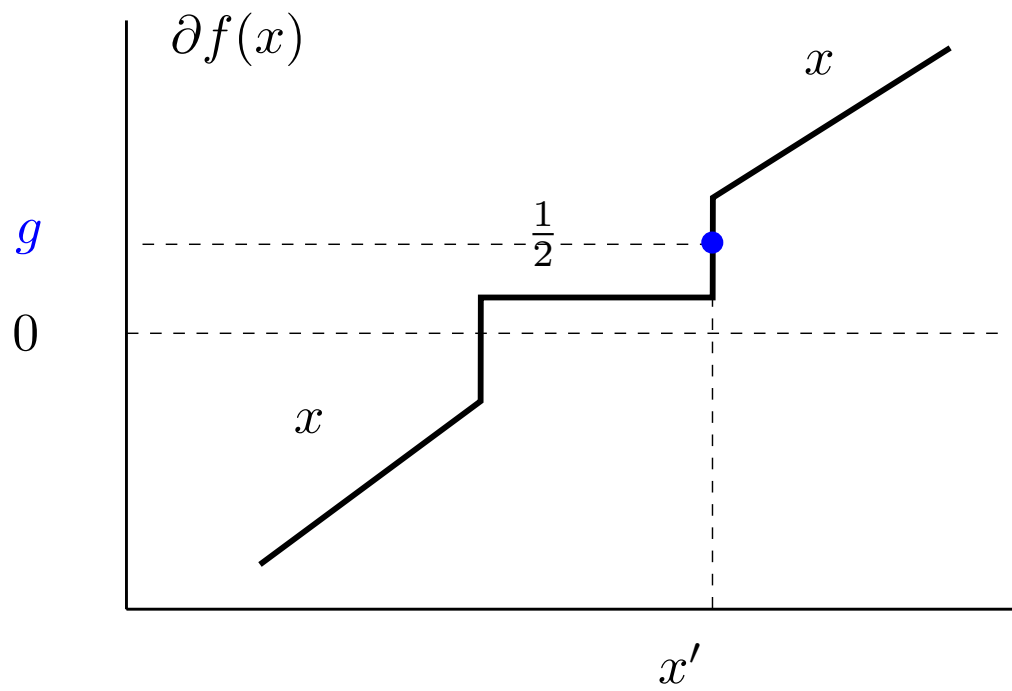
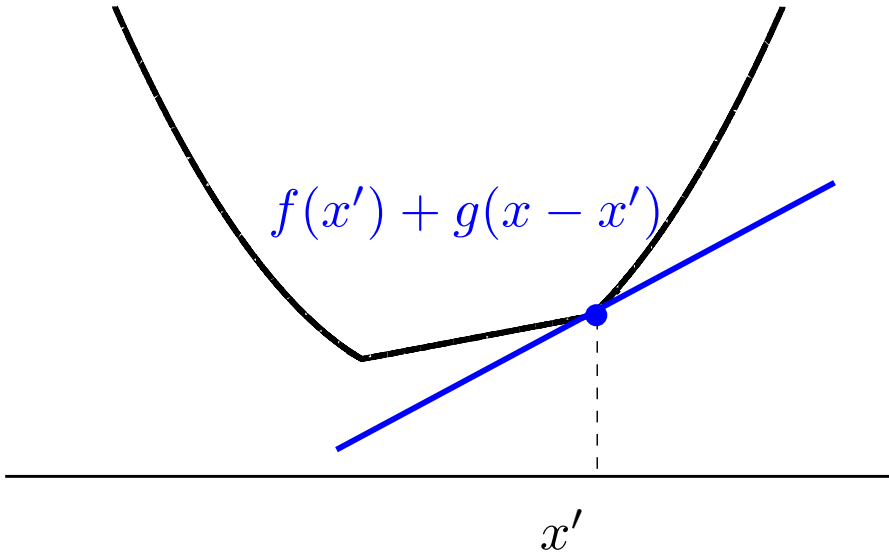
# An example: non-differentiable function and its subgradients

$$f(x) = \max\left\{\frac{1}{2}x^2, \frac{1}{2}x + 2\right\}$$



# An example: non-differentiable function and its subgradients

$$f(x) = \max\left\{\frac{1}{2}x^2, \frac{1}{2}x + 2\right\}$$



## Basic subgradient calculus



- ◆  $\partial f(x)$  is a closed convex set.

## Basic subgradient calculus

- ◆  $\partial f(\mathbf{x})$  is a closed convex set.
- ◆  $\partial f(\mathbf{x}) = \{\mathbf{g}\} \iff f$  is differentiable and  $\mathbf{g} = \nabla f(\mathbf{x})$ .

## Basic subgradient calculus

- ◆  $\partial f(\mathbf{x})$  is a closed convex set.
- ◆  $\partial f(\mathbf{x}) = \{\mathbf{g}\} \iff f$  is differentiable and  $\mathbf{g} = \nabla f(\mathbf{x})$ .
- ◆ **scaling:**  $\partial(\alpha f(\mathbf{x})) = \alpha \partial f(\mathbf{x})$  if  $\alpha > 0$ .

## Basic subgradient calculus

- ◆  $\partial f(\mathbf{x})$  is a closed convex set.
- ◆  $\partial f(\mathbf{x}) = \{\mathbf{g}\} \iff f$  is differentiable and  $\mathbf{g} = \nabla f(\mathbf{x})$ .
- ◆ **scaling:**  $\partial(\alpha f(\mathbf{x})) = \alpha \partial f(\mathbf{x})$  if  $\alpha > 0$ .
- ◆ **addition:**  $\partial(f_1(\mathbf{x}) + f_2(\mathbf{x})) = \partial f_1(\mathbf{x}) + \partial f_2(\mathbf{x})$ .

## Basic subgradient calculus

- ◆  $\partial f(\mathbf{x})$  is a closed convex set.
- ◆  $\partial f(\mathbf{x}) = \{\mathbf{g}\} \iff f$  is differentiable and  $\mathbf{g} = \nabla f(\mathbf{x})$ .
- ◆ **scaling:**  $\partial(\alpha f(\mathbf{x})) = \alpha \partial f(\mathbf{x})$  if  $\alpha > 0$ .
- ◆ **addition:**  $\partial(f_1(\mathbf{x}) + f_2(\mathbf{x})) = \partial f_1(\mathbf{x}) + \partial f_2(\mathbf{x})$ .
- ◆ **point-wise maximum:**  $f(\mathbf{x}) = \max_{i=1, \dots, m} f_i(\mathbf{x})$  when  $f_i(\mathbf{x})$  are differentiable then

$$\partial f(\mathbf{x}) = \mathbf{Co}\{\nabla f_i(\mathbf{x}) \mid f_i(\mathbf{x}) = f(\mathbf{x})\},$$

i.e., convex hull of gradients of active functions at  $\mathbf{x}$ .



## Basic subgradient calculus

- ◆  $\partial f(\mathbf{x})$  is a closed convex set.
- ◆  $\partial f(\mathbf{x}) = \{\mathbf{g}\} \iff f$  is differentiable and  $\mathbf{g} = \nabla f(\mathbf{x})$ .
- ◆ **scaling:**  $\partial(\alpha f(\mathbf{x})) = \alpha \partial f(\mathbf{x})$  if  $\alpha > 0$ .
- ◆ **addition:**  $\partial(f_1(\mathbf{x}) + f_2(\mathbf{x})) = \partial f_1(\mathbf{x}) + \partial f_2(\mathbf{x})$ .
- ◆ **point-wise maximum:**  $f(\mathbf{x}) = \max_{i=1,\dots,m} f_i(\mathbf{x})$  when  $f_i(\mathbf{x})$  are differentiable then

$$\partial f(\mathbf{x}) = \mathbf{Co}\{\nabla f_i(\mathbf{x}) \mid f_i(\mathbf{x}) = f(\mathbf{x})\},$$

i.e., convex hull of gradients of active functions at  $\mathbf{x}$ .

- ◆ **Optimality condition** for a convex  $f$ :

$$f(\mathbf{x}^*) = \inf_{\mathbf{x}} f(\mathbf{x}) \iff \mathbf{0} \in \partial f(\mathbf{x}^*)$$

## Example: cutting plane model for SVM

- ◆ The convex problem to solve

$$\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \frac{\lambda}{2} \|\mathbf{w}\|^2 + \underbrace{\frac{1}{m} \sum_{i=1}^m \max \{0, 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle\}}_{R(\mathbf{w})} \right)$$

## Example: cutting plane model for SVM

- ◆ The convex problem to solve

$$\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \max \{0, 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle\}}_{R(\mathbf{w})} \right)$$

- ◆ The cutting plane model of  $R(\mathbf{w})$  reads

$$R_t(\mathbf{w}) = \max_{i=0, \dots, t-1} \left[ R(\mathbf{w}_i) + \langle \mathbf{g}_i, \mathbf{w} - \mathbf{w}_i \rangle \right]$$

- ◆ The subgradient of  $R(\mathbf{w})$  at  $\mathbf{w}$

$$\mathbf{g}_i = -\frac{1}{m} \sum_{i=1}^m y^i \mathbf{x}^i \llbracket \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 1 \rrbracket$$

# Example: cutting plane model for SO-SVM with margin-rescaling loss



- ◆ The convex problem to solve

$$\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \frac{\lambda}{2} \|\mathbf{w}\|^2 + \underbrace{\frac{1}{m} \sum_{i=1}^m \max_{y \in \mathcal{Y}} \{ \ell_i(y) + \langle \mathbf{w}, \phi_i(y) \rangle \}}_{R(\mathbf{w})} \right)$$

# Example: cutting plane model for SO-SVM with margin-rescaling loss



- ◆ The convex problem to solve

$$\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \frac{\lambda}{2} \|\mathbf{w}\|^2 + \underbrace{\frac{1}{m} \sum_{i=1}^m \max_{y \in \mathcal{Y}} \{ \ell_i(y) + \langle \mathbf{w}, \phi_i(y) \rangle \}}_{R(\mathbf{w})} \right)$$

- ◆ The cutting plane model of  $R(\mathbf{w})$  reads

$$R_t(\mathbf{w}) = \max_{i=0, \dots, t-1} \left[ R(\mathbf{w}_i) + \langle \mathbf{g}_i, \mathbf{w} - \mathbf{w}_i \rangle \right]$$

- ◆ The subgradient of  $R(\mathbf{w})$  at  $\mathbf{w}$

$$\mathbf{g}_i = \frac{1}{m} \sum_{i=1}^m \phi_i(\hat{y}^i) \quad \text{where} \quad \hat{y}^i = \operatorname{argmax}_{y \in \mathcal{Y}} \{ \ell_i(y) + \langle \mathbf{w}_i, \phi_i(y) \rangle \}$$

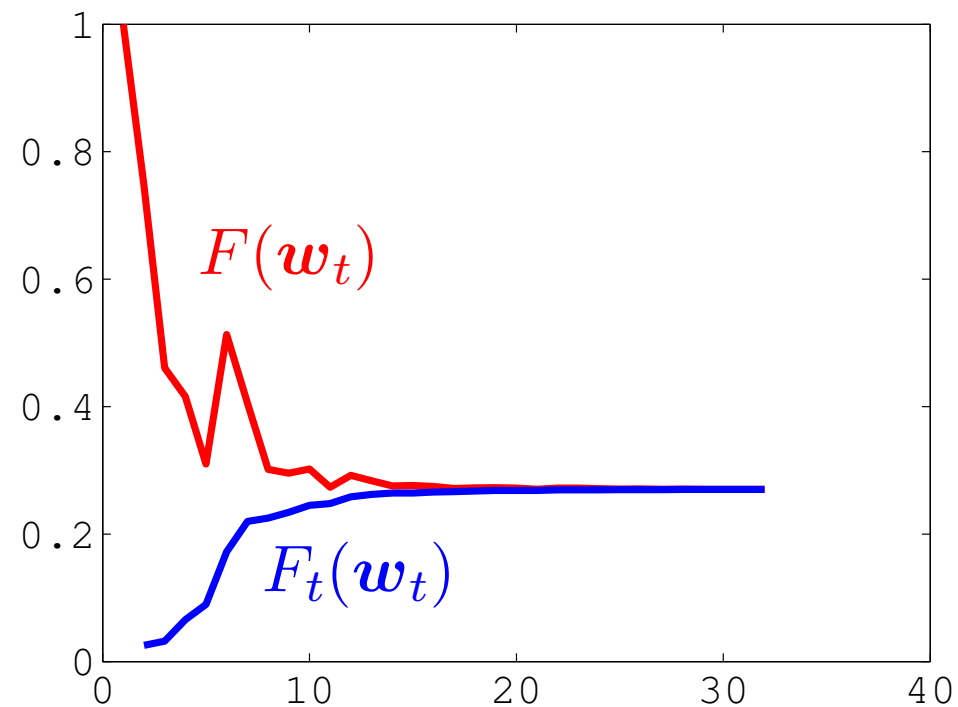
## Bundle Method for Risk Minimization

1. Init:  $t \leftarrow 0, \mathbf{w}_0 \in \mathbb{R}^n$
2. Compute  $R(\mathbf{w}_t)$  and  $\mathbf{g}_t \in \partial R(\mathbf{w}_t)$
3.  $\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \frac{\lambda}{2} \|\mathbf{w}\|^2 + R_t(\mathbf{w}) \right)$

where

$$R_t(\mathbf{w}) = \max_{i=0, \dots, t} \left[ R(\mathbf{w}_i) + \langle \mathbf{g}_i, \mathbf{w} - \mathbf{w}_i \rangle \right]$$

4. if  $\min_{i=1, \dots, t} F(\mathbf{w}_i) - F_t(\mathbf{w}_{t+1}) \leq \varepsilon$  stop  
 else  $t \leftarrow t + 1$  go to 2.



## How to solve the reduced problem

- ◆ Let us define a matrix  $\mathbf{A} = [\mathbf{g}_0, \dots, \mathbf{g}_t] \in \mathbb{R}^{n \times t}$  and a vector  $\mathbf{b} = [b_0, \dots, b_{t-1}]$  with components  $b_i = R(\mathbf{w}_i) - \langle \mathbf{g}_i, \mathbf{w}_i \rangle$ .
- ◆ The **reduced problem** can be expressed as

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n, \xi \in \mathbb{R}} \left[ \frac{\lambda}{2} \|\mathbf{w}\|^2 + \xi \right] \quad \text{s.t.} \quad \xi \geq \langle \mathbf{w}, \mathbf{g}_i \rangle + b_i, i \in \{0, \dots, t\}$$

## How to solve the reduced problem

- ◆ Let us define a matrix  $\mathbf{A} = [\mathbf{g}_0, \dots, \mathbf{g}_t] \in \mathbb{R}^{n \times t}$  and a vector  $\mathbf{b} = [b_0, \dots, b_{t-1}]$  with components  $b_i = R(\mathbf{w}_i) - \langle \mathbf{g}_i, \mathbf{w}_i \rangle$ .
- ◆ The **reduced problem** can be expressed as

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n, \xi \in \mathbb{R}} \left[ \frac{\lambda}{2} \|\mathbf{w}\|^2 + \xi \right] \quad \text{s.t.} \quad \xi \geq \langle \mathbf{w}, \mathbf{g}_i \rangle + b_i, i \in \{0, \dots, t\}$$

- ◆ The Lagrange **dual of the reduced problem** reads

$$\boldsymbol{\alpha}_{t+1} = \operatorname{argmax}_{\boldsymbol{\alpha} \in \mathbb{R}^t} \left[ \langle \boldsymbol{\alpha}, \mathbf{b} \rangle - \frac{1}{2\lambda} \langle \boldsymbol{\alpha}, \mathbf{A}^T \mathbf{A} \boldsymbol{\alpha} \rangle \right] \quad \text{s.t.} \quad \|\boldsymbol{\alpha}\|_1 = 1, \boldsymbol{\alpha} \geq \mathbf{0}$$



## How to solve the reduced problem

- ◆ Let us define a matrix  $\mathbf{A} = [\mathbf{g}_0, \dots, \mathbf{g}_t] \in \mathbb{R}^{n \times t}$  and a vector  $\mathbf{b} = [b_0, \dots, b_{t-1}]$  with components  $b_i = R(\mathbf{w}_i) - \langle \mathbf{g}_i, \mathbf{w}_i \rangle$ .
- ◆ The **reduced problem** can be expressed as

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n, \xi \in \mathbb{R}} \left[ \frac{\lambda}{2} \|\mathbf{w}\|^2 + \xi \right] \quad \text{s.t.} \quad \xi \geq \langle \mathbf{w}, \mathbf{g}_i \rangle + b_i, i \in \{0, \dots, t\}$$

- ◆ The Lagrange **dual of the reduced problem** reads

$$\boldsymbol{\alpha}_{t+1} = \operatorname{argmax}_{\boldsymbol{\alpha} \in \mathbb{R}^t} \left[ \langle \boldsymbol{\alpha}, \mathbf{b} \rangle - \frac{1}{2\lambda} \langle \boldsymbol{\alpha}, \mathbf{A}^T \mathbf{A} \boldsymbol{\alpha} \rangle \right] \quad \text{s.t.} \quad \|\boldsymbol{\alpha}\|_1 = 1, \boldsymbol{\alpha} \geq \mathbf{0}$$

- ◆ The primal solution is  $\mathbf{w}_{t+1} = -\lambda^{-1} \mathbf{A} \boldsymbol{\alpha}_{t+1}$

## Bundle Method for Risk Minimization

(+) Provides a certificate of optimality:

$$F(\mathbf{w}_t) \leq F(\mathbf{w}^*) + \varepsilon$$

## Bundle Method for Risk Minimization

(+) Provides a certificate of optimality:

$$F(\mathbf{w}_t) \leq F(\mathbf{w}^*) + \varepsilon$$

(+) Converges for arbitrary  $\varepsilon > 0$  in

$$\log_2 \frac{\lambda F(\mathbf{0})}{G^2} + \frac{4G^2}{\lambda \varepsilon} - 1$$

iteration at most where  $G \geq \|\mathbf{g}\|_2, \forall \mathbf{g} \in \partial R(\mathbf{w}), \mathbf{w} \in \mathbb{R}^n$ .

## Bundle Method for Risk Minimization

(+) Provides a certificate of optimality:

$$F(\mathbf{w}_t) \leq F(\mathbf{w}^*) + \varepsilon$$

(+) Converges for arbitrary  $\varepsilon > 0$  in

$$\log_2 \frac{\lambda F(\mathbf{0})}{G^2} + \frac{4G^2}{\lambda \varepsilon} - 1$$

iteration at most where  $G \geq \|\mathbf{g}\|_2, \forall \mathbf{g} \in \partial R(\mathbf{w}), \mathbf{w} \in \mathbb{R}^n$ .

(+) Requires only the first order oracle computing  $R(\mathbf{w})$  and  $\mathbf{g} \in \partial R(\mathbf{w})$ .

## Bundle Method for Risk Minimization

(+) Provides a certificate of optimality:

$$F(\mathbf{w}_t) \leq F(\mathbf{w}^*) + \varepsilon$$

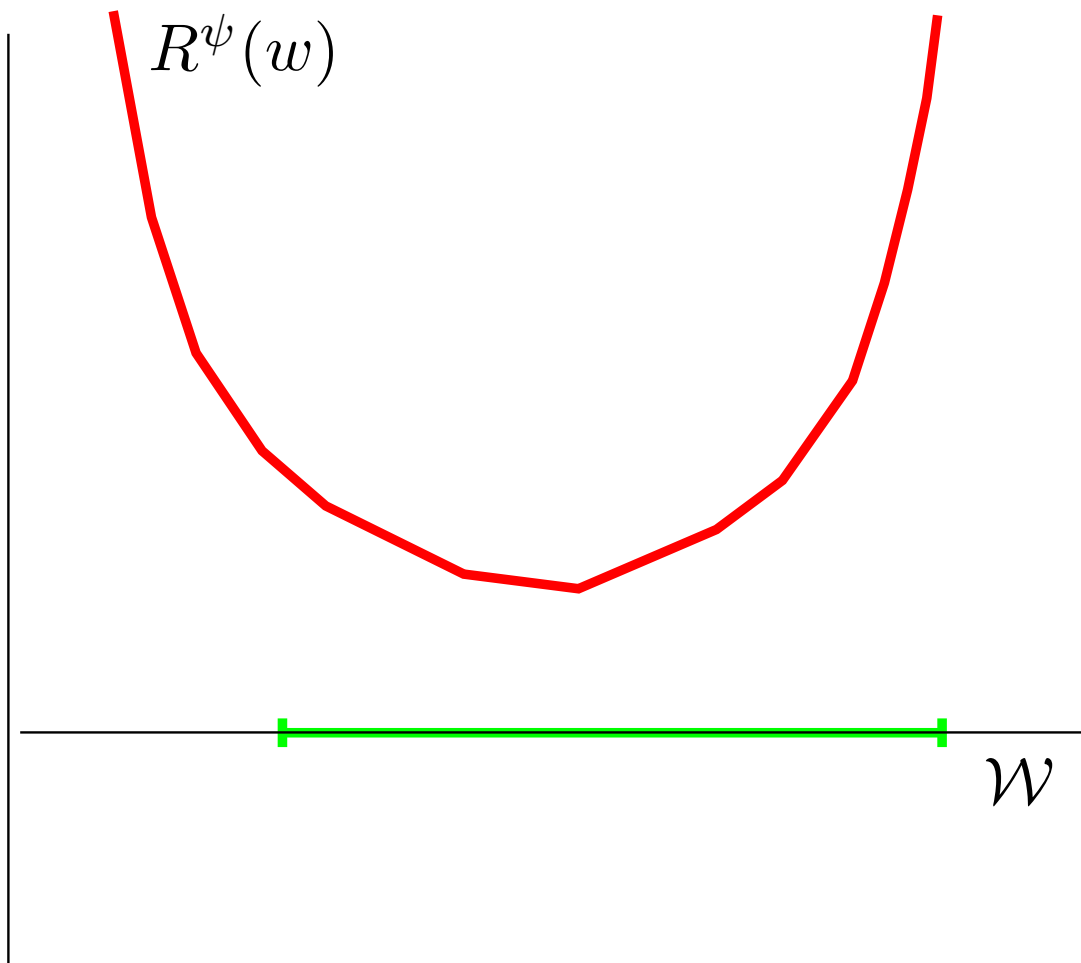
(+) Converges for arbitrary  $\varepsilon > 0$  in

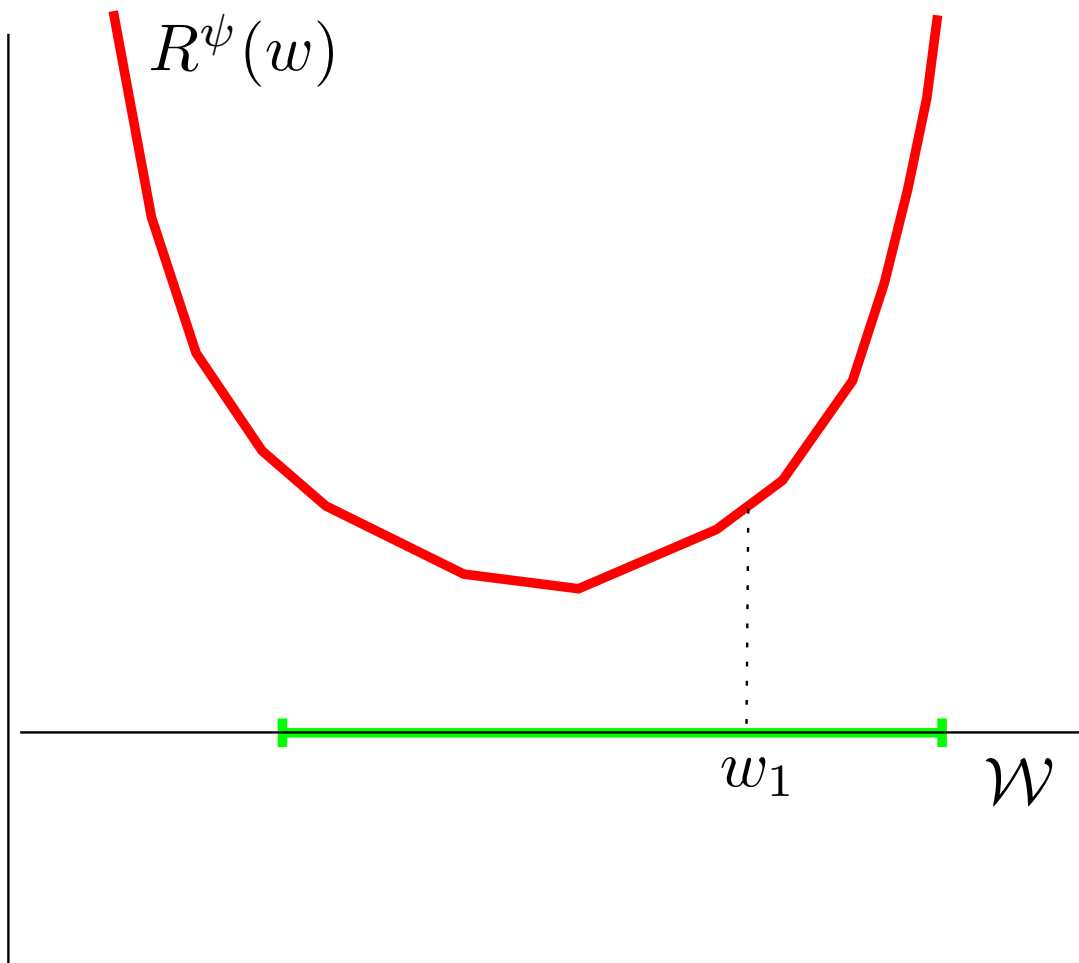
$$\log_2 \frac{\lambda F(\mathbf{0})}{G^2} + \frac{4G^2}{\lambda \varepsilon} - 1$$

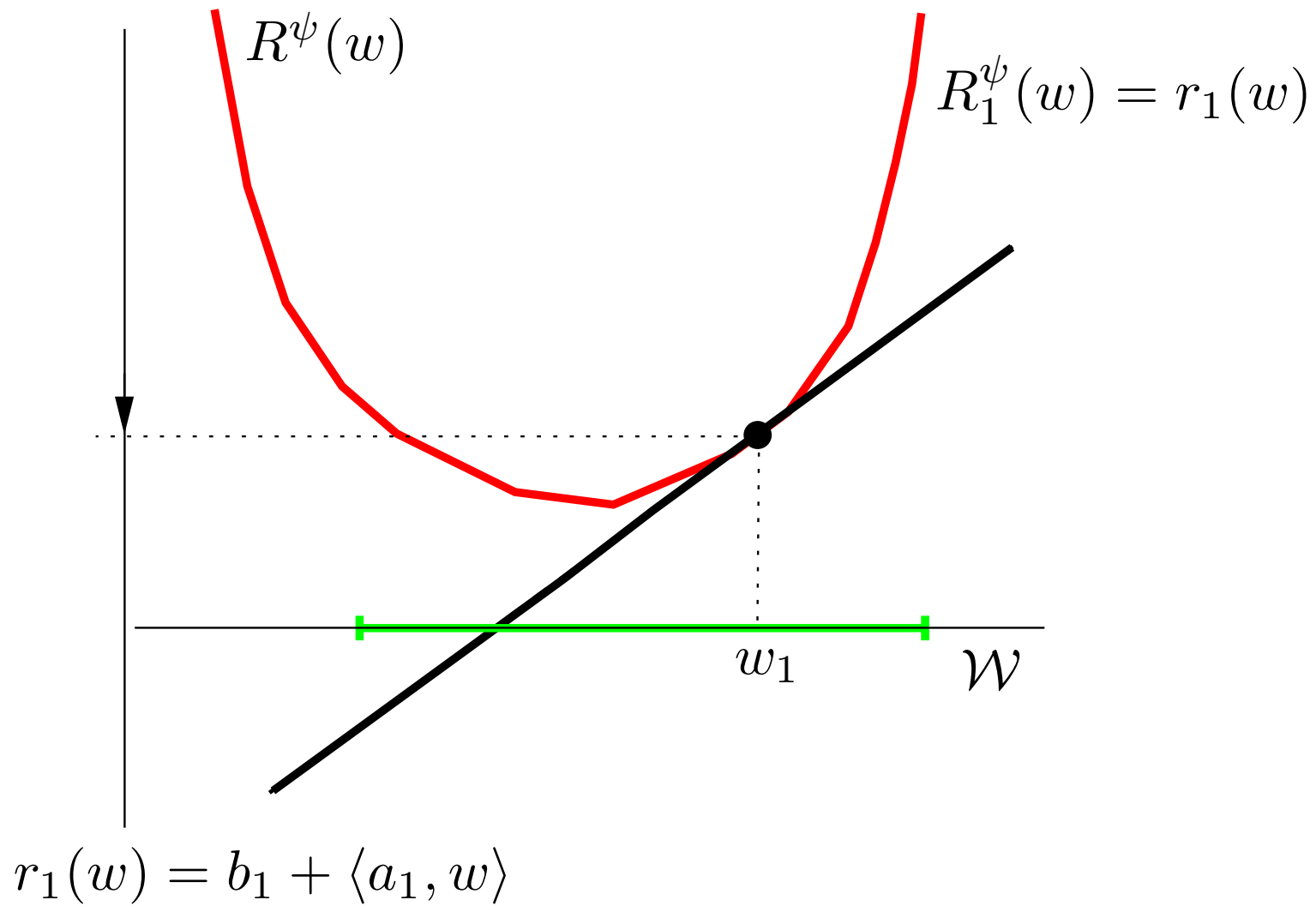
iteration at most where  $G \geq \|\mathbf{g}\|_2, \forall \mathbf{g} \in \partial R(\mathbf{w}), \mathbf{w} \in \mathbb{R}^n$ .

(+) Requires only the first order oracle computing  $R(\mathbf{w})$  and  $\mathbf{g} \in \partial R(\mathbf{w})$ .

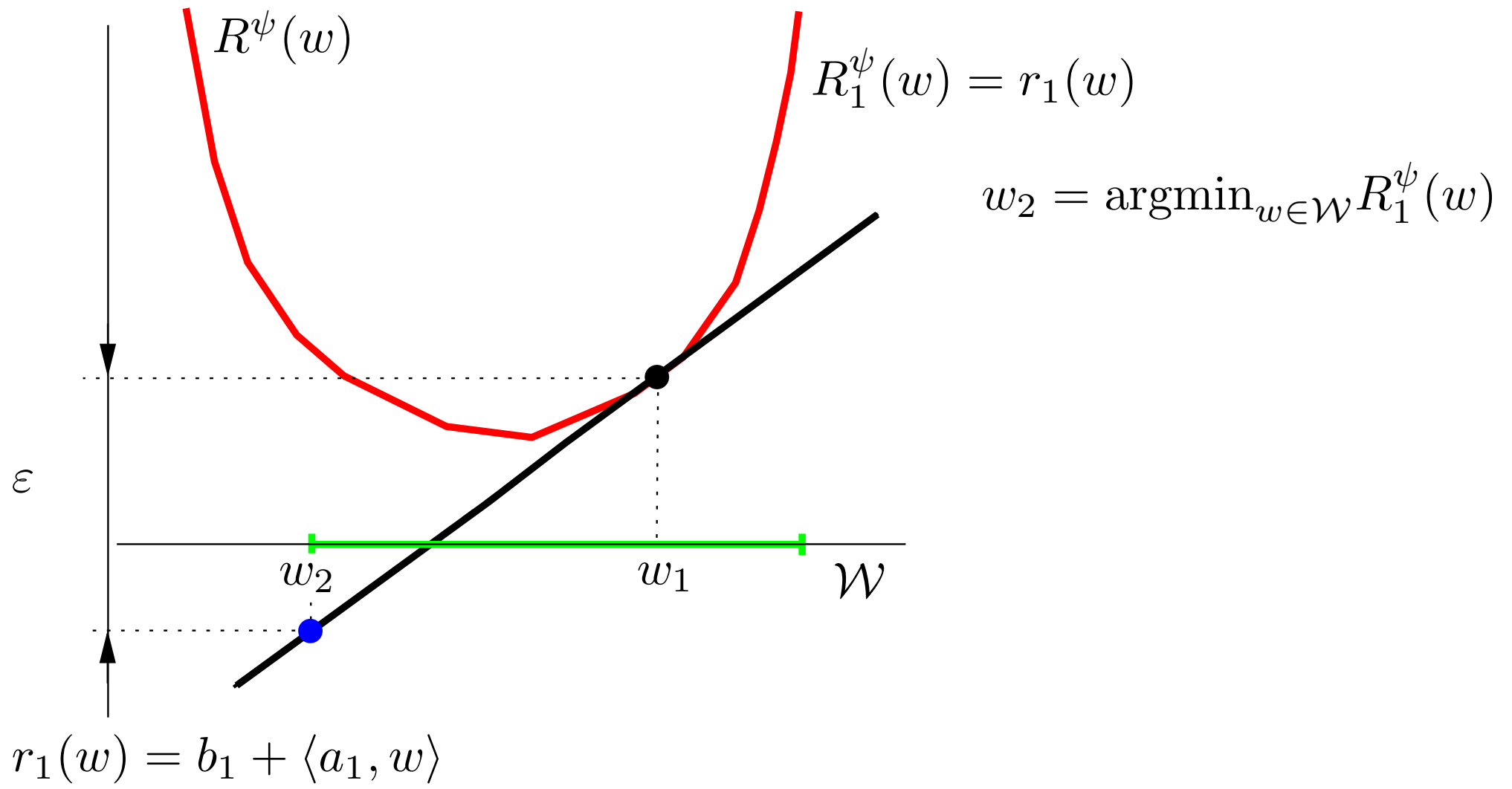
(-) Slow convergence for  $\lambda \rightarrow 0$ .

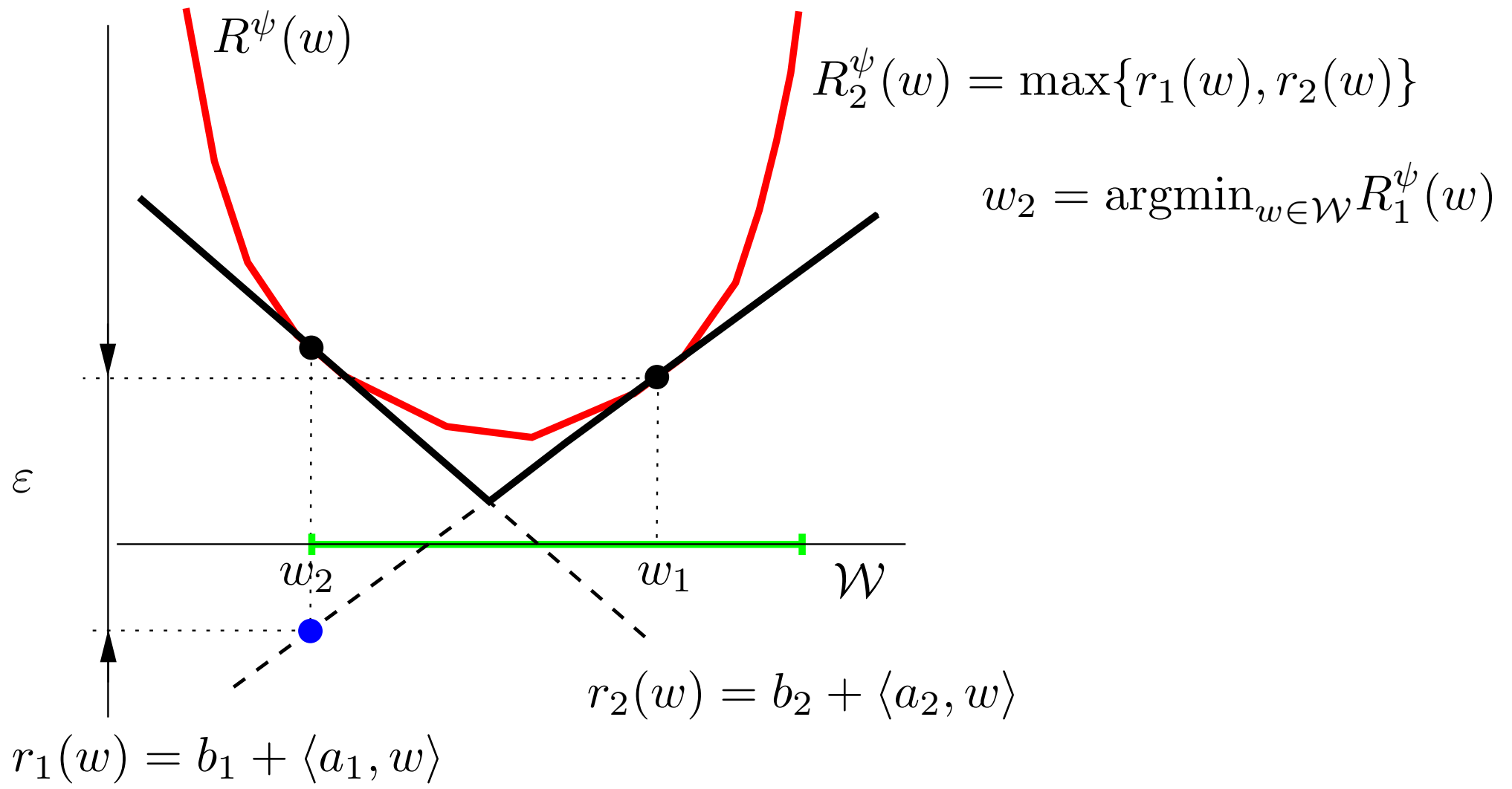


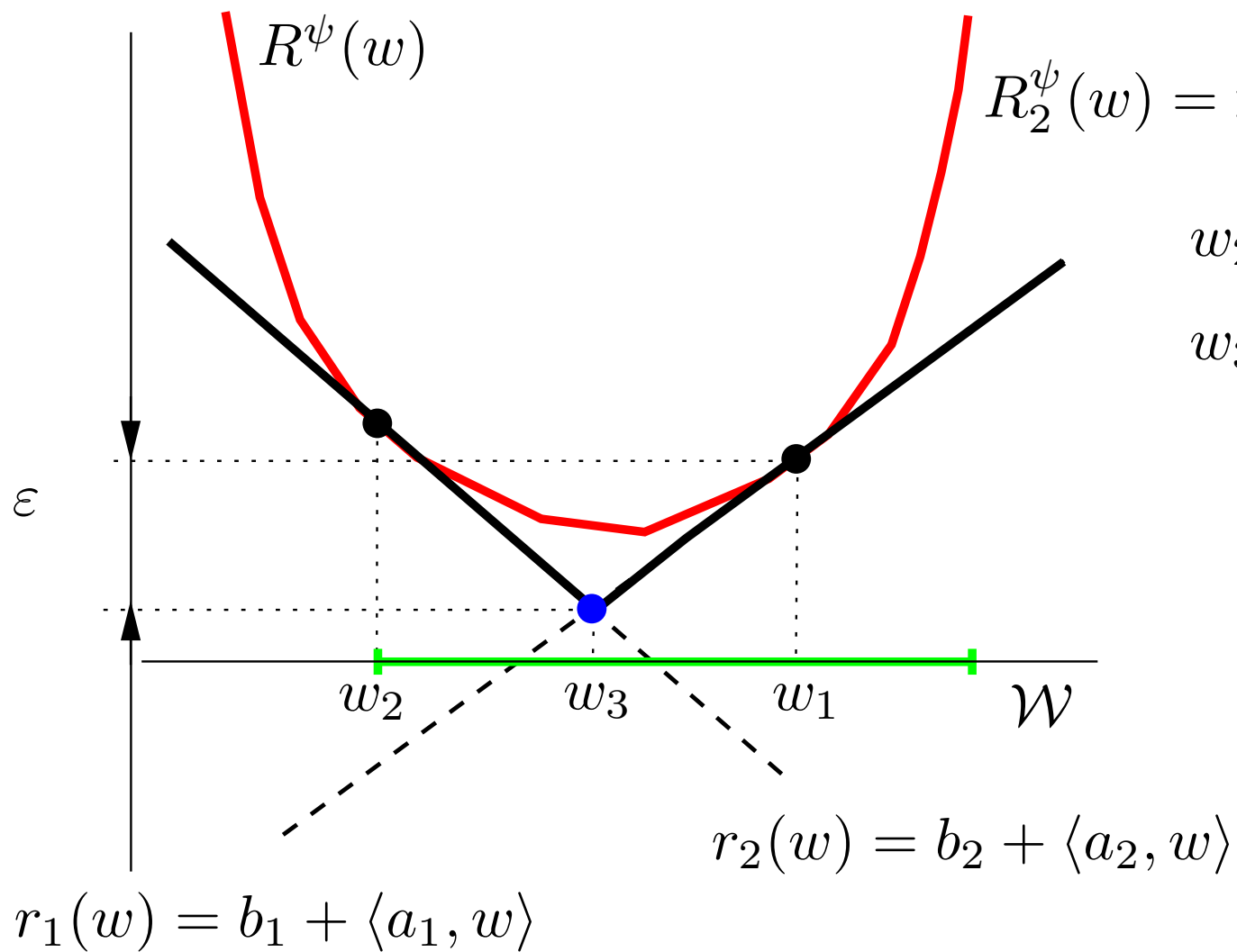


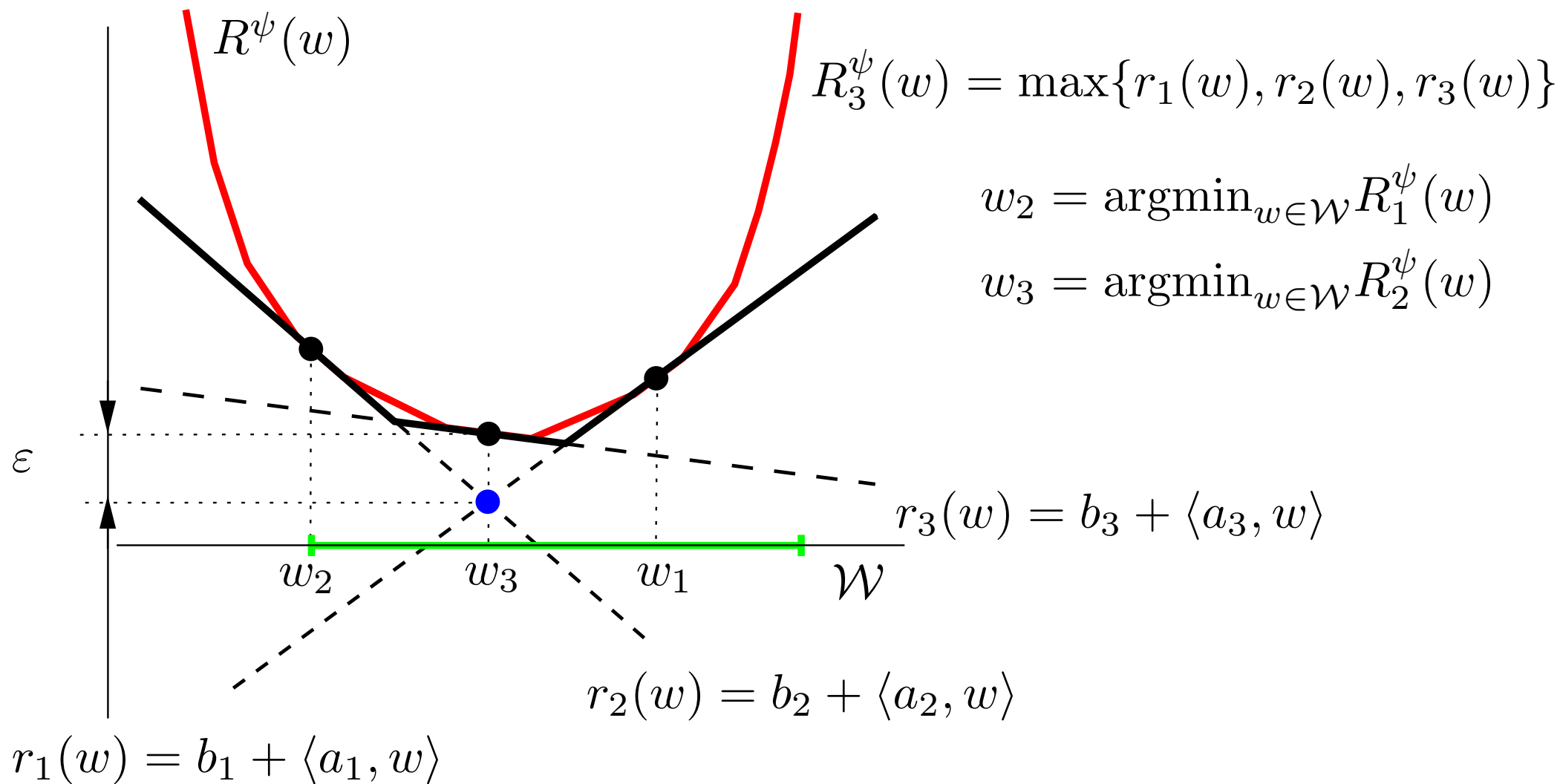


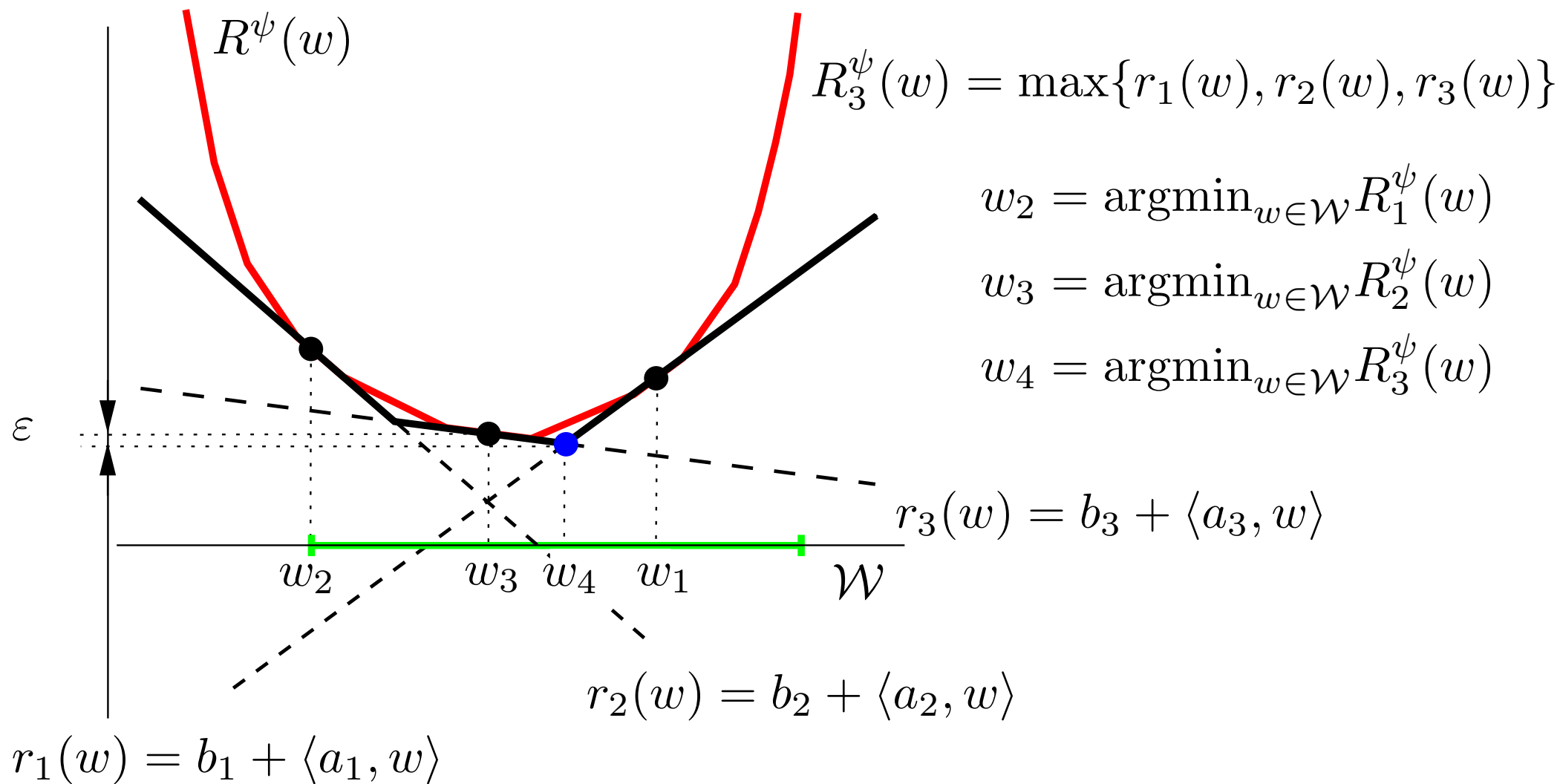




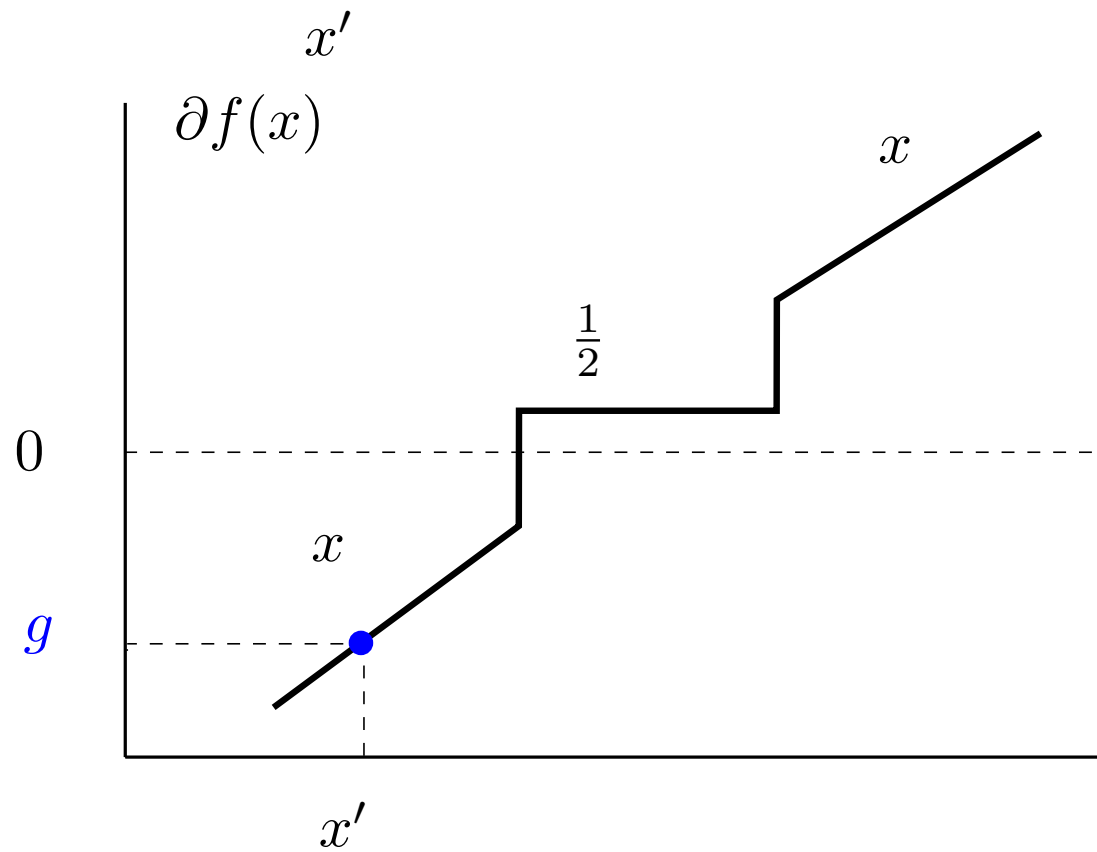
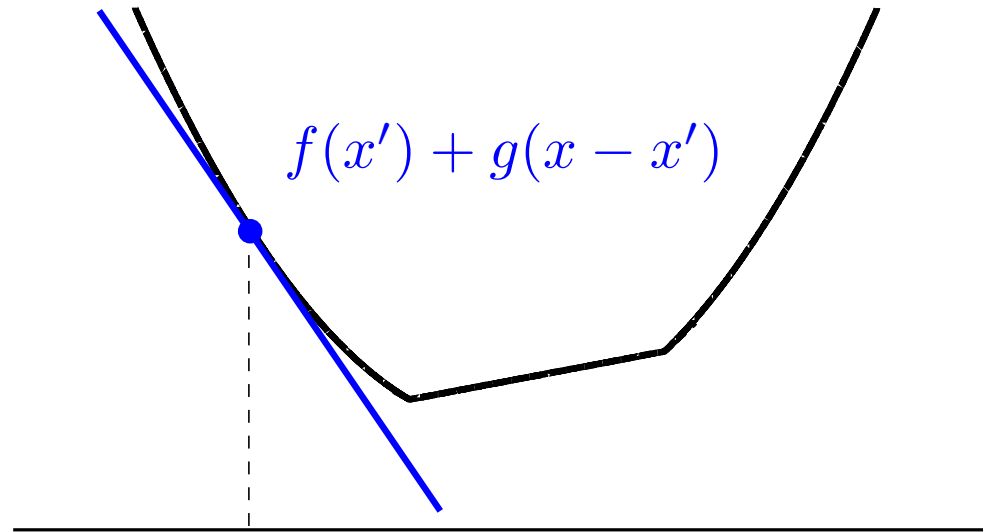




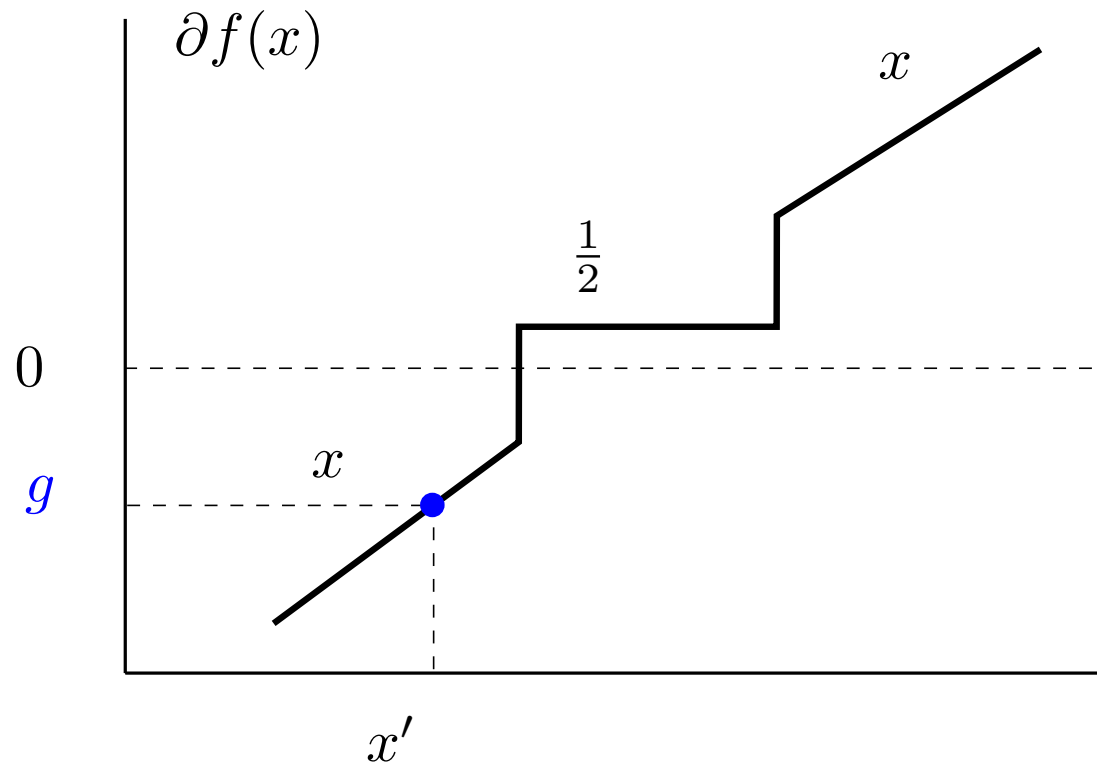
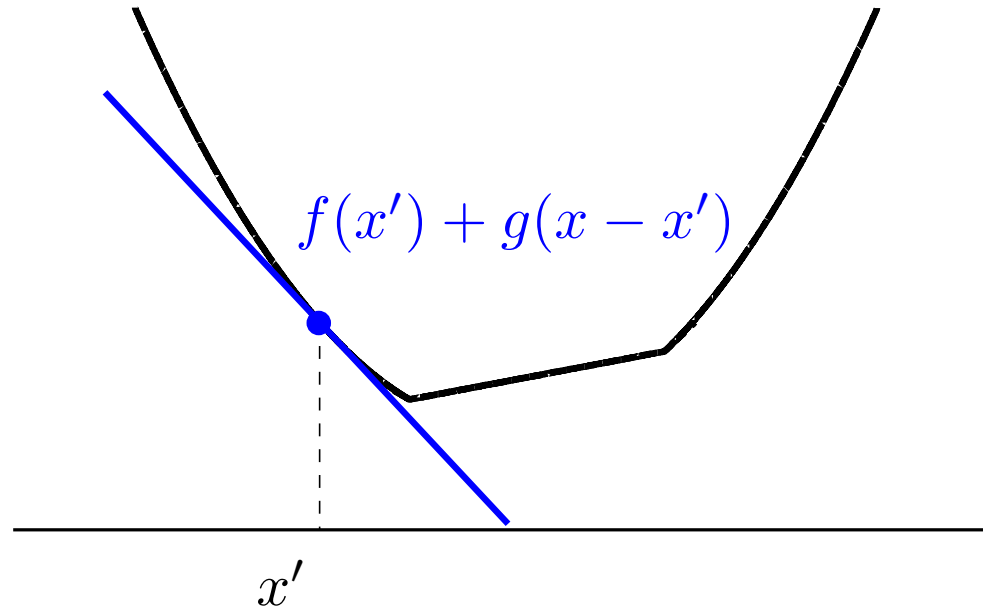




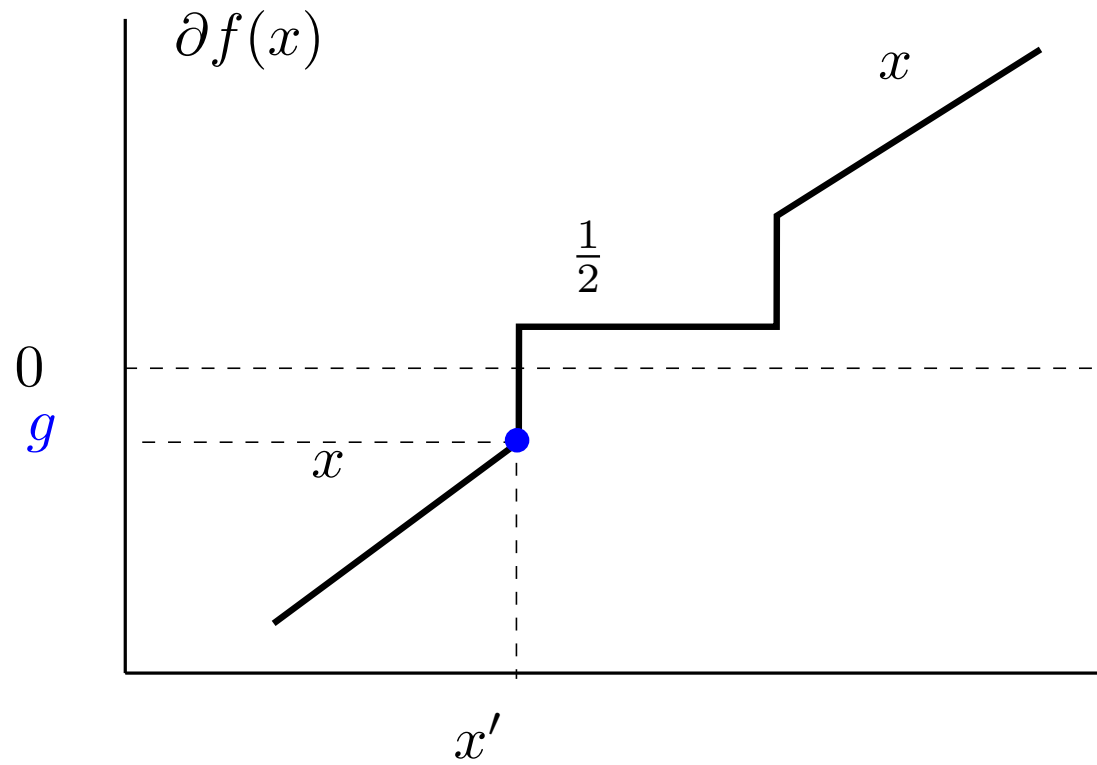
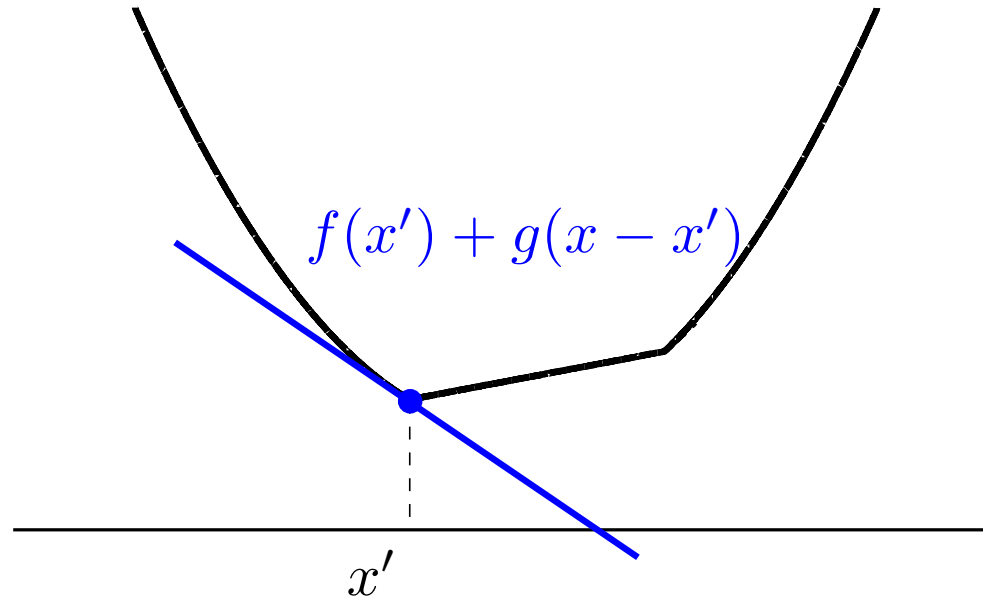
$$f(x) = \max\left\{\frac{1}{2}x^2, \frac{1}{2}x + 2\right\}$$



$$f(x) = \max\left\{\frac{1}{2}x^2, \frac{1}{2}x + 2\right\}$$

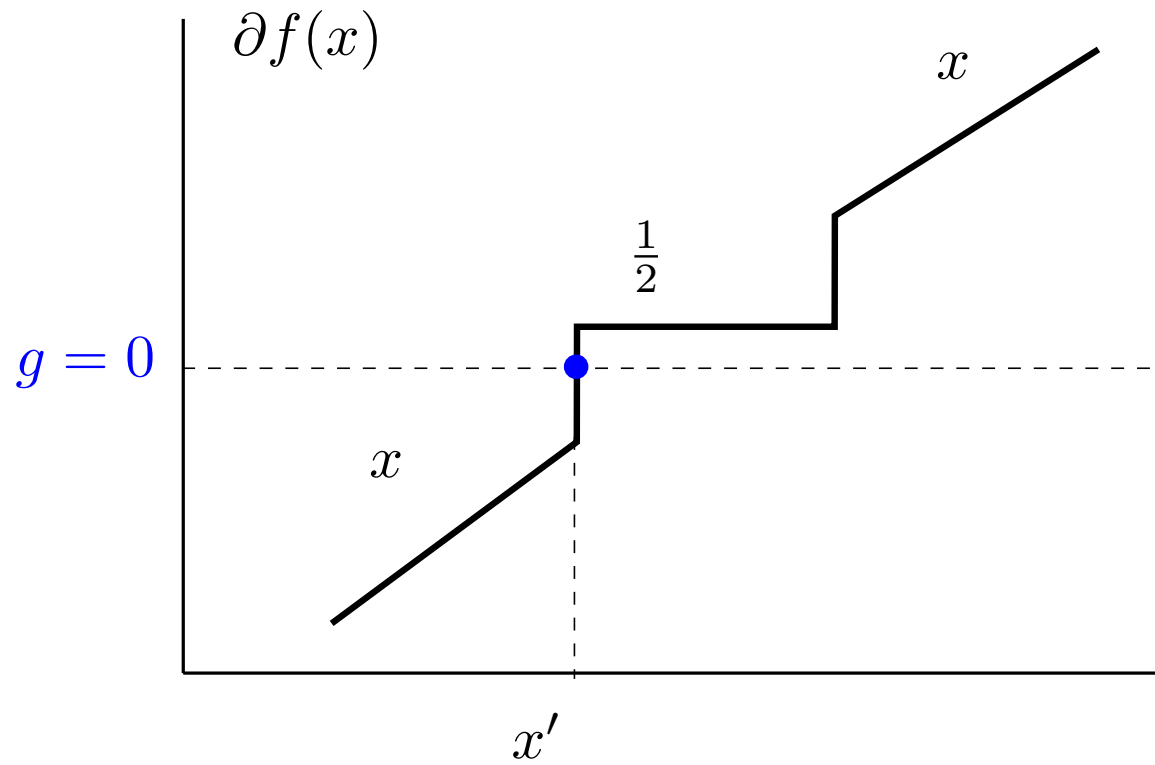
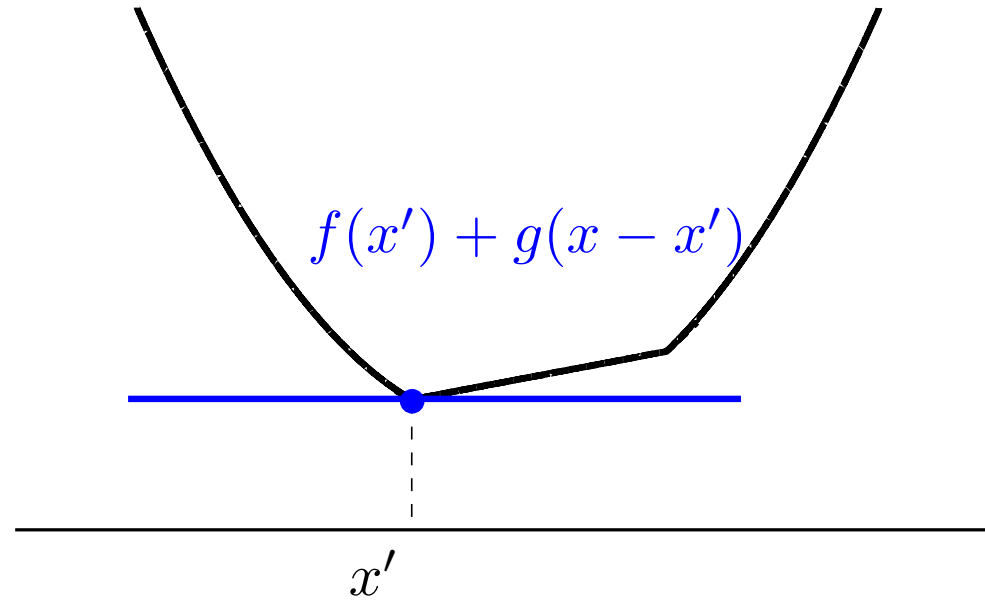


$$f(x) = \max\left\{\frac{1}{2}x^2, \frac{1}{2}x + 2\right\}$$

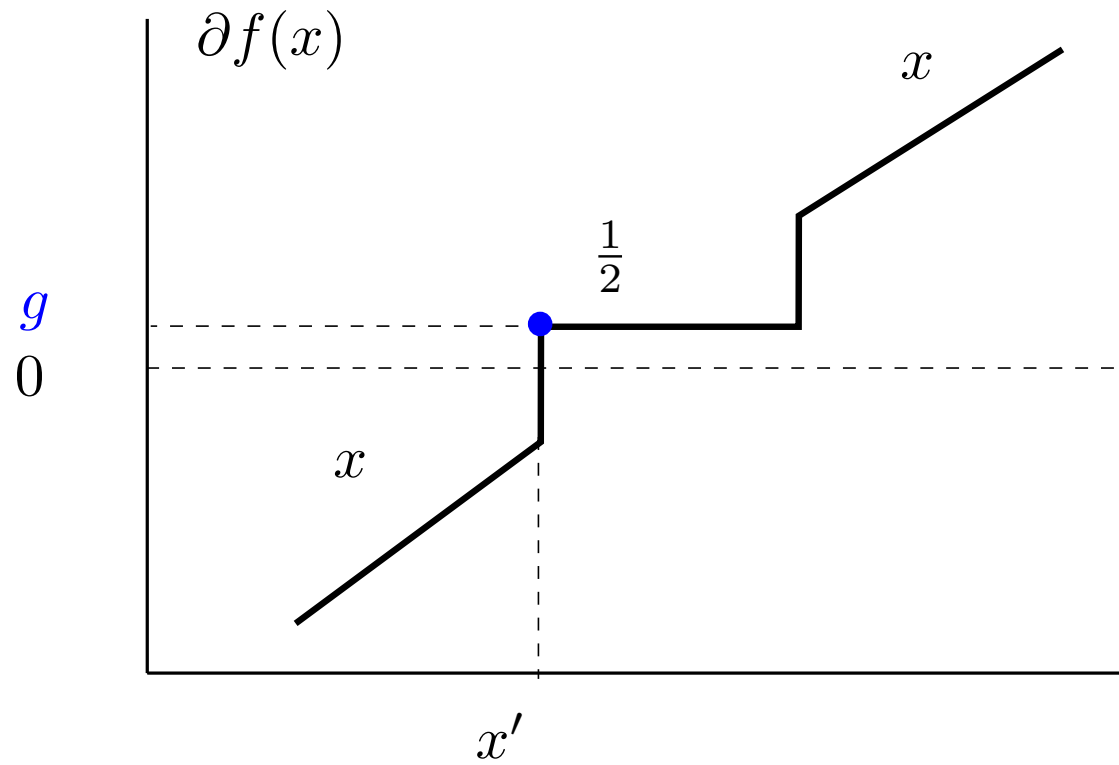
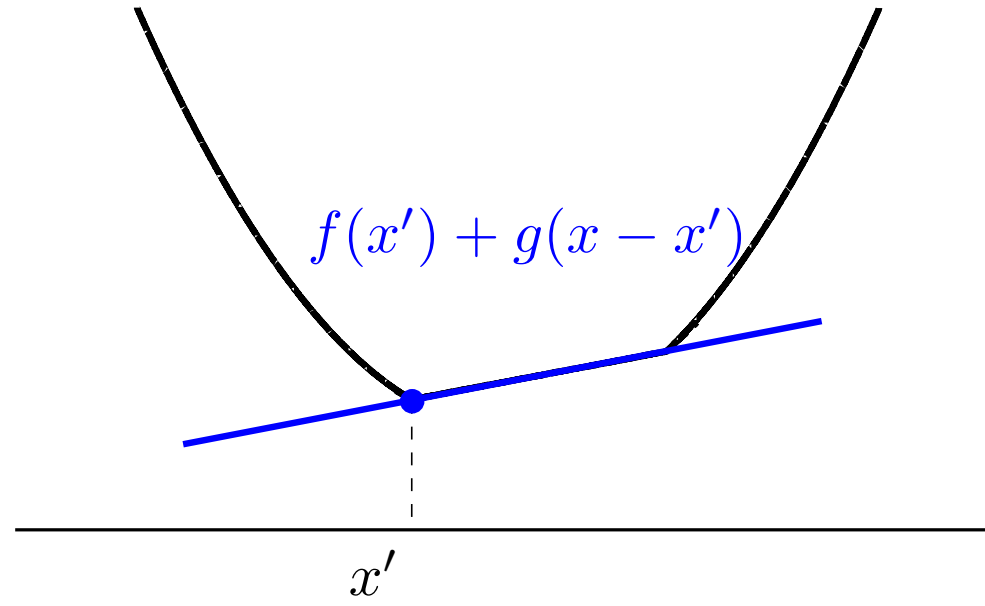




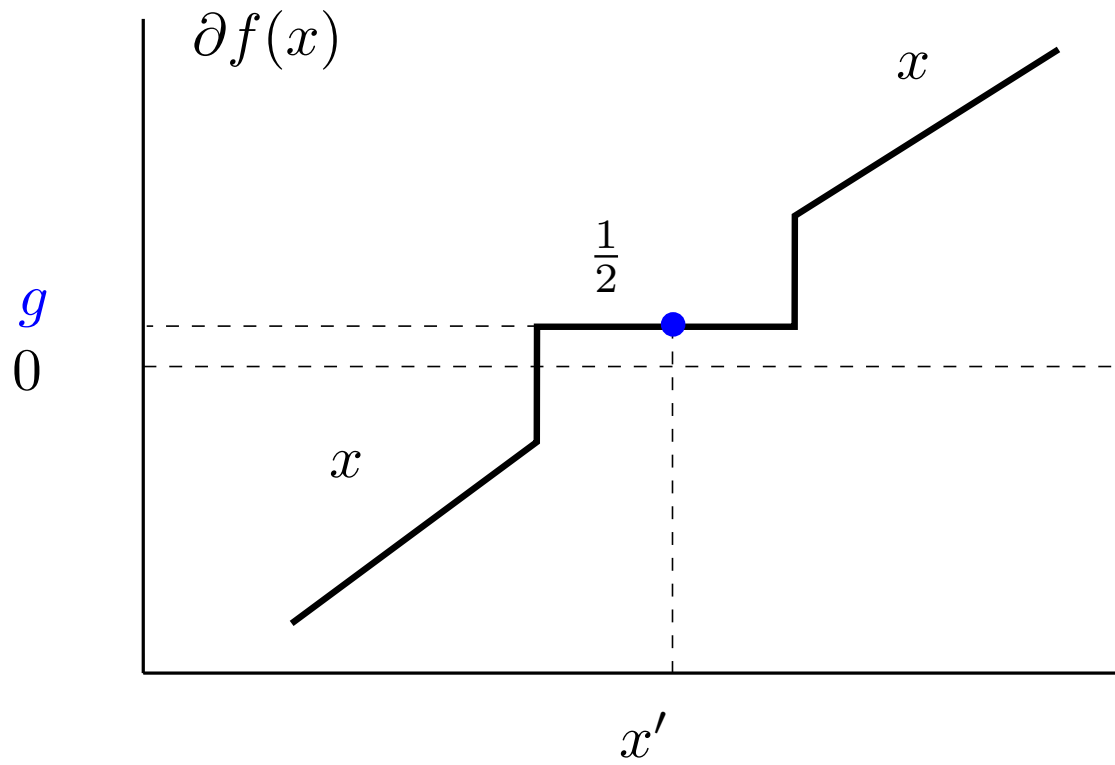
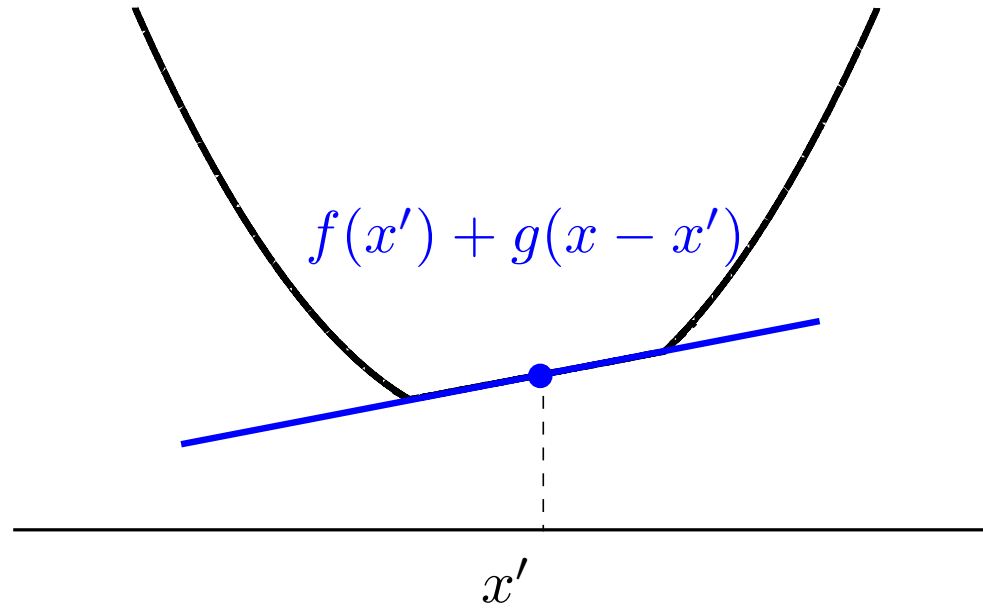
$$f(x) = \max\left\{\frac{1}{2}x^2, \frac{1}{2}x + 2\right\}$$



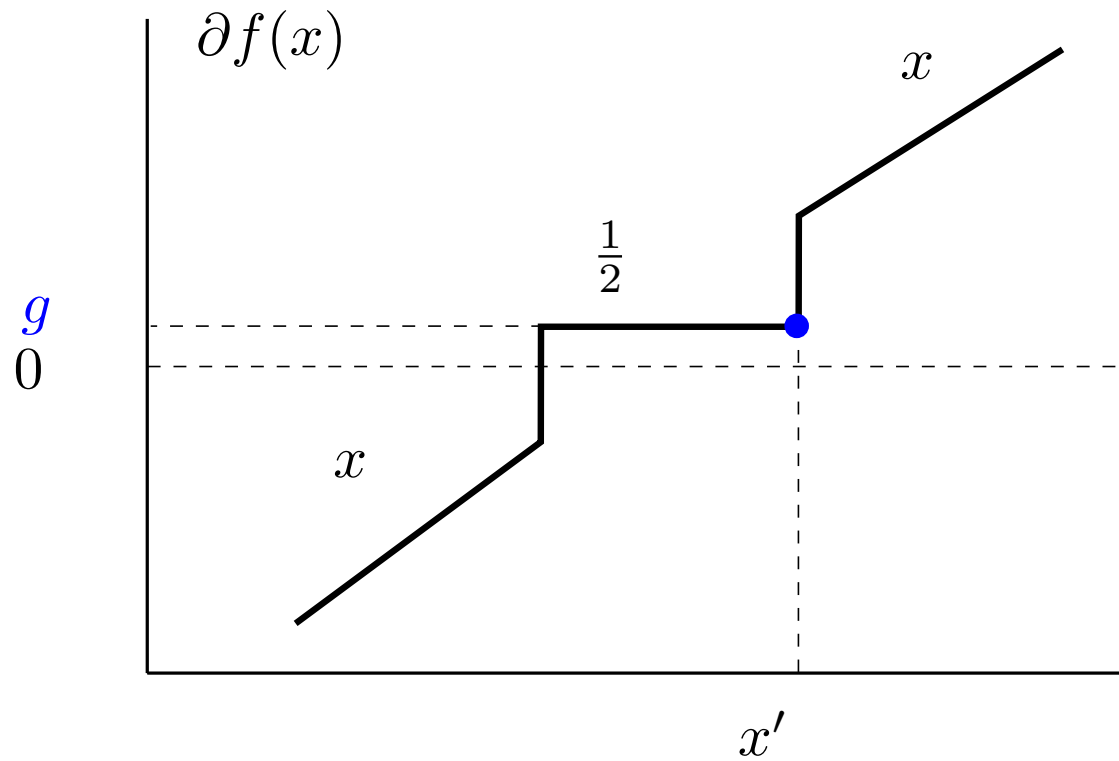
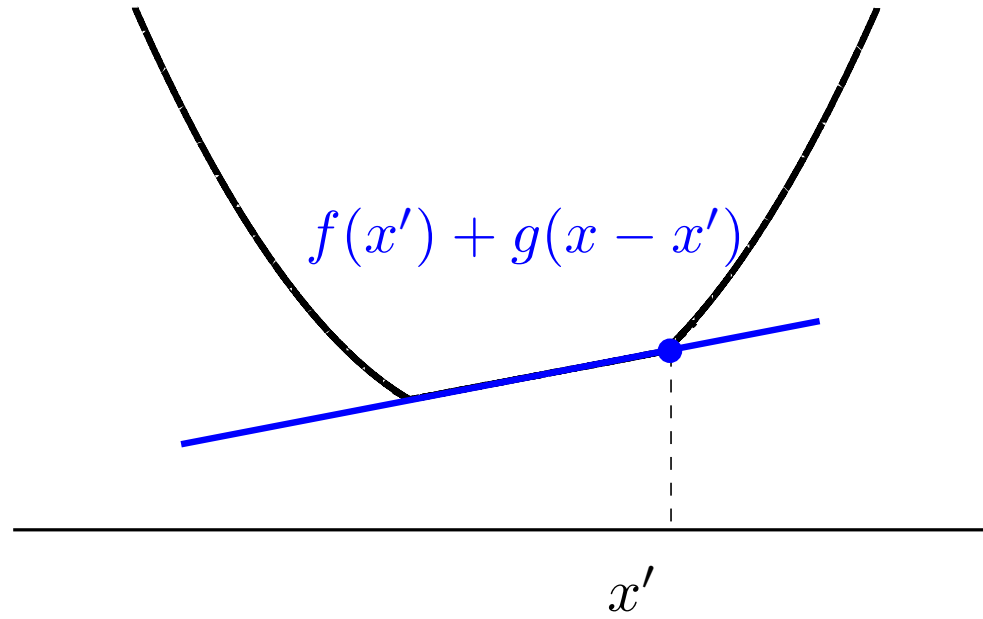
$$f(x) = \max\left\{\frac{1}{2}x^2, \frac{1}{2}x + 2\right\}$$



$$f(x) = \max\left\{\frac{1}{2}x^2, \frac{1}{2}x + 2\right\}$$



$$f(x) = \max\left\{\frac{1}{2}x^2, \frac{1}{2}x + 2\right\}$$



$$f(x) = \max\left\{\frac{1}{2}x^2, \frac{1}{2}x + 2\right\}$$

