

# COMS 4771 Lecture 10

1. Solving convex optimization problems

# SOLVING CONVEX OPTIMIZATION PROBLEMS

## Standard form of convex optimization problem

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^d} & f_0(\mathbf{x}) \\ \text{s.t.} & f_i(\mathbf{x}) \leq 0 \quad i = 1, 2, \dots, n \end{array}$$

(for convex functions  $f_0, f_1, \dots, f_n: \mathbb{R}^d \rightarrow \mathbb{R}$ ).

## Unconstrained convex optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$$

( $f$  is the convex objective function).

## Unconstrained convex optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$$

( $f$  is the convex objective function).

## Unconstrained convex optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$$

( $f$  is the convex objective function).

## Optimality condition for differentiable convex objectives

$\mathbf{x}$  is a global minimizer if and only if  $\nabla f(\mathbf{x}) = \mathbf{0}$ .

## Unconstrained convex optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$$

( $f$  is the convex objective function).

## Optimality condition for differentiable convex objectives

$\mathbf{x}$  is a global minimizer if and only if  $\nabla f(\mathbf{x}) = \mathbf{0}$ .

Unfortunately, can't always find closed-form solution to system of equations  $\nabla f(\mathbf{x}) = \mathbf{0}$ .  $\rightarrow$  Resort to iterative methods to find a solution.

## Iterative local optimization

**Main idea:** locally change  $\mathbf{x} \rightarrow \mathbf{x} + \boldsymbol{\delta}$  to improve its objective value  
 $f(\mathbf{x}) \rightarrow f(\mathbf{x} + \boldsymbol{\delta})$ .

## Iterative local optimization

**Main idea:** locally change  $\mathbf{x} \rightarrow \mathbf{x} + \delta$  to improve its objective value  $f(\mathbf{x}) \rightarrow f(\mathbf{x} + \delta)$ . **What should  $\delta$  be?**



## Iterative local optimization

**Main idea:** locally change  $\mathbf{x} \rightarrow \mathbf{x} + \boldsymbol{\delta}$  to improve its objective value  $f(\mathbf{x}) \rightarrow f(\mathbf{x} + \boldsymbol{\delta})$ . **What should  $\boldsymbol{\delta}$  be?**

**By convexity of  $f$ :**  $f(\mathbf{x} + \boldsymbol{\delta}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle$ .

## Iterative local optimization

**Main idea:** locally change  $\mathbf{x} \rightarrow \mathbf{x} + \boldsymbol{\delta}$  to improve its objective value  $f(\mathbf{x}) \rightarrow f(\mathbf{x} + \boldsymbol{\delta})$ . **What should  $\boldsymbol{\delta}$  be?**

**By convexity of  $f$ :**  $f(\mathbf{x} + \boldsymbol{\delta}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle$ .

If  $\langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle \geq 0$ , then

$$f(\mathbf{x} + \boldsymbol{\delta}) \geq f(\mathbf{x}).$$

# LOCAL OPTIMIZATION FOR CONVEX OBJECTIVES

## Iterative local optimization

**Main idea:** locally change  $\mathbf{x} \rightarrow \mathbf{x} + \boldsymbol{\delta}$  to improve its objective value  $f(\mathbf{x}) \rightarrow f(\mathbf{x} + \boldsymbol{\delta})$ . **What should  $\boldsymbol{\delta}$  be?**

**By convexity of  $f$ :**  $f(\mathbf{x} + \boldsymbol{\delta}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle$ .

If  $\langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle \geq 0$ , then

$$f(\mathbf{x} + \boldsymbol{\delta}) \geq f(\mathbf{x}). \quad \text{Clearly a bad direction.}$$

# LOCAL OPTIMIZATION FOR CONVEX OBJECTIVES

## Iterative local optimization

**Main idea:** locally change  $\mathbf{x} \rightarrow \mathbf{x} + \boldsymbol{\delta}$  to improve its objective value  $f(\mathbf{x}) \rightarrow f(\mathbf{x} + \boldsymbol{\delta})$ . **What should  $\boldsymbol{\delta}$  be?**

**By convexity of  $f$ :**  $f(\mathbf{x} + \boldsymbol{\delta}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle$ .

If  $\langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle \geq 0$ , then

$$f(\mathbf{x} + \boldsymbol{\delta}) \geq f(\mathbf{x}). \quad \text{Clearly a bad direction.}$$

**Moral:** to be useful, the change  $\boldsymbol{\delta}$  must satisfy

$$\langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle < 0.$$

For example,  $\boldsymbol{\delta} := -\eta \nabla f(\mathbf{x})$  for some  $\eta > 0$ :

$$\langle \nabla f(\mathbf{x}), -\eta \nabla f(\mathbf{x}) \rangle = -\eta \|\nabla f(\mathbf{x})\|_2^2 < 0$$

as long as  $\nabla f(\mathbf{x}) \neq \mathbf{0}$ .

## Gradient descent for differentiable objectives

- ▶ Start with some initial  $\mathbf{x}^{(1)} \in \mathbb{R}^d$ .
- ▶ For  $t = 1, 2, \dots$  until some stopping condition is satisfied.
  - ▶ Compute gradient of  $f$  at  $\mathbf{x}^{(t)}$ :

$$\boldsymbol{\lambda}^{(t)} := \nabla f(\mathbf{x}^{(t)}).$$

- ▶ Update:

$$\mathbf{x}^{(t+1)} := \mathbf{x}^{(t)} - \eta_t \boldsymbol{\lambda}^{(t)}.$$

## Gradient descent for differentiable objectives

- ▶ Start with some initial  $\mathbf{x}^{(1)} \in \mathbb{R}^d$ .
- ▶ For  $t = 1, 2, \dots$  until some stopping condition is satisfied.
  - ▶ Compute gradient of  $f$  at  $\mathbf{x}^{(t)}$ :

$$\boldsymbol{\lambda}^{(t)} := \nabla f(\mathbf{x}^{(t)}).$$

- ▶ Update:

$$\mathbf{x}^{(t+1)} := \mathbf{x}^{(t)} - \eta_t \boldsymbol{\lambda}^{(t)}.$$

Here,  $\eta_1, \eta_2, \dots > 0$  are the **step sizes**.

## Gradient descent for differentiable objectives

- ▶ Start with some initial  $\mathbf{x}^{(1)} \in \mathbb{R}^d$ .
- ▶ For  $t = 1, 2, \dots$  until some stopping condition is satisfied.
  - ▶ Compute gradient of  $f$  at  $\mathbf{x}^{(t)}$ :

$$\boldsymbol{\lambda}^{(t)} := \nabla f(\mathbf{x}^{(t)}).$$

- ▶ Update:

$$\mathbf{x}^{(t+1)} := \mathbf{x}^{(t)} - \eta_t \boldsymbol{\lambda}^{(t)}.$$

Here,  $\eta_1, \eta_2, \dots > 0$  are the **step sizes**. Common choices include:

1. Set  $\eta_t := c$  for some constant  $c > 0$ .
2. Set  $\eta_t := c/\sqrt{t}$  for some constant  $c > 0$ .
3. Set  $\eta_t$  using a line search procedure.

## Backtracking line search

**Goal:** given  $\mathbf{x} \in \mathbb{R}^d$  and  $\boldsymbol{\lambda} = \nabla f(\mathbf{x}) \in \mathbb{R}^d$ , find  $\eta > 0$  so that  $f(\mathbf{x} - \eta\boldsymbol{\lambda}) < f(\mathbf{x})$  by a reasonable amount.

- ▶ Start with  $\eta := 1$ .
- ▶ While  $f(\mathbf{x} - \eta\boldsymbol{\lambda}) > f(\mathbf{x}) - \frac{1}{2}\eta\|\boldsymbol{\lambda}\|_2^2$ : Set  $\eta := \frac{1}{2}\eta$ .



# STEP SIZES VIA LINE SEARCH

## Backtracking line search

**Goal:** given  $\mathbf{x} \in \mathbb{R}^d$  and  $\boldsymbol{\lambda} = \nabla f(\mathbf{x}) \in \mathbb{R}^d$ , find  $\eta > 0$  so that  $f(\mathbf{x} - \eta\boldsymbol{\lambda}) < f(\mathbf{x})$  by a reasonable amount.

- ▶ Start with  $\eta := 1$ .
- ▶ While  $f(\mathbf{x} - \eta\boldsymbol{\lambda}) > f(\mathbf{x}) - \frac{1}{2}\eta\|\boldsymbol{\lambda}\|_2^2$ : Set  $\eta := \frac{1}{2}\eta$ .

**Main idea:**  $f(\mathbf{x} - \eta\boldsymbol{\lambda}) \approx f(\mathbf{x}) - \eta\|\boldsymbol{\lambda}\|_2^2$  when  $\eta$  is small, so can optimistically hope to decrease value by about  $\eta\|\boldsymbol{\lambda}\|_2^2$ .

# STEP SIZES VIA LINE SEARCH

## Backtracking line search

**Goal:** given  $\mathbf{x} \in \mathbb{R}^d$  and  $\boldsymbol{\lambda} = \nabla f(\mathbf{x}) \in \mathbb{R}^d$ , find  $\eta > 0$  so that  $f(\mathbf{x} - \eta\boldsymbol{\lambda}) < f(\mathbf{x})$  by a reasonable amount.

- ▶ Start with  $\eta := 1$ .
- ▶ While  $f(\mathbf{x} - \eta\boldsymbol{\lambda}) > f(\mathbf{x}) - \frac{1}{2}\eta\|\boldsymbol{\lambda}\|_2^2$ : Set  $\eta := \frac{1}{2}\eta$ .

**Main idea:**  $f(\mathbf{x} - \eta\boldsymbol{\lambda}) \approx f(\mathbf{x}) - \eta\|\boldsymbol{\lambda}\|_2^2$  when  $\eta$  is small, so can optimistically hope to decrease value by about  $\eta\|\boldsymbol{\lambda}\|_2^2$ .

Settle for decreasing by  $\frac{1}{2}\eta\|\boldsymbol{\lambda}\|_2^2$ : upon termination of while-loop,

$$f(\mathbf{x} - \eta\boldsymbol{\lambda}) \leq f(\mathbf{x}) - \frac{1}{2}\eta\|\boldsymbol{\lambda}\|_2^2.$$

# STEP SIZES VIA LINE SEARCH

## Backtracking line search

**Goal:** given  $\mathbf{x} \in \mathbb{R}^d$  and  $\boldsymbol{\lambda} = \nabla f(\mathbf{x}) \in \mathbb{R}^d$ , find  $\eta > 0$  so that  $f(\mathbf{x} - \eta\boldsymbol{\lambda}) < f(\mathbf{x})$  by a reasonable amount.

- ▶ Start with  $\eta := 1$ .
- ▶ While  $f(\mathbf{x} - \eta\boldsymbol{\lambda}) > f(\mathbf{x}) - \frac{1}{2}\eta\|\boldsymbol{\lambda}\|_2^2$ : Set  $\eta := \frac{1}{2}\eta$ .

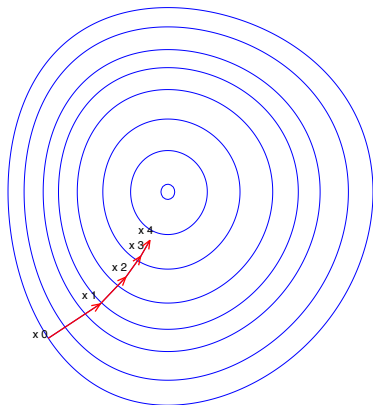
**Main idea:**  $f(\mathbf{x} - \eta\boldsymbol{\lambda}) \approx f(\mathbf{x}) - \eta\|\boldsymbol{\lambda}\|_2^2$  when  $\eta$  is small, so can optimistically hope to decrease value by about  $\eta\|\boldsymbol{\lambda}\|_2^2$ .

Settle for decreasing by  $\frac{1}{2}\eta\|\boldsymbol{\lambda}\|_2^2$ : upon termination of while-loop,

$$f(\mathbf{x} - \eta\boldsymbol{\lambda}) \leq f(\mathbf{x}) - \frac{1}{2}\eta\|\boldsymbol{\lambda}\|_2^2.$$

**Many other line search methods are possible.**

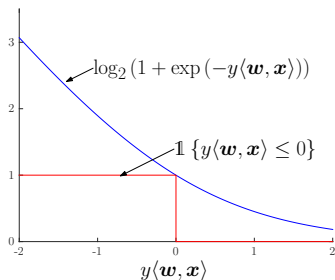
# ILLUSTRATION OF GRADIENT DESCENT



If  $f$  is convex (and satisfies some other smoothness and curvature conditions), then  $f(x^{(t)})$  converges to the optimal value at a geometric rate.

# EXAMPLE: (UNCONSTRAINED) LOGISTIC REGRESSION

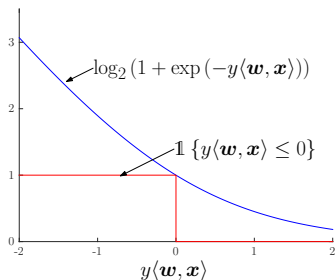
$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) := \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \ln(1 + \exp(-y\langle \mathbf{w}, \mathbf{x} \rangle))$$



We've already established that objective  $f(\mathbf{w})$  is convex.

# EXAMPLE: (UNCONSTRAINED) LOGISTIC REGRESSION

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) := \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \ln(1 + \exp(-y\langle \mathbf{w}, \mathbf{x} \rangle))$$



We've already established that objective  $f(\mathbf{w})$  is convex.

**Question:** How do we compute its gradient at a given point  $\mathbf{w} \in \mathbb{R}^d$ ?

# EXAMPLE: (UNCONSTRAINED) LOGISTIC REGRESSION

Gradient of  $f$  at  $\mathbf{w}$ :

$$\nabla f(\mathbf{w}) = -\frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \frac{1}{1 + e^{y\langle \mathbf{w}, \mathbf{x} \rangle}} y \mathbf{x}$$

# EXAMPLE: (UNCONSTRAINED) LOGISTIC REGRESSION

Gradient of  $f$  at  $\mathbf{w}$ :

$$\nabla f(\mathbf{w}) = -\frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \frac{1}{1 + e^{y\langle \mathbf{w}, \mathbf{x} \rangle}} y \mathbf{x}$$

**Gradient descent algorithm for logistic regression:**

- ▶ Start with some initial  $\mathbf{w}^{(1)} \in \mathbb{R}^d$ .
- ▶ For  $t = 1, 2, \dots$  until some stopping condition is satisfied.

$$\begin{aligned} \mathbf{w}^{(t+1)} &:= \mathbf{w}^{(t)} - \eta_t \nabla f(\mathbf{w}^{(t)}) \\ &= \mathbf{w}^{(t)} + \eta_t \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \frac{1}{1 + e^{y\langle \mathbf{w}^{(t)}, \mathbf{x} \rangle}} y \mathbf{x}. \end{aligned}$$



# STOPPING CONDITION

- ▶ In many applications of (convex) optimization, care about solving problems to very high precision.

**Example:** stop when gradient is close enough to zero ( $\|\nabla f(\mathbf{x})\|_2 \leq \epsilon$  for some small parameter  $\epsilon > 0$ ).

# STOPPING CONDITION

- ▶ In many applications of (convex) optimization, care about solving problems to very high precision.

**Example:** stop when gradient is close enough to zero ( $\|\nabla f(\mathbf{x})\|_2 \leq \epsilon$  for some small parameter  $\epsilon > 0$ ).

- ▶ For machine learning applications: optimization problem based on training data often just a means-to-an-end.

# STOPPING CONDITION

- ▶ In many applications of (convex) optimization, care about solving problems to very high precision.

**Example:** stop when gradient is close enough to zero ( $\|\nabla f(\mathbf{x})\|_2 \leq \epsilon$  for some small parameter  $\epsilon > 0$ ).

- ▶ For machine learning applications: optimization problem based on training data often just a means-to-an-end.

**We really just care about true error.**

# STOPPING CONDITION

- ▶ In many applications of (convex) optimization, care about solving problems to very high precision.

**Example:** stop when gradient is close enough to zero ( $\|\nabla f(\mathbf{x})\|_2 \leq \epsilon$  for some small parameter  $\epsilon > 0$ ).

- ▶ For machine learning applications: optimization problem based on training data often just a means-to-an-end.

**We really just care about true error.**

- ▶ Running gradient descent to convergence not strictly necessary: **may be beneficial to stop early (e.g., when hold-out error starts to increase significantly).**

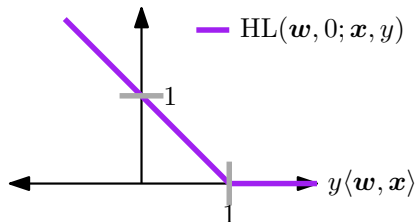
# NON-DIFFERENTIABILITY

## Non-differentiable convex objectives

Some convex functions  $f$  are not differentiable everywhere;  
**gradient descent not even well-specified for these problems.**

Example: hinge loss

$$f(\mathbf{w}) = \ell_{\text{hl}}(\mathbf{w}, 0; \mathbf{x}, y) = [1 - y\langle \mathbf{w}, \mathbf{x} \rangle]_+.$$



Not differentiable at  $\mathbf{w} \in \mathbb{R}^d$  where  $y\langle \mathbf{w}, \mathbf{x} \rangle = 1$ .

## Subgradients

Although not every function  $f$  is differentiable everywhere, every **convex function**  $f$  has **subgradients everywhere**<sup>†</sup>.

We say  $\boldsymbol{\lambda} \in \mathbb{R}^d$  is a **subgradient** of a function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  at  $\mathbf{x}_0 \in \mathbb{R}^d$  if

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \langle \boldsymbol{\lambda}, \mathbf{x} - \mathbf{x}_0 \rangle \quad \forall \mathbf{x} \in \mathbb{R}^d.$$

<sup>†</sup>some technical conditions apply.

## Subgradients

Although not every function  $f$  is differentiable everywhere, every **convex function**  $f$  has **subgradients everywhere**<sup>†</sup>.

We say  $\boldsymbol{\lambda} \in \mathbb{R}^d$  is a **subgradient** of a function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  at  $\boldsymbol{x}_0 \in \mathbb{R}^d$  if

$$f(\boldsymbol{x}) \geq f(\boldsymbol{x}_0) + \langle \boldsymbol{\lambda}, \boldsymbol{x} - \boldsymbol{x}_0 \rangle \quad \forall \boldsymbol{x} \in \mathbb{R}^d.$$

In other words, a subgradient of a convex function  $f$  at a point  $\boldsymbol{x}_0$  specifies an **affine lower bound** on the function . . .

<sup>†</sup>some technical conditions apply.

## Subgradients

Although not every function  $f$  is differentiable everywhere, every **convex function**  $f$  has **subgradients** *everywhere*<sup>†</sup>.

We say  $\boldsymbol{\lambda} \in \mathbb{R}^d$  is a **subgradient** of a function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  at  $\mathbf{x}_0 \in \mathbb{R}^d$  if

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \langle \boldsymbol{\lambda}, \mathbf{x} - \mathbf{x}_0 \rangle \quad \forall \mathbf{x} \in \mathbb{R}^d.$$

In other words, a subgradient of a convex function  $f$  at a point  $\mathbf{x}_0$  specifies an **affine lower bound** on the function ...

... **just like the gradient in the case of a differentiable function:**

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0 \rangle \quad \forall \mathbf{x} \in \mathbb{R}^d.$$

<sup>†</sup>some technical conditions apply.



## Subgradients

Although not every function  $f$  is differentiable everywhere, every **convex function**  $f$  has **subgradients everywhere**<sup>†</sup>.

We say  $\lambda \in \mathbb{R}^d$  is a **subgradient** of a function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  at  $x_0 \in \mathbb{R}^d$  if

$$f(x) \geq f(x_0) + \langle \lambda, x - x_0 \rangle \quad \forall x \in \mathbb{R}^d.$$

In other words, a subgradient of a convex function  $f$  at a point  $x_0$  specifies an **affine lower bound** on the function ...

... **just like the gradient in the case of a differentiable function:**

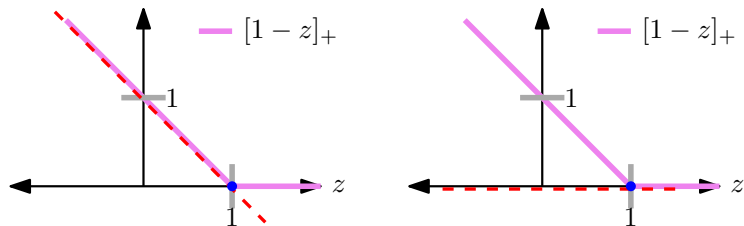
$$f(x) \geq f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle \quad \forall x \in \mathbb{R}^d.$$

**There might be many subgradients at a given point  $x_0$ —i.e., many affine lower bounds:** call the entire set the **subdifferential of  $f$  at  $x_0$** ,  $\partial f(x_0)$ .

<sup>†</sup>some technical conditions apply.

# EXAMPLE: SUBGRADIENT OF HINGE LOSS (SORTA)

Consider one-dimensional function  $f(z) := [1 - z]_+ = \max\{0, 1 - z\}$ .



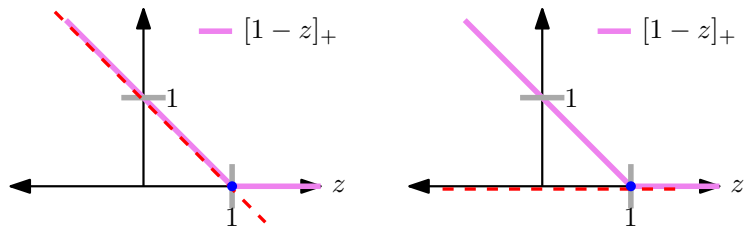
Two subgradients of  $f$  at  $z = 1$ :  $-1$  and  $0$ .

$$f(z) \geq f(1) + (-1) \cdot (z - 1) = 1 - z;$$

$$f(z) \geq f(1) + (0) \cdot (z - 1) = 0.$$

# EXAMPLE: SUBGRADIENT OF HINGE LOSS (SORTA)

Consider one-dimensional function  $f(z) := [1 - z]_+ = \max\{0, 1 - z\}$ .



Two subgradients of  $f$  at  $z = 1$ :  $-1$  and  $0$ .

$$f(z) \geq f(1) + (-1) \cdot (z - 1) = 1 - z;$$

$$f(z) \geq f(1) + (0) \cdot (z - 1) = 0.$$

Actually, **infinitely-many subgradients** at  $z = 1$ : all  $\lambda \in [-1, 0]$  satisfy

$$f(z) \geq f(1) + \lambda \cdot (z - 1).$$

# SUBGRADIENT CALCULUS

Suppose  $g, g_1, g_2$  are convex functions.

Below, sufficient conditions under which  $f$  is convex, and corresponding subdifferential:

- ▶ **Addition:** If  $f(\mathbf{x}) = g_1(\mathbf{x}) + g_2(\mathbf{x})$ , then

$$\partial f(\mathbf{x}) = \{\boldsymbol{\lambda} + \boldsymbol{\nu} : \boldsymbol{\lambda} \in \partial g_1(\mathbf{x}), \boldsymbol{\nu} \in \partial g_2(\mathbf{x})\}.$$

# SUBGRADIENT CALCULUS

Suppose  $g, g_1, g_2$  are convex functions.

Below, sufficient conditions under which  $f$  is convex, and corresponding subdifferential:

- ▶ **Addition:** If  $f(\mathbf{x}) = g_1(\mathbf{x}) + g_2(\mathbf{x})$ , then

$$\partial f(\mathbf{x}) = \{\boldsymbol{\lambda} + \boldsymbol{\nu} : \boldsymbol{\lambda} \in \partial g_1(\mathbf{x}), \boldsymbol{\nu} \in \partial g_2(\mathbf{x})\}.$$

- ▶ **Positive scaling:** If  $f(\mathbf{x}) = \alpha \cdot g(\mathbf{x})$  for some  $\alpha > 0$ , then

$$\partial f(\mathbf{x}) = \{\alpha \boldsymbol{\lambda} : \boldsymbol{\lambda} \in \partial g(\mathbf{x})\}.$$

# SUBGRADIENT CALCULUS

Suppose  $g, g_1, g_2$  are convex functions.

Below, sufficient conditions under which  $f$  is convex, and corresponding subdifferential:

- ▶ **Addition:** If  $f(\mathbf{x}) = g_1(\mathbf{x}) + g_2(\mathbf{x})$ , then

$$\partial f(\mathbf{x}) = \{\boldsymbol{\lambda} + \boldsymbol{\nu} : \boldsymbol{\lambda} \in \partial g_1(\mathbf{x}), \boldsymbol{\nu} \in \partial g_2(\mathbf{x})\}.$$

- ▶ **Positive scaling:** If  $f(\mathbf{x}) = \alpha \cdot g(\mathbf{x})$  for some  $\alpha > 0$ , then

$$\partial f(\mathbf{x}) = \{\alpha \boldsymbol{\lambda} : \boldsymbol{\lambda} \in \partial g(\mathbf{x})\}.$$

- ▶ **Affine composition:** If  $f(\mathbf{x}) = g(\mathbf{A}\mathbf{x} + \mathbf{b})$ , then

$$\partial f(\mathbf{x}) = \{\mathbf{A}^\top \boldsymbol{\lambda} : \boldsymbol{\lambda} \in \partial g(\mathbf{A}\mathbf{x} + \mathbf{b})\}.$$

# SUBGRADIENT CALCULUS (CONTINUED)

Suppose  $g_1, g_2$  are convex functions.

Below, sufficient conditions under which  $f$  is convex, and corresponding subdifferential:

- ▶ **Max of convex functions:** If  $f(\mathbf{x}) = \max\{g_1(\mathbf{x}), g_2(\mathbf{x})\}$ , then

$$\partial f(\mathbf{x}) = \begin{cases} \partial g_1(\mathbf{x}) & \text{if } g_1(\mathbf{x}) > g_2(\mathbf{x}); \\ \partial g_2(\mathbf{x}) & \text{if } g_1(\mathbf{x}) < g_2(\mathbf{x}); \\ \text{conv}(\partial g_1(\mathbf{x}) \cup \partial g_2(\mathbf{x})) & \text{if } g_1(\mathbf{x}) = g_2(\mathbf{x}). \end{cases}$$

## EXAMPLE: SUBGRADIENT OF HINGE LOSS (REALLY)

**Hinge loss function:**  $f(\mathbf{w}) := [1 - y\langle \mathbf{w}, \mathbf{x} \rangle]_+ = \max\{0, 1 - y\langle \mathbf{w}, \mathbf{x} \rangle\}$ .



## EXAMPLE: SUBGRADIENT OF HINGE LOSS (REALLY)

**Hinge loss function:**  $f(\mathbf{w}) := [1 - y\langle \mathbf{w}, \mathbf{x} \rangle]_+ = \max\{0, 1 - y\langle \mathbf{w}, \mathbf{x} \rangle\}$ .

$f(\mathbf{w}) = \max\{g_1(\mathbf{w}), g_2(\mathbf{w})\}$  where  $g_1(\mathbf{w}) = 0$  and  $g_2(\mathbf{w}) = 1 - y\langle \mathbf{w}, \mathbf{x} \rangle$ .

## EXAMPLE: SUBGRADIENT OF HINGE LOSS (REALLY)

**Hinge loss function:**  $f(\mathbf{w}) := [1 - y\langle \mathbf{w}, \mathbf{x} \rangle]_+ = \max\{0, 1 - y\langle \mathbf{w}, \mathbf{x} \rangle\}$ .

$f(\mathbf{w}) = \max\{g_1(\mathbf{w}), g_2(\mathbf{w})\}$  where  $g_1(\mathbf{w}) = 0$  and  $g_2(\mathbf{w}) = 1 - y\langle \mathbf{w}, \mathbf{x} \rangle$ .

►  $\partial g_1(\mathbf{w}) = \{\mathbf{0}\}$  and  $\partial g_2(\mathbf{w}) = \{-y\mathbf{x}\}$ .

# EXAMPLE: SUBGRADIENT OF HINGE LOSS (REALLY)

**Hinge loss function:**  $f(\mathbf{w}) := [1 - y\langle \mathbf{w}, \mathbf{x} \rangle]_+ = \max\{0, 1 - y\langle \mathbf{w}, \mathbf{x} \rangle\}$ .

$f(\mathbf{w}) = \max\{g_1(\mathbf{w}), g_2(\mathbf{w})\}$  where  $g_1(\mathbf{w}) = 0$  and  $g_2(\mathbf{w}) = 1 - y\langle \mathbf{w}, \mathbf{x} \rangle$ .

- ▶  $\partial g_1(\mathbf{w}) = \{\mathbf{0}\}$  and  $\partial g_2(\mathbf{w}) = \{-y\mathbf{x}\}$ .
- ▶ If  $f(\mathbf{w}) = g_1(\mathbf{w}) = 0 > 1 - y\langle \mathbf{w}, \mathbf{x} \rangle = g_2(\mathbf{w})$  (i.e., if  $y\langle \mathbf{w}, \mathbf{x} \rangle > 1$ ), then

$$\partial f(\mathbf{w}) = \{\mathbf{0}\}.$$

# EXAMPLE: SUBGRADIENT OF HINGE LOSS (REALLY)

**Hinge loss function:**  $f(\mathbf{w}) := [1 - y\langle \mathbf{w}, \mathbf{x} \rangle]_+ = \max\{0, 1 - y\langle \mathbf{w}, \mathbf{x} \rangle\}$ .

$f(\mathbf{w}) = \max\{g_1(\mathbf{w}), g_2(\mathbf{w})\}$  where  $g_1(\mathbf{w}) = 0$  and  $g_2(\mathbf{w}) = 1 - y\langle \mathbf{w}, \mathbf{x} \rangle$ .

▶  $\partial g_1(\mathbf{w}) = \{\mathbf{0}\}$  and  $\partial g_2(\mathbf{w}) = \{-y\mathbf{x}\}$ .

▶ If  $f(\mathbf{w}) = g_1(\mathbf{w}) = 0 > 1 - y\langle \mathbf{w}, \mathbf{x} \rangle = g_2(\mathbf{w})$  (i.e., if  $y\langle \mathbf{w}, \mathbf{x} \rangle > 1$ ), then

$$\partial f(\mathbf{w}) = \{\mathbf{0}\}.$$

▶ If  $f(\mathbf{w}) = g_2(\mathbf{w}) = 1 - y\langle \mathbf{w}, \mathbf{x} \rangle > 0 = g_1(\mathbf{w})$  (i.e., if  $y\langle \mathbf{w}, \mathbf{x} \rangle < 1$ ), then

$$\partial f(\mathbf{w}) = \{-y\mathbf{x}\}.$$

# EXAMPLE: SUBGRADIENT OF HINGE LOSS (REALLY)

**Hinge loss function:**  $f(\mathbf{w}) := [1 - y\langle \mathbf{w}, \mathbf{x} \rangle]_+ = \max\{0, 1 - y\langle \mathbf{w}, \mathbf{x} \rangle\}$ .

$f(\mathbf{w}) = \max\{g_1(\mathbf{w}), g_2(\mathbf{w})\}$  where  $g_1(\mathbf{w}) = 0$  and  $g_2(\mathbf{w}) = 1 - y\langle \mathbf{w}, \mathbf{x} \rangle$ .

▶  $\partial g_1(\mathbf{w}) = \{\mathbf{0}\}$  and  $\partial g_2(\mathbf{w}) = \{-y\mathbf{x}\}$ .

▶ If  $f(\mathbf{w}) = g_1(\mathbf{w}) = 0 > 1 - y\langle \mathbf{w}, \mathbf{x} \rangle = g_2(\mathbf{w})$  (i.e., if  $y\langle \mathbf{w}, \mathbf{x} \rangle > 1$ ), then

$$\partial f(\mathbf{w}) = \{\mathbf{0}\}.$$

▶ If  $f(\mathbf{w}) = g_2(\mathbf{w}) = 1 - y\langle \mathbf{w}, \mathbf{x} \rangle > 0 = g_1(\mathbf{w})$  (i.e., if  $y\langle \mathbf{w}, \mathbf{x} \rangle < 1$ ), then

$$\partial f(\mathbf{w}) = \{-y\mathbf{x}\}.$$

▶ If  $f(\mathbf{w}) = g_1(\mathbf{w}) = g_2(\mathbf{w})$  (i.e., if  $y\langle \mathbf{w}, \mathbf{x} \rangle = 1$ ), then

$$\partial f(\mathbf{w}) = \text{conv}\{\mathbf{0}, -y\mathbf{x}\}.$$

# EXAMPLE: SUBGRADIENT OF HINGE LOSS (REALLY)

**Hinge loss function:**  $f(\mathbf{w}) := [1 - y\langle \mathbf{w}, \mathbf{x} \rangle]_+ = \max\{0, 1 - y\langle \mathbf{w}, \mathbf{x} \rangle\}$ .

$f(\mathbf{w}) = \max\{g_1(\mathbf{w}), g_2(\mathbf{w})\}$  where  $g_1(\mathbf{w}) = 0$  and  $g_2(\mathbf{w}) = 1 - y\langle \mathbf{w}, \mathbf{x} \rangle$ .

▶  $\partial g_1(\mathbf{w}) = \{\mathbf{0}\}$  and  $\partial g_2(\mathbf{w}) = \{-y\mathbf{x}\}$ .

▶ If  $f(\mathbf{w}) = g_1(\mathbf{w}) = 0 > 1 - y\langle \mathbf{w}, \mathbf{x} \rangle = g_2(\mathbf{w})$  (i.e., if  $y\langle \mathbf{w}, \mathbf{x} \rangle > 1$ ), then

$$\partial f(\mathbf{w}) = \{\mathbf{0}\}.$$

▶ If  $f(\mathbf{w}) = g_2(\mathbf{w}) = 1 - y\langle \mathbf{w}, \mathbf{x} \rangle > 0 = g_1(\mathbf{w})$  (i.e., if  $y\langle \mathbf{w}, \mathbf{x} \rangle < 1$ ), then

$$\partial f(\mathbf{w}) = \{-y\mathbf{x}\}.$$

▶ If  $f(\mathbf{w}) = g_1(\mathbf{w}) = g_2(\mathbf{w})$  (i.e., if  $y\langle \mathbf{w}, \mathbf{x} \rangle = 1$ ), then

$$\partial f(\mathbf{w}) = \text{conv}\{\mathbf{0}, -y\mathbf{x}\}.$$

Can also derive using affine composition rule.

## Subgradient descent for general convex objectives

- ▶ Start with some initial  $\mathbf{x}^{(1)} \in \mathbb{R}^d$ .
- ▶ For  $t = 1, 2, \dots$  until some stopping condition is satisfied.
  - ▶ Compute *any* subgradient  $\boldsymbol{\lambda}^{(t)} \in \partial f(\mathbf{x}^{(t)})$ .
  - ▶ Update:

$$\mathbf{x}^{(t+1)} := \mathbf{x}^{(t)} - \eta_t \boldsymbol{\lambda}^{(t)}.$$

## EXAMPLE: SOFT-MARGIN SVM (WITHOUT OFFSET)

$$\min_{\mathbf{w} \in \mathbb{R}^d} \quad f(\mathbf{w}) := \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} [1 - y\langle \mathbf{w}, \mathbf{x} \rangle]_+$$

$f$  is the sum of convex functions, and hence is convex.



## EXAMPLE: SOFT-MARGIN SVM (WITHOUT OFFSET)

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) := \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} [1 - y\langle \mathbf{w}, \mathbf{x} \rangle]_+$$

$f$  is the sum of convex functions, and hence is convex.

**Question:** How do we compute a subgradient  $\mathbf{g}$  of  $f$  at a given point  $\mathbf{w} \in \mathbb{R}^d$ ?

## EXAMPLE: SOFT-MARGIN SVM (WITHOUT OFFSET)

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) := \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} [1 - y\langle \mathbf{w}, \mathbf{x} \rangle]_+$$

$f$  is the sum of convex functions, and hence is convex.

**Question:** How do we compute a subgradient  $\mathbf{g}$  of  $f$  at a given point  $\mathbf{w} \in \mathbb{R}^d$ ?

$$\partial f(\mathbf{w}) \ni \mathbf{g} = \lambda \mathbf{w} + \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \begin{cases} \mathbf{0} & \text{if } 1 - y\langle \mathbf{w}, \mathbf{x} \rangle < 0; \\ -y\mathbf{x} & \text{if } 1 - y\langle \mathbf{w}, \mathbf{x} \rangle \geq 0. \end{cases}$$

# EXAMPLE: SOFT-MARGIN SVM (WITHOUT OFFSET)

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) := \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} [1 - y\langle \mathbf{w}, \mathbf{x} \rangle]_+$$

$f$  is the sum of convex functions, and hence is convex.

**Question:** How do we compute a subgradient  $\mathbf{g}$  of  $f$  at a given point  $\mathbf{w} \in \mathbb{R}^d$ ?

$$\partial f(\mathbf{w}) \ni \mathbf{g} = \lambda \mathbf{w} + \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \begin{cases} \mathbf{0} & \text{if } 1 - y\langle \mathbf{w}, \mathbf{x} \rangle < 0; \\ -y\mathbf{x} & \text{if } 1 - y\langle \mathbf{w}, \mathbf{x} \rangle \geq 0. \end{cases}$$

Also fine:

$$\partial f(\mathbf{w}) \ni \mathbf{g} = \lambda \mathbf{w} + \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \begin{cases} \mathbf{0} & \text{if } 1 - y\langle \mathbf{w}, \mathbf{x} \rangle < 0; \\ -\frac{1}{3}y\mathbf{x} & \text{if } 1 - y\langle \mathbf{w}, \mathbf{x} \rangle = 0; \\ -y\mathbf{x} & \text{if } 1 - y\langle \mathbf{w}, \mathbf{x} \rangle > 0. \end{cases}$$

# EXAMPLE: SOFT-MARGIN SVM (WITHOUT OFFSET)

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) := \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} [1 - y\langle \mathbf{w}, \mathbf{x} \rangle]_+$$

$f$  is the sum of convex functions, and hence is convex.

**Question:** How do we compute a subgradient  $\mathbf{g}$  of  $f$  at a given point  $\mathbf{w} \in \mathbb{R}^d$ ?

$$\partial f(\mathbf{w}) \ni \mathbf{g} = \lambda \mathbf{w} + \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \begin{cases} \mathbf{0} & \text{if } 1 - y\langle \mathbf{w}, \mathbf{x} \rangle < 0; \\ -y\mathbf{x} & \text{if } 1 - y\langle \mathbf{w}, \mathbf{x} \rangle \geq 0. \end{cases}$$

Also fine:

$$\partial f(\mathbf{w}) \ni \mathbf{g} = \lambda \mathbf{w} + \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \begin{cases} \mathbf{0} & \text{if } 1 - y\langle \mathbf{w}, \mathbf{x} \rangle < 0; \\ -\frac{1}{3}y\mathbf{x} & \text{if } 1 - y\langle \mathbf{w}, \mathbf{x} \rangle = 0; \\ -y\mathbf{x} & \text{if } 1 - y\langle \mathbf{w}, \mathbf{x} \rangle > 0. \end{cases}$$

In practice, usually don't have examples with  $y\langle \mathbf{w}, \mathbf{x} \rangle = 1$  *exactly* anyway.

# EXAMPLE: SOFT-MARGIN SVM (WITHOUT OFFSET)

## Subgradient descent algorithm for soft-margin SVM:

- ▶ Start with some initial  $\mathbf{w}^{(1)} \in \mathbb{R}^d$ .
- ▶ For  $t = 1, 2, \dots$  until some stopping condition is satisfied.

$$\begin{aligned}\mathbf{w}^{(t+1)} &:= \mathbf{w}^{(t)} - \eta_t \left( \lambda \mathbf{w}^{(t)} + \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \begin{cases} \mathbf{0} & \text{if } 1 - y \langle \mathbf{w}^{(t)}, \mathbf{x} \rangle < 0; \\ -y\mathbf{x} & \text{if } 1 - y \langle \mathbf{w}^{(t)}, \mathbf{x} \rangle \geq 0 \end{cases} \right) \\ &= (1 - \lambda \eta_t) \mathbf{w}^{(t)} + \eta_t \frac{1}{|S|} \sum_{\substack{(\mathbf{x}, y) \in S: \\ y \langle \mathbf{w}^{(t)}, \mathbf{x} \rangle \leq 1}} y \mathbf{x}.\end{aligned}$$

# EXAMPLE: SOFT-MARGIN SVM (WITHOUT OFFSET)

## Subgradient descent algorithm for soft-margin SVM:

- ▶ Start with some initial  $\mathbf{w}^{(1)} \in \mathbb{R}^d$ .
- ▶ For  $t = 1, 2, \dots$  until some stopping condition is satisfied.

$$\begin{aligned}\mathbf{w}^{(t+1)} &:= \mathbf{w}^{(t)} - \eta_t \left( \lambda \mathbf{w}^{(t)} + \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \begin{cases} \mathbf{0} & \text{if } 1 - y \langle \mathbf{w}^{(t)}, \mathbf{x} \rangle < 0; \\ -y \mathbf{x} & \text{if } 1 - y \langle \mathbf{w}^{(t)}, \mathbf{x} \rangle \geq 0 \end{cases} \right) \\ &= (1 - \lambda \eta_t) \mathbf{w}^{(t)} + \eta_t \frac{1}{|S|} \sum_{\substack{(\mathbf{x}, y) \in S: \\ y \langle \mathbf{w}^{(t)}, \mathbf{x} \rangle \leq 1}} y \mathbf{x}.\end{aligned}$$

**Note effect of regularization term  $\frac{\lambda}{2} \|\mathbf{w}\|_2^2$  (whenever  $\eta_t < 1/\lambda$ ):**

Shrink  $\mathbf{w}^{(t)}$  by a factor  $1 - \lambda \eta_t$  before updating with subgradient of loss term – tries to **prevent length of  $\mathbf{w}^{(t)}$  from becoming too large.**

## Convergence of subgradient descent

In general, subgradient descent (with possibly non-differentiable convex objectives) **can converge relatively slowly**.

## Convergence of subgradient descent

In general, subgradient descent (with possibly non-differentiable convex objectives) **can converge relatively slowly**.

## Optimality condition for general convex objectives

$\mathbf{x}$  is a global minimizer if and only if  $\mathbf{0} \in \partial f(\mathbf{x})$ .



## Convergence of subgradient descent

In general, subgradient descent (with possibly non-differentiable convex objectives) **can converge relatively slowly**.

## Optimality condition for general convex objectives

$x$  is a global minimizer if and only if  $\mathbf{0} \in \partial f(x)$ .

Unfortunately, does not suggest a good stopping criterion for subgradient descent (since difficult to check all possible subgradients of  $f$ ).

# CONVERGENCE AND OPTIMALITY CONDITIONS

## Convergence of subgradient descent

In general, subgradient descent (with possibly non-differentiable convex objectives) **can converge relatively slowly**.

## Optimality condition for general convex objectives

$x$  is a global minimizer if and only if  $\mathbf{0} \in \partial f(x)$ .

Unfortunately, does not suggest a good stopping criterion for subgradient descent (since difficult to check all possible subgradients of  $f$ ).

But for machine learning purposes, we can use alternative criteria for stopping (e.g., hold-out error).

# CONSTRAINED CONVEX OPTIMIZATION

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^d} & f(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} \in A \end{array}$$

for convex function  $f$  and convex feasible set  $A$  (possibly a strict subset of  $\mathbb{R}^d$ ).

# CONSTRAINED CONVEX OPTIMIZATION

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^d} & f(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} \in A \end{array}$$

for convex function  $f$  and convex feasible set  $A$  (possibly a strict subset of  $\mathbb{R}^d$ ).

## Projected subgradient descent

- ▶ Start with some initial  $\mathbf{x}^{(1)} \in A$ .
- ▶ For  $t = 1, 2, \dots$  until some stopping condition is satisfied.
  - ▶ Compute *any* subgradient  $\boldsymbol{\lambda}^{(t)} \in \partial f(\mathbf{x}^{(t)})$ .
  - ▶ Update:

$$\mathbf{x}^{(t+0.5)} := \mathbf{x}^{(t)} - \eta_t \boldsymbol{\lambda}^{(t)}.$$

- ▶ Project:

$$\mathbf{x}^{(t+1)} := \text{Proj}_A(\mathbf{x}^{(t+0.5)}).$$

**Projection onto convex set  $A$ :**  $\text{Proj}_A(\mathbf{x}) := \arg \min_{\mathbf{x}' \in A} \|\mathbf{x} - \mathbf{x}'\|_2^2$   
(another convex optimization problem, but hopefully simpler objective!).

# OTHER SOLVERS

**Many** other algorithms for solving convex optimization problems

- ▶ *Newton-Raphson*: use Hessian to pick better descent directions.
- ▶ *Quasi-Newton methods* (e.g., conjugate gradient, “BFGS”, “L-BFGS”): use efficient approximations of Hessians.
- ▶ Techniques for dealing with constraints:
  - ▶ *Barrier methods*: add penalties for constraint violations, slowly relax.
  - ▶ *Primal-dual methods*: start with dual-feasible point, iteratively improve until corresponding primal point is feasible.
  - ▶ ...
- ▶ *Stochastic gradient methods*: a lot like Perceptron
- ▶ ...

**But remember**: end goal in machine learning is *not* to minimize training error (let alone training surrogate loss).