

## COMS 4771 Lecture 9

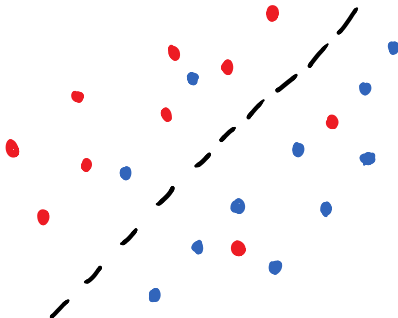
1. Soft-margin SVMs and surrogate losses
2. Convex optimization

# SOFT-MARGIN SVMs AND SURROGATE LOSSES

# NON-SEPARABLE CASES

## Non-separable cases:

No linear classifier has zero training error on  $S$ .



But if non-separability is only due to a handful of points,  
can we still find a good linear classifier?

no need to achieve  
0 training error rate

# SOFT-MARGIN SVMs (CORTES AND VAPNIK, 1995)

When  $S = ((\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)}))$  is not linearly separable, the (primal) SVM optimization problem

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y^{(i)} \left( \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - \theta \right) \geq 1 \quad \text{for } i = 1, 2, \dots, n \end{aligned}$$

**has no solution.**

# SOFT-MARGIN SVMs (CORTES AND VAPNIK, 1995)

When  $S = ((\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)}))$  is not linearly separable, the (primal) SVM optimization problem

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y^{(i)} \left( \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - \theta \right) \geq 1 \quad \text{for } i = 1, 2, \dots, n \end{aligned}$$

**has no solution.**

Introduce **slack variables**  $\xi_1, \xi_2, \dots, \xi_n \geq 0$ , and a trade-off parameter  $C > 0$ :

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}, \xi \in \mathbb{R}^n} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y^{(i)} \left( \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - \theta \right) \geq 1 - \xi_i \quad \text{for } i = 1, 2, \dots, n \\ & \xi_i \geq 0 \quad \text{for } i = 1, 2, \dots, n \end{aligned}$$

which is **always feasible**—“soft margin” SVM.

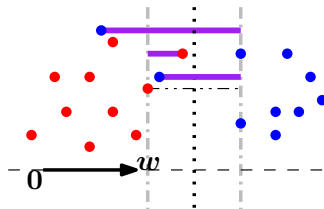
## **Winner of 2008 ACM Paris Kanellakis Award:**

*For “their revolutionary development of a highly effective algorithm known as Support Vector Machines (SVM), a set of related supervised learning methods used for data classification and regression”, which is “one of the most frequently used algorithms in machine learning, and is used in medical diagnosis, weather forecasting, and intrusion detection among many other practical applications”.*

**Other winners include:** public key cryptography, Lempel-Ziv compression, Splay Trees, interior point method for linear programming, ... and AdaBoost (discussed later in the course).

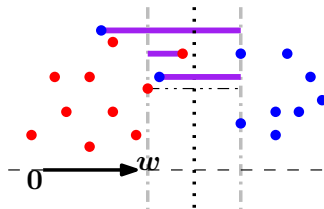
# SLACK INTERPRETATION

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^n} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y^{(i)} (\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - \theta) \geq 1 - \xi_i \quad \text{for } i = 1, 2, \dots, n \\ & \xi_i \geq 0 \quad \text{for } i = 1, 2, \dots, n \end{aligned}$$



# SLACK INTERPRETATION

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}, \xi \in \mathbb{R}^n} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y^{(i)} (\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - \theta) \geq 1 - \xi_i \quad \text{for } i = 1, 2, \dots, n \\ & \xi_i \geq 0 \quad \text{for } i = 1, 2, \dots, n \end{aligned}$$

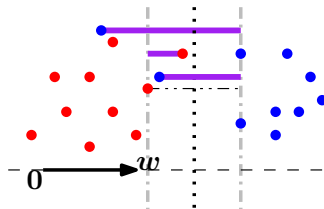


- For a given  $(\mathbf{w}, \theta)$ ,  $\xi_i / \|\mathbf{w}\|_2$  measures distance that  $\mathbf{x}^{(i)}$  must be moved so that  $y^{(i)} \langle \mathbf{w}, \mathbf{x}^{(i)} - \theta \rangle \geq 1$ .



# SLACK INTERPRETATION

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^n} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y^{(i)} (\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - \theta) \geq 1 - \xi_i \quad \text{for } i = 1, 2, \dots, n \\ & \xi_i \geq 0 \quad \text{for } i = 1, 2, \dots, n \end{aligned}$$



- ▶ For a given  $(\mathbf{w}, \theta)$ ,  $\xi_i / \|\mathbf{w}\|_2$  measures distance that  $\mathbf{x}^{(i)}$  must be moved so that  $y^{(i)} \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - \theta \geq 1$ .
- ▶  $C$  controls trade-off between slack penalties and size of margin (which is  $1 / \|\mathbf{w}\|_2$ ).

# A DIFFERENT INTERPRETATION OF SLACK

**Constraints with non-negative slack variables:** (using  $\lambda = 1/(nC)$ )

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^n} \quad & \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{n} \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y^{(i)} \left( \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - \theta \right) \geq 1 - \xi_i \quad \text{for } i = 1, 2, \dots, n \\ & \xi_i \geq 0 \quad \text{for } i = 1, 2, \dots, n \end{aligned}$$

# A DIFFERENT INTERPRETATION OF SLACK

**Constraints with non-negative slack variables:** (using  $\lambda = 1/(nC)$ )

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^n} \quad & \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{n} \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y^{(i)} \left( \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - \theta \right) \geq 1 - \xi_i \quad \text{for } i = 1, 2, \dots, n \\ & \xi_i \geq 0 \quad \text{for } i = 1, 2, \dots, n \end{aligned}$$

**Equivalent unconstrained form:**

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

*Notation:*  $[a]_+ := \max\{0, a\}$ .

# A DIFFERENT INTERPRETATION OF SLACK

**Constraints with non-negative slack variables:** (using  $\lambda = 1/(nC)$ )

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^n} \quad & \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{n} \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y^{(i)} \left( \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - \theta \right) \geq 1 - \xi_i \quad \text{for } i = 1, 2, \dots, n \\ & \xi_i \geq 0 \quad \text{for } i = 1, 2, \dots, n \end{aligned}$$

**Equivalent unconstrained form:**

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

*Notation:*  $[a]_+ := \max\{0, a\}$ .

The **hinge loss** of a linear classifier  $f_{\mathbf{w}, \theta}$  on an example  $(\mathbf{x}, y)$  is defined to be

$$\text{HL}(\mathbf{w}, \theta; \mathbf{x}, y) := \left[ 1 - y(\langle \mathbf{w}, \mathbf{x} \rangle - \theta) \right]_+.$$

# A DIFFERENT INTERPRETATION OF SLACK

**Constraints with non-negative slack variables:** (using  $\lambda = 1/(nC)$ )

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^n} \quad & \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{n} \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y^{(i)} \left( \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - \theta \right) \geq 1 - \xi_i \quad \text{for } i = 1, 2, \dots, n \\ & \xi_i \geq 0 \quad \text{for } i = 1, 2, \dots, n \end{aligned}$$

**Equivalent unconstrained form:**

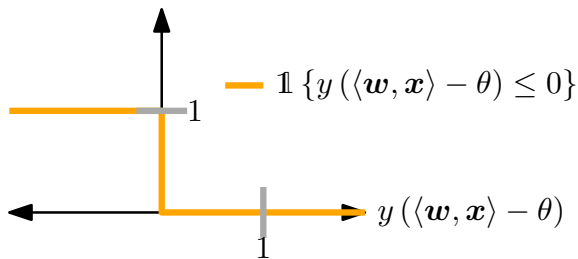
$$\min_{\mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}} \quad \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{n} \sum_{i=1}^n \text{HL}(\mathbf{w}, \theta; \mathbf{x}^{(i)}, y^{(i)})$$

*Notation:*  $[a]_+ := \max\{0, a\}$ . trade off : larger margin ? lower hinge

The **hinge loss** of a linear classifier  $f_{\mathbf{w}, \theta}$  on an example  $(\mathbf{x}, y)$  is defined to be

$$\text{HL}(\mathbf{w}, \theta; \mathbf{x}, y) := \left[ 1 - y(\langle \mathbf{w}, \mathbf{x} \rangle - \theta) \right]_+.$$

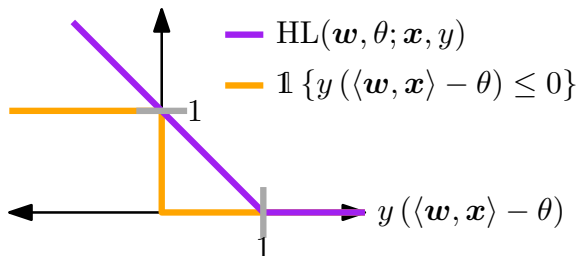
# ZERO-ONE LOSS VS. HINGE LOSS



**Zero-one loss:** count if  $f_{\mathbf{w},\theta}(\mathbf{x}) \neq y$ .

$$\mathbb{1}\{y(\langle \mathbf{w}, \mathbf{x} \rangle - \theta) \leq 0\} \leq \left[1 - y(\langle \mathbf{w}, \mathbf{x} \rangle - \theta)\right]_+ = \text{HL}(\mathbf{w}, \theta; \mathbf{x}, y).$$

# ZERO-ONE LOSS VS. HINGE LOSS



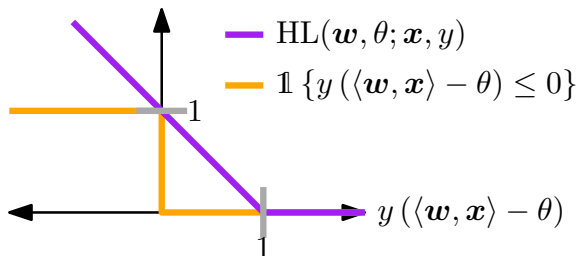
**Hinge loss:** an upper-bound on zero-one loss.

$$\mathbb{1} \{y(\langle \mathbf{w}, \mathbf{x} \rangle - \theta) \leq 0\} \leq \left[ 1 - y(\langle \mathbf{w}, \mathbf{x} \rangle - \theta) \right]_+ = HL(\mathbf{w}, \theta; \mathbf{x}, y).$$

wrong predication!

When this smaller than one, a penalty would be incurred!

# ZERO-ONE LOSS VS. HINGE LOSS



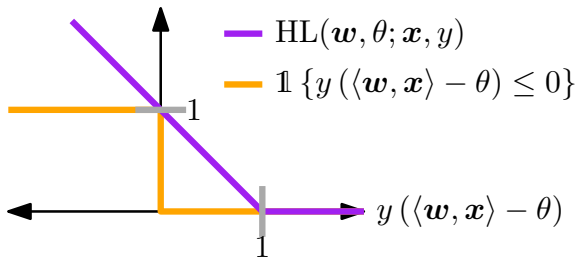
**Hinge loss:** an upper-bound on **zero-one loss**.

$$\mathbb{1} \{y (\langle \mathbf{w}, \mathbf{x} \rangle - \theta) \leq 0\} \leq \left[ 1 - y (\langle \mathbf{w}, \mathbf{x} \rangle - \theta) \right]_+ = \text{HL}(\mathbf{w}, \theta; \mathbf{x}, y).$$

**Soft-margin SVM** minimizes an upper-bound on the training error, plus a term that favors large margins.



# ZERO-ONE LOSS VS. HINGE LOSS



**Hinge loss:** an upper-bound on **zero-one loss**.

$$\mathbb{1} \{y (\langle \mathbf{w}, \mathbf{x} \rangle - \theta) \leq 0\} \leq \left[ 1 - y (\langle \mathbf{w}, \mathbf{x} \rangle - \theta) \right]_+ = \text{HL}(\mathbf{w}, \theta; \mathbf{x}, y).$$

**Soft-margin SVM** minimizes an upper-bound on the training error, plus a term that favors large margins.

This is **computationally tractable** (unlike minimizing training error) because the hinge loss is a **convex function** of  $(\mathbf{w}, \theta)$ , and so is  $\frac{\lambda}{2} \|\mathbf{w}\|_2^2$ .

# GENERAL FORM

**Empirical risk minimization** (i.e., minimize training error):

$$\min_{\mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}} \quad \frac{1}{n} \sum_{i=1}^n \mathbb{1} \left\{ y^{(i)} \left( \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - \theta \right) \leq 0 \right\}$$

**Soft-margin SVM:**

$$\min_{\mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}} \quad \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{n} \sum_{i=1}^n \text{HL}(\mathbf{w}, \theta; \mathbf{x}^{(i)}, y^{(i)})$$

# GENERAL FORM

**Empirical risk minimization** (i.e., minimize training error):

$$\min_{\mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}} \quad \frac{1}{n} \sum_{i=1}^n \mathbb{1} \left\{ y^{(i)} \left( \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - \theta \right) \leq 0 \right\}$$

**Soft-margin SVM:**

$$\min_{\mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}} \quad \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{n} \sum_{i=1}^n \text{HL}(\mathbf{w}, \theta; \mathbf{x}^{(i)}, y^{(i)})$$

**Generic learning objective:**

$$\min_{\mathbf{w} \in \mathbb{R}^d} \quad R(\mathbf{w}) + \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}; \mathbf{x}^{(i)}, y^{(i)})$$

- ▶ **Regularization:** encodes “learning bias” (e.g., preference for large margins), sometimes promotes stability.
- ▶ **Data fitting/empirical loss:** how poorly does the classifier “fit” the data.

# SIMILARITY TO MAXIMUM LIKELIHOOD

Generic learning objective:

$$\min_{\mathbf{w} \in \mathbb{R}^d} R(\mathbf{w}) + \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}; \mathbf{x}^{(i)}, y^{(i)})$$

**Maximum likelihood estimation** for a parameteric family  $\{p(\cdot; \boldsymbol{\theta}) : \boldsymbol{\theta} \in \mathcal{T}\}$  (assuming data  $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(n)}$  are iid):

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} -\frac{1}{n} \sum_{i=1}^n \log p(\mathbf{z}^{(i)}; \boldsymbol{\theta})$$

(i.e., minimize  $1/n$  times negative log-likelihood of parameter  $\boldsymbol{\theta}$ ).

# SIMILARITY TO MAXIMUM LIKELIHOOD

**Generic learning objective:**

$$\min_{\mathbf{w} \in \mathbb{R}^d} R(\mathbf{w}) + \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}; \mathbf{x}^{(i)}, y^{(i)})$$

**Maximum likelihood estimation** for a parameteric family  $\{p(\cdot; \boldsymbol{\theta}) : \boldsymbol{\theta} \in \mathcal{T}\}$  (assuming data  $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(n)}$  are iid):

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} -\frac{1}{n} \sum_{i=1}^n \log p(\mathbf{z}^{(i)}; \boldsymbol{\theta})$$

(i.e., minimize  $1/n$  times negative log-likelihood of parameter  $\boldsymbol{\theta}$ ).

Sometimes generic learning objective is called “generalized (and penalized, because of  $R(\mathbf{w})$ ) maximum likelihood estimation”.

# LEARNING VIA OPTIMIZATION

- ▶ Many different choices for regularization and loss.
  - ▶  $R(\mathbf{w}) \propto \|\mathbf{w}\|_2^2$ : encourage large margins
  - ▶  $R(\mathbf{w}) \propto \|\mathbf{w}\|_1$ : encourage  $\mathbf{w}$  to be sparse
  - ▶  $R(\mathbf{w}) \propto \sum_{i=1}^n w_i \ln w_i$ : “maximum entropy” interpretation
  - ▶  $\ell(\mathbf{w}, \theta; \mathbf{x}, y) = [\langle \mathbf{w}, \mathbf{x} \rangle - \theta - y]_+^2$
  - ▶  $\ell(\mathbf{w}, \theta; \mathbf{x}, y) = \ln(1 + \exp(-y(\langle \mathbf{w}, \mathbf{x} \rangle - \theta)))$  (“logistic regression”)
  - ▶ ...
  - ▶ Also used beyond classification (e.g., regression, parameter estimation, clustering)

# LEARNING VIA OPTIMIZATION

- ▶ Many different choices for regularization and loss.
  - ▶  $R(\mathbf{w}) \propto \|\mathbf{w}\|_2^2$ : encourage large margins
  - ▶  $R(\mathbf{w}) \propto \|\mathbf{w}\|_1$ : encourage  $\mathbf{w}$  to be sparse
  - ▶  $R(\mathbf{w}) \propto \sum_{i=1}^n w_i \ln w_i$ : “maximum entropy” interpretation
  - ▶  $\ell(\mathbf{w}, \theta; \mathbf{x}, y) = [\langle \mathbf{w}, \mathbf{x} \rangle - \theta - y]_+^2$
  - ▶  $\ell(\mathbf{w}, \theta; \mathbf{x}, y) = \ln(1 + \exp(-y(\langle \mathbf{w}, \mathbf{x} \rangle - \theta)))$  (“logistic regression”)
  - ▶ ...
  - ▶ Also used beyond classification (e.g., regression, parameter estimation, clustering)
- ▶ Often want  $\ell$  be an upper-bound on zero-one loss—i.e., a **surrogate loss**.

# LEARNING VIA OPTIMIZATION

- ▶ Many different choices for regularization and loss.
  - ▶  $R(\mathbf{w}) \propto \|\mathbf{w}\|_2^2$ : encourage large margins
  - ▶  $R(\mathbf{w}) \propto \|\mathbf{w}\|_1$ : encourage  $\mathbf{w}$  to be sparse
  - ▶  $R(\mathbf{w}) \propto \sum_{i=1}^n w_i \ln w_i$ : “maximum entropy” interpretation
  - ▶  $\ell(\mathbf{w}, \theta; \mathbf{x}, y) = [\langle \mathbf{w}, \mathbf{x} \rangle - \theta - y]_+^2$
  - ▶  $\ell(\mathbf{w}, \theta; \mathbf{x}, y) = \ln(1 + \exp(-y(\langle \mathbf{w}, \mathbf{x} \rangle - \theta)))$  (“logistic regression”)
  - ▶ ...
  - ▶ Also used beyond classification (e.g., regression, parameter estimation, clustering)
- ▶ Often want  $\ell$  be an upper-bound on zero-one loss—i.e., a **surrogate loss**.
- ▶ Trade-off parameter  $\lambda$ : usually determine using hold out error or cross validation error.



# LEARNING VIA OPTIMIZATION

- ▶ Many different choices for regularization and loss.
  - ▶  $R(\mathbf{w}) \propto \|\mathbf{w}\|_2^2$ : encourage large margins
  - ▶  $R(\mathbf{w}) \propto \|\mathbf{w}\|_1$ : encourage  $\mathbf{w}$  to be sparse
  - ▶  $R(\mathbf{w}) \propto \sum_{i=1}^n w_i \ln w_i$ : “maximum entropy” interpretation
  - ▶  $\ell(\mathbf{w}, \theta; \mathbf{x}, y) = [\langle \mathbf{w}, \mathbf{x} \rangle - \theta - y]_+^2$
  - ▶  $\ell(\mathbf{w}, \theta; \mathbf{x}, y) = \ln(1 + \exp(-y(\langle \mathbf{w}, \mathbf{x} \rangle - \theta)))$  (“logistic regression”)
  - ▶ ...
  - ▶ Also used beyond classification (e.g., regression, parameter estimation, clustering)
- ▶ Often want  $\ell$  be an upper-bound on zero-one loss—i.e., a **surrogate loss**.
- ▶ Trade-off parameter  $\lambda$ : usually determine using hold out error or cross validation error.
- ▶ Computationally easier when overall objective function is **convex**: possible to efficiently find global minimizer in polynomial time.

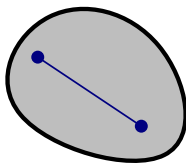
# LEARNING VIA OPTIMIZATION

- ▶ Many different choices for regularization and loss.
  - ▶  $R(\mathbf{w}) \propto \|\mathbf{w}\|_2^2$ : encourage large margins
  - ▶  $R(\mathbf{w}) \propto \|\mathbf{w}\|_1$ : encourage  $\mathbf{w}$  to be sparse **absolute value**
  - ▶  $R(\mathbf{w}) \propto \sum_{i=1}^n w_i \ln w_i$ : “maximum entropy” interpretation
  - ▶  $\ell(\mathbf{w}, \theta; \mathbf{x}, y) = [\langle \mathbf{w}, \mathbf{x} \rangle - \theta - y]_+^2$
  - ▶  $\ell(\mathbf{w}, \theta; \mathbf{x}, y) = \ln(1 + \exp(-y(\langle \mathbf{w}, \mathbf{x} \rangle - \theta)))$  (“logistic regression”)
  - ▶ ...
  - ▶ Also used beyond classification (e.g., regression, parameter estimation, clustering)
- ▶ Often want  $\ell$  be an upper-bound on zero-one loss—i.e., a **surrogate loss**.
- ▶ Trade-off parameter  $\lambda$ : usually determine using hold out error or cross validation error.
- ▶ Computationally easier when overall objective function is **convex**: possible to efficiently find global minimizer in polynomial time.
- ▶ **Next**: techniques for analyzing and solving these optimization problems.

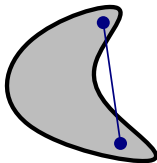
# INTRODUCTION TO CONVEXITY

# CONVEX SETS

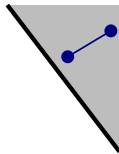
We say a set  $A$  is **convex** if, for every pair of points  $\{x, x'\}$  in the set  $A$ , the line segment between the points  $x$  and  $x'$  is also contained in the set  $A$ .



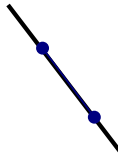
convex



not convex



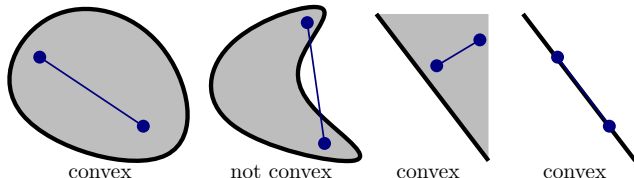
convex



convex

# CONVEX SETS

We say a set  $A$  is **convex** if, for every pair of points  $\{x, x'\}$  in the set  $A$ , the line segment between the points  $x$  and  $x'$  is also contained in the set  $A$ .



## Examples:

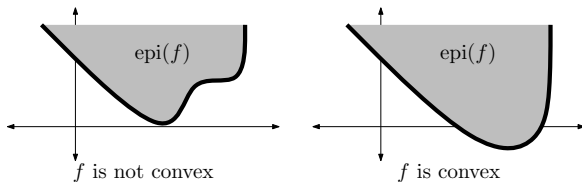
- ▶ All of  $\mathbb{R}^d$ .
- ▶ Empty set.
- ▶ Affine hyperplanes.
- ▶ Half-spaces:  $\{a \in \mathbb{R}^d : \langle a, x \rangle - b \leq 0\}$ .
- ▶ Intersections of convex sets.
- ▶ Convex hulls of points.

# CONVEX FUNCTIONS

For any function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ , the **epigraph** of  $f$ , denoted  $\text{epi}(f)$ , is the set

$$\text{epi}(f) := \{(\mathbf{x}, b) \in \mathbb{R}^{d+1} : f(\mathbf{x}) \leq b\}.$$

We say a function  $f$  is **convex** if  $\text{epi}(f)$  is a convex set in  $\mathbb{R}^{d+1}$ .

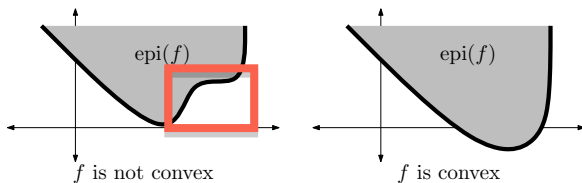


# CONVEX FUNCTIONS

For any function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ , the **epigraph** of  $f$ , denoted  $\text{epi}(f)$ , is the set

$$\text{epi}(f) := \{(\mathbf{x}, b) \in \mathbb{R}^{d+1} : f(\mathbf{x}) \leq b\}.$$

We say a function  $f$  is **convex** if  $\text{epi}(f)$  is a convex set in  $\mathbb{R}^{d+1}$ .



## Examples:

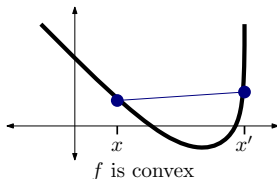
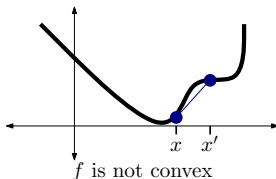
- ▶  $f(x) = c^x$  for any  $c > 0$  (on  $\mathbb{R}$ )
- ▶  $f(x) = |x|^c$  for any  $c \geq 1$  (on  $\mathbb{R}$ )
- ▶  $f(\mathbf{x}) = c$  for any constant  $c \in \mathbb{R}$ .
- ▶  $f(\mathbf{x}) = \langle \mathbf{a}, \mathbf{x} \rangle$  for any  $\mathbf{a} \in \mathbb{R}^d$ .
- ▶  $f(\mathbf{x}) = \|\mathbf{x}\|_p$  for any  $1 \leq p \leq \infty$ .
- ▶  $f(\mathbf{x}) = \langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle$  for symmetric positive semidefinite  $\mathbf{A}$ .

# JENSEN'S INEQUALITY

## Equivalent definition of convex functions

A function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex if and only if, for any  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$  and  $\alpha \in [0, 1]$ ,

$$f((1 - \alpha)\mathbf{x} + \alpha\mathbf{x}') \leq (1 - \alpha) \cdot f(\mathbf{x}) + \alpha \cdot f(\mathbf{x}').$$



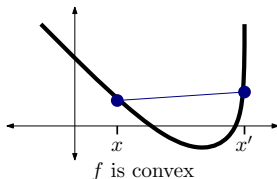
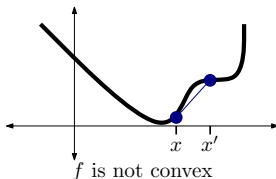


# JENSEN'S INEQUALITY

## Equivalent definition of convex functions

A function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex if and only if, for any  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$  and  $\alpha \in [0, 1]$ ,

$$f((1 - \alpha)\mathbf{x} + \alpha\mathbf{x}') \leq (1 - \alpha) \cdot f(\mathbf{x}) + \alpha \cdot f(\mathbf{x}').$$



Covex points

## Jensen's inequality

If  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex, then for any  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$  and  $\alpha_1, \alpha_2, \dots, \alpha_n \in [0, 1]$  such that  $\sum_{i=1}^n \alpha_i = 1$ ,

$$f\left(\sum_{i=1}^n \alpha_i \mathbf{x}_i\right) \leq \sum_{i=1}^n \alpha_i \cdot f(\mathbf{x}_i).$$

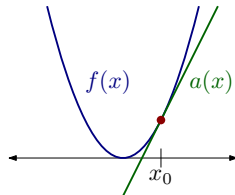
# CONVEXITY OF DIFFERENTIABLE FUNCTIONS

## Differentiable functions

If  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is differentiable, then  $f$  is convex if and only if

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0 \rangle$$

for all  $\mathbf{x}, \mathbf{x}_0 \in \mathbb{R}^d$ .



$$a(x) = f(x_0) + f'(x_0)(x - x_0)$$

Text

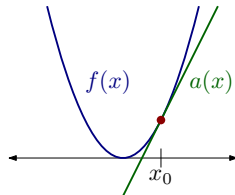
# CONVEXITY OF DIFFERENTIABLE FUNCTIONS

## Differentiable functions

If  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is differentiable, then  $f$  is convex if and only if

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0 \rangle$$

for all  $\mathbf{x}, \mathbf{x}_0 \in \mathbb{R}^d$ .



$$a(x) = f(x_0) + f'(x_0)(x - x_0)$$

## Twice-differentiable functions

If  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is twice-differentiable, then  $f$  is convex if and only if

$$\nabla^2 f(\mathbf{x}) \succeq \mathbf{0}$$

for all  $\mathbf{x} \in \mathbb{R}^d$  (i.e., the Hessian, or matrix of second-derivatives, is positive semidefinite for all  $\mathbf{x}$ ).

# MORE CONVEX FUNCTIONS

## Building new convex functions from old ones

- ▶  $f(\mathbf{x}) = c \cdot h(\mathbf{x}) + g(\mathbf{x})$  for convex functions  $h, g$  and scalar  $c \geq 0$ .

Example:  $f(\mathbf{x}) = -2\langle \mathbf{a}, \mathbf{x} \rangle + 1$ .

- ▶  $f(\mathbf{x}) = \max\{h(\mathbf{x}), g(\mathbf{x})\}$  for convex functions  $h, g$ .

Example:  $f(\mathbf{x}) = \max\{0, 1 - \langle \mathbf{a}, \mathbf{x} \rangle\}$ .

- ▶  $f(\mathbf{x}) = g(\mathbf{A}\mathbf{x} + \mathbf{b})$  for any  $\mathbf{A} \in \mathbb{R}^{m \times d}$  and  $\mathbf{b} \in \mathbb{R}^m$  and convex function  $g: \mathbb{R}^m \rightarrow \mathbb{R}$ .

Examples:  $f(\mathbf{x}) = \exp(\langle \mathbf{a}, \mathbf{x} \rangle)$ ,  $f(\mathbf{x}) = \langle \mathbf{a}, \mathbf{x} \rangle^2$ .

- ▶  $f(\mathbf{x}) = g(h_1(\mathbf{x}), \dots, h_m(\mathbf{x}))$  for any convex  $h_1, \dots, h_m$  and convex  $g: \mathbb{R}^m \rightarrow \mathbb{R}$  such that  $g$  is non-decreasing in each argument.

Example:  $f(\mathbf{x}) = \exp(\langle \mathbf{a}, \mathbf{x} \rangle^2)$  (but not  $\exp(-\langle \mathbf{a}, \mathbf{x} \rangle^2)$ ).

- ▶ Many other composition rules.

# HOW TO TEST FOR CONVEXITY

- ▶ First principles (via definitions).
- ▶ First- or second-derivative tests (assuming derivatives exist).
- ▶ Valid transformation of existing convex functions.

# CONVEX OPTIMIZATION PROBLEMS

# OPTIMIZATION PROBLEMS

A typical optimization problem is written as

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

where  $f_0: \mathbb{R}^d \rightarrow \mathbb{R}$  is the **objective function** and  $f_1, f_2, \dots, f_n: \mathbb{R}^d \rightarrow \mathbb{R}$  are the **constraint functions**.

# OPTIMIZATION PROBLEMS

A typical optimization problem is written as

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

where  $f_0: \mathbb{R}^d \rightarrow \mathbb{R}$  is the **objective function** and  $f_1, f_2, \dots, f_n: \mathbb{R}^d \rightarrow \mathbb{R}$  are the **constraint functions**.

The inequalities  $f_i(\mathbf{x}) \leq 0$  are **constraints**.



# OPTIMIZATION PROBLEMS

A typical optimization problem is written as

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

where  $f_0: \mathbb{R}^d \rightarrow \mathbb{R}$  is the **objective function** and  $f_1, f_2, \dots, f_n: \mathbb{R}^d \rightarrow \mathbb{R}$  are the **constraint functions**.

The inequalities  $f_i(\mathbf{x}) \leq 0$  are **constraints**.

The set  $A := \{\mathbf{x} \in \mathbb{R}^d : f_i(\mathbf{x}) \leq 0 \text{ for all } i = 1, 2, \dots, n\}$  is the **feasible set**.

# OPTIMIZATION PROBLEMS

A typical optimization problem is written as

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

where  $f_0: \mathbb{R}^d \rightarrow \mathbb{R}$  is the **objective function** and  $f_1, f_2, \dots, f_n: \mathbb{R}^d \rightarrow \mathbb{R}$  are the **constraint functions**.

The inequalities  $f_i(\mathbf{x}) \leq 0$  are **constraints**.

The set  $A := \{\mathbf{x} \in \mathbb{R}^d : f_i(\mathbf{x}) \leq 0 \text{ for all } i = 1, 2, \dots, n\}$  is the **feasible set**.

**The goal:** Find  $\mathbf{x} \in A$  so that  $f_0(\mathbf{x})$  is as small as possible.

# OPTIMIZATION PROBLEMS

A typical optimization problem is written as

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

where  $f_0: \mathbb{R}^d \rightarrow \mathbb{R}$  is the **objective function** and  $f_1, f_2, \dots, f_n: \mathbb{R}^d \rightarrow \mathbb{R}$  are the **constraint functions**.

The inequalities  $f_i(\mathbf{x}) \leq 0$  are **constraints**.

The set  $A := \{\mathbf{x} \in \mathbb{R}^d : f_i(\mathbf{x}) \leq 0 \text{ for all } i = 1, 2, \dots, n\}$  is the **feasible set**.

**The goal:** Find  $\mathbf{x} \in A$  so that  $f_0(\mathbf{x})$  is as small as possible.

The **(optimal) value** of the optimization problem is the smallest such value of  $f_0(\mathbf{x})$  achieved by a feasible point  $\mathbf{x} \in A$ .

# OPTIMIZATION PROBLEMS

A typical optimization problem is written as

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

where  $f_0: \mathbb{R}^d \rightarrow \mathbb{R}$  is the **objective function** and  $f_1, f_2, \dots, f_n: \mathbb{R}^d \rightarrow \mathbb{R}$  are the **constraint functions**.

The inequalities  $f_i(\mathbf{x}) \leq 0$  are **constraints**.

The set  $A := \{\mathbf{x} \in \mathbb{R}^d : f_i(\mathbf{x}) \leq 0 \text{ for all } i = 1, 2, \dots, n\}$  is the **feasible set**.

**The goal:** Find  $\mathbf{x} \in A$  so that  $f_0(\mathbf{x})$  is as small as possible.

The **(optimal) value** of the optimization problem is the smallest such value of  $f_0(\mathbf{x})$  achieved by a feasible point  $\mathbf{x} \in A$ .

A point  $\mathbf{x} \in A$  achieving the optimal value is a **(global) minimizer** of the problem.

# CONVEX OPTIMIZATION PROBLEMS

Standard form of a **convex optimization problem**:

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

where  $f_0, f_1, \dots, f_n: \mathbb{R}^d \rightarrow \mathbb{R}$  are *convex functions*.

# CONVEX OPTIMIZATION PROBLEMS

Standard form of a **convex optimization problem**:

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

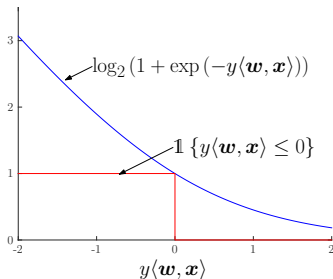
where  $f_0, f_1, \dots, f_n: \mathbb{R}^d \rightarrow \mathbb{R}$  are *convex functions*.

**Fact:** the feasible set  $A := \{\mathbf{x} \in \mathbb{R}^d : f_i(\mathbf{x}) \leq 0 \text{ for all } i = 1, 2, \dots, n\}$  is a convex set.

# EXAMPLE: NON-NEGATIVE LOGISTIC REGRESSION

Training data  $S := ((\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)}))$  with  $(\mathbf{x}^{(i)}, y^{(i)}) \in \mathbb{R}^d \times \{\pm 1\}$ .

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d} \quad & f(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \ln \left( 1 + \exp \left( -y^{(i)} \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle \right) \right) \\ \text{s.t.} \quad & w_i \geq 0 \quad \text{for } i = 1, 2, \dots, d. \end{aligned}$$



Like SVM, but using a different surrogate loss, no regularization term, and only want non-negative weights.

## EXAMPLE: NON-NEGATIVE LOGISTIC REGRESSION

**Objective function:**

$$f_0(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \ln \left( 1 + \exp \left( -y^{(i)} \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle \right) \right).$$



# EXAMPLE: NON-NEGATIVE LOGISTIC REGRESSION

**Objective function:**

$$f_0(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \ln \left( 1 + \exp \left( -y^{(i)} \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle \right) \right).$$

**Second-derivative test of convexity:** Is  $\nabla^2 f_0(\mathbf{w})$  positive semidefinite?

$$\nabla^2 f_0(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + e^{\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle}} \cdot \frac{e^{\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle}}{1 + e^{\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle}} \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^\top.$$

# EXAMPLE: NON-NEGATIVE LOGISTIC REGRESSION

**Objective function:**

$$f_0(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \ln \left( 1 + \exp \left( -y^{(i)} \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle \right) \right).$$

**Second-derivative test of convexity:** Is  $\nabla^2 f_0(\mathbf{w})$  positive semidefinite?

$$\nabla^2 f_0(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + e^{\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle}} \cdot \frac{e^{\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle}}{1 + e^{\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle}} \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^\top.$$

Yes: for any  $\mathbf{v} \in \mathbb{R}^d$ ,

$$\langle \nabla^2 f_0(\mathbf{w}) \mathbf{v}, \mathbf{v} \rangle = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + e^{\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle}} \cdot \frac{e^{\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle}}{1 + e^{\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle}} \langle \mathbf{v}, \mathbf{x}^{(i)} \rangle^2 \geq 0.$$

# EXAMPLE: NON-NEGATIVE LOGISTIC REGRESSION

**Objective function:**

$$f_0(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \ln \left( 1 + \exp \left( -y^{(i)} \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle \right) \right).$$

**Second-derivative test of convexity:** Is  $\nabla^2 f_0(\mathbf{w})$  positive semidefinite?

$$\nabla^2 f_0(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + e^{\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle}} \cdot \frac{e^{\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle}}{1 + e^{\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle}} \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^\top.$$

Yes: for any  $\mathbf{v} \in \mathbb{R}^d$ ,

$$\langle \nabla^2 f_0(\mathbf{w}) \mathbf{v}, \mathbf{v} \rangle = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + e^{\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle}} \cdot \frac{e^{\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle}}{1 + e^{\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle}} \langle \mathbf{v}, \mathbf{x}^{(i)} \rangle^2 \geq 0.$$

**Conclusion:** objective function  $f_0(\mathbf{w})$  is convex.

# EXAMPLE: NON-NEGATIVE LOGISTIC REGRESSION

**Constraints:**

$$w_i \geq 0 \quad \text{for } i = 1, 2, \dots, d.$$

# EXAMPLE: NON-NEGATIVE LOGISTIC REGRESSION

**Constraints:**

$$w_i \geq 0 \quad \text{for } i = 1, 2, \dots, d.$$

**In standard form:**  $f_i(\mathbf{w}) \leq 0$ , where

$$f_i(\mathbf{w}) := -w_i \quad \text{for } i = 1, 2, \dots, d.$$

# EXAMPLE: NON-NEGATIVE LOGISTIC REGRESSION

**Constraints:**

$$w_i \geq 0 \quad \text{for } i = 1, 2, \dots, d.$$

**In standard form:**  $f_i(\mathbf{w}) \leq 0$ , where

$$f_i(\mathbf{w}) := -w_i \quad \text{for } i = 1, 2, \dots, d.$$

Each  $f_i$  is a linear function of  $\mathbf{w}$ , and therefore is convex.

# EXAMPLE: NON-NEGATIVE LOGISTIC REGRESSION

**Constraints:**

$$w_i \geq 0 \quad \text{for } i = 1, 2, \dots, d.$$

**In standard form:**  $f_i(\mathbf{w}) \leq 0$ , where

$$f_i(\mathbf{w}) := -w_i \quad \text{for } i = 1, 2, \dots, d.$$

Each  $f_i$  is a linear function of  $\mathbf{w}$ , and therefore is convex.

**Conclusion:** constraint functions  $f_i(\mathbf{w})$  are convex.

# EXAMPLE: NON-NEGATIVE LOGISTIC REGRESSION

**Constraints:**

$$w_i \geq 0 \quad \text{for } i = 1, 2, \dots, d.$$

**In standard form:**  $f_i(\mathbf{w}) \leq 0$ , where

$$f_i(\mathbf{w}) := -w_i \quad \text{for } i = 1, 2, \dots, d.$$

Each  $f_i$  is a linear function of  $\mathbf{w}$ , and therefore is convex.

**Conclusion:** constraint functions  $f_i(\mathbf{w})$  are convex.

**Overall problem is a convex optimization problem.**





# LOCAL MINIMIZERS

Consider an optimization problem (not necessarily convex):

$$\begin{array}{ll}\min_{\boldsymbol{x} \in \mathbb{R}^d} & f_0(\boldsymbol{x}) \\ \text{s.t.} & \boldsymbol{x} \in A.\end{array}$$

# LOCAL MINIMIZERS

Consider an optimization problem (not necessarily convex):

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^d} & f_0(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} \in A. \end{array}$$

We say  $\tilde{\mathbf{x}} \in A$  is a **local minimizer** if there is an open ball

$$U := \left\{ \mathbf{x} \in \mathbb{R}^d : \|\mathbf{x} - \tilde{\mathbf{x}}\|_2 < r \right\}$$

of positive radius  $r > 0$  such that  $\tilde{\mathbf{x}}$  is a global minimizer for

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^d} & f_0(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} \in A \cap U. \end{array}$$

# LOCAL MINIMIZERS

Consider an optimization problem (not necessarily convex):

$$\begin{array}{ll}\min_{\mathbf{x} \in \mathbb{R}^d} & f_0(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} \in A.\end{array}$$

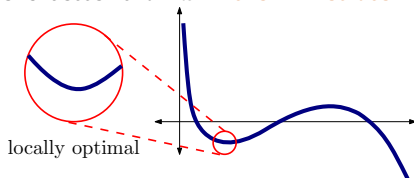
We say  $\tilde{\mathbf{x}} \in A$  is a **local minimizer** if there is an open ball

$$U := \left\{ \mathbf{x} \in \mathbb{R}^d : \|\mathbf{x} - \tilde{\mathbf{x}}\|_2 < r \right\}$$

of positive radius  $r > 0$  such that  $\tilde{\mathbf{x}}$  is a global minimizer for

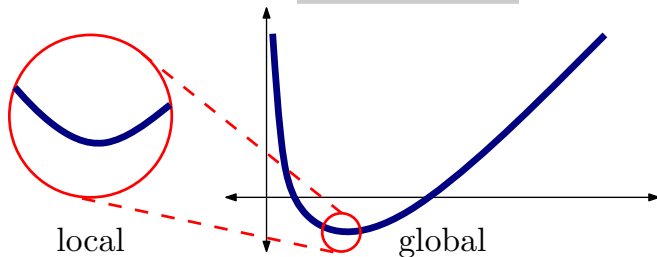
$$\begin{array}{ll}\min_{\mathbf{x} \in \mathbb{R}^d} & f_0(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} \in A \cap U.\end{array}$$

Nothing looks better than  $\tilde{\mathbf{x}}$  **in the immediate vicinity** of  $\tilde{\mathbf{x}}$ .



# LOCAL-TO-GLOBAL PHENOMENON

If the optimization problem is **convex**, and  $\tilde{x} \in A$  is a **local minimizer**,  
then it is also **a global minimizer**.



# SOLVING CONVEX OPTIMIZATION PROBLEMS

## Unconstrained convex optimization problems

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

( $f$  is the convex objective function).

# UNCONSTRAINED CONVEX OPTIMIZATION

## Unconstrained convex optimization problems

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

( $f$  is the convex objective function).

Optimality condition for differentiable convex objectives

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

guarantee convex

# LOCAL OPTIMIZATION FOR CONVEX OBJECTIVES

## Local optimization

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



# LOCAL OPTIMIZATION FOR CONVEX OBJECTIVES

## Local optimization

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

# LOCAL OPTIMIZATION FOR CONVEX OBJECTIVES

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

# LOCAL OPTIMIZATION FOR CONVEX OBJECTIVES

## 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

## 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

$$f(\mathbf{x} + \boldsymbol{\delta}) \geq f(\mathbf{x}). \quad \text{Oops!}$$

# LOCAL OPTIMIZATION FOR CONVEX OBJECTIVES

## 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

$$f(\mathbf{x} + \boldsymbol{\delta}) \geq f(\mathbf{x}). \quad \text{Oops!}$$

**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

## 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.



# 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  $\eta := \frac{1}{2}\eta$ .

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

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

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

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

# 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  $\eta := \frac{1}{2}\eta$ .

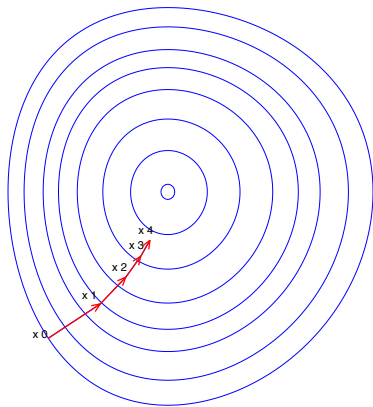
**Main idea:**  $f(\mathbf{x} - \eta\boldsymbol{\lambda}) \approx f(\mathbf{x}) - \eta\|\boldsymbol{\lambda}\|_2^2$ , so optimistically hope to decrease value by about  $\eta\|\boldsymbol{\lambda}\|_2^2$  when  $\eta$  is small enough.

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

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

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

# ILLUSTRATION OF GRADIENT DESCENT

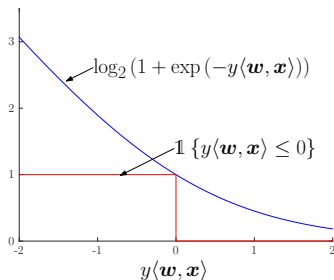


If  $f$  is convex (and satisfies some other smoothness and curvature conditions), then  $f(x^{(t)})$  converges to the optimal value at a geometric rate.

# EXAMPLE: (UNCONSTRAINED) LOGISTIC REGRESSION

Training data  $S := ((\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)}))$  with  $(\mathbf{x}^{(i)}, y^{(i)}) \in \mathbb{R}^d \times \{\pm 1\}$ .

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

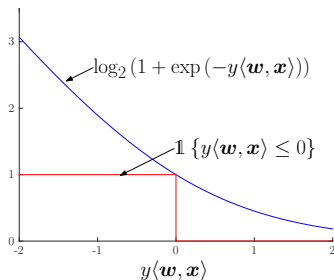


We've already established that objective  $f(\mathbf{w})$  is convex.

# EXAMPLE: (UNCONSTRAINED) LOGISTIC REGRESSION

Training data  $S := ((\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)}))$  with  $(\mathbf{x}^{(i)}, y^{(i)}) \in \mathbb{R}^d \times \{\pm 1\}$ .

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



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  $\mathbf{w}$ :

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



# EXAMPLE: (UNCONSTRAINED) LOGISTIC REGRESSION

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

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

Gradient descent algorithm:

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

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

# 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 ( $\|\nabla f(\boldsymbol{x})\|_2 \leq \epsilon$  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 ( $\|\nabla f(\boldsymbol{x})\|_2 \leq \epsilon$  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., based on hold-out error).**

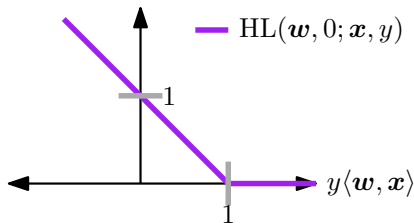
# NON-DIFFERENTIABILITY

## Non-differentiable convex objectives

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

Example: hinge loss

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



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

# COPING WITH NON-DIFFERENTIABILITY

## 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

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

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

# COPING WITH NON-DIFFERENTIABILITY

## 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 . . .

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



# COPING WITH NON-DIFFERENTIABILITY

## 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:**

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

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

# COPING WITH NON-DIFFERENTIABILITY

## 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:**

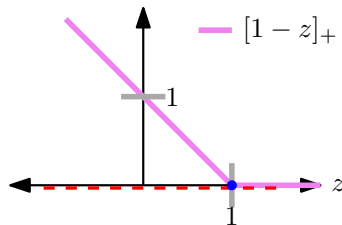
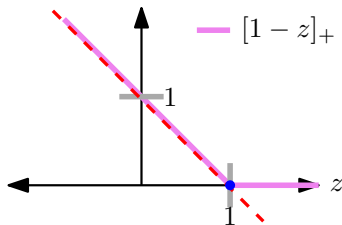
$$f(x) \geq f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle \quad \forall x \in \mathbb{R}^d.$$

**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

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



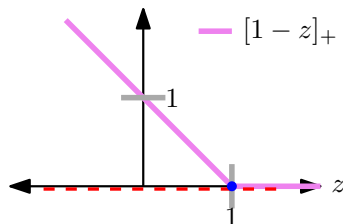
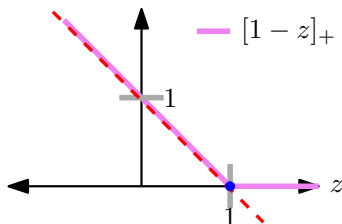
Two subgradients of  $f$  at  $z = 1$ :  $-1$  and  $0$ .

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

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

# EXAMPLE: SUBGRADIENT OF HINGE LOSS

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**: all  $\lambda \in [-1, 0]$  satisfy

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

# SUBGRADIENT CALCULUS

Suppose  $g, g_1, g_2, \dots, g_n$  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  $\partial f(\mathbf{x}) = \partial g_1(\mathbf{x}) + \partial g_2(\mathbf{x})$ .
- ▶ **Positive scaling:** If  $f(\mathbf{x}) = \alpha \cdot g(\mathbf{x})$  for some  $\alpha > 0$ , then  $\partial f(\mathbf{x}) = \alpha \cdot \partial g(\mathbf{x})$ .
- ▶ **Affine composition:** If  $f(\mathbf{x}) = g(\mathbf{A}\mathbf{x} + \mathbf{b})$ , then  $\partial f(\mathbf{x}) = \mathbf{A}^\top \partial g(\mathbf{A}\mathbf{x} + \mathbf{b})$ .
- ▶ **Finite pointwise maximum:** If  $f(\mathbf{x}) = \max_{i \in [n]} g_i(\mathbf{x})$ , then

$$\partial f(\mathbf{x}) = \text{conv} \left( \bigcup_{i \in [n]: g_i(\mathbf{x}) = f(\mathbf{x})} \partial g_i(\mathbf{x}) \right).$$

- ▶ **More general composition:** If  $h: \mathbb{R}^n \rightarrow \mathbb{R}$  is convex and non-decreasing in each argument, and  $f(\mathbf{x}) := h(g_1(\mathbf{x}), g_2(\mathbf{x}), \dots, g_n(\mathbf{x}))$ , then

$$\partial f(\mathbf{x}) = \bigcup_{(\lambda_1, \lambda_2, \dots, \lambda_n) \in \partial h(g_1(\mathbf{x}), g_2(\mathbf{x}), \dots, g_n(\mathbf{x}))} \partial \left( \sum_{i=1}^n \lambda_i g_i(\mathbf{x}) \right).$$

## 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  $\boldsymbol{\lambda}^{(t)}$  of  $f$  at  $\mathbf{x}^{(t)}$ .
  - ▶ Update:

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