



# DS-GA 3001.007

## Introduction to Machine Learning

### Lecture 10

### Support Vector Machines - Classifying with Hinge Loss



Extending Perceptron  
Algorithm by  
Incorporating Margins

# DS-GA 3001.007

## Introduction to Machine Learning

Lecture 10

Support Vector Machines - Classifying with Hinge Loss

# Announcements

- ▶ Homework 4 extended to **Wednesday November 13** at 11:59pm
- ▶ Survey 3 due **Sunday November 10** at 11:59pm
- ▶ Project
  - ▶ Milestone due **Thursday November 28** at 11:59pm
  - ▶ Background
  - ▶ Plans
    - ▶ Description of Methodology
    - ▶ Proposed Experiments
    - ▶ Some Relevant Datasets



# Review: Loss Functions for Classification

## ► Notation

Outcome space  $\mathcal{Y} = \{-1, 1\}$

Action space  $\mathcal{A} = \{-1, 1\}$

## ► 0-1 Loss

$$\ell(f(x), y) = 1(f(x) \neq y)$$

# Review: Loss Functions for Classification

## ► Notation

Outcome space  $\mathcal{Y} = \{-1, 1\}$

Action space  $\mathcal{A} = \{-1, 1\}$

## ► 0-1 Loss

$$\ell(f(x), y) = 1(f(x) \neq y)$$

Does not capture certainty about the classification



## ► Notation

Output space  $\mathcal{Y} = \{-1, 1\}$

Action space  $\mathcal{A} = \mathbf{R}$

# Review: Loss Functions for Classification

## ► Notation

Outcome space  $\mathcal{Y} = \{-1, 1\}$

Action space  $\mathcal{A} = \{-1, 1\}$

## ► 0-1 Loss

$$\ell(f(x), y) = 1(f(x) \neq y)$$

Does not capture certainty about the classification



## ► Notation

Output space  $\mathcal{Y} = \{-1, 1\}$

Action space  $\mathcal{A} = \mathbb{R}$

## ► Margin

► For prediction  $f(x)$  and label  $y \in \{-1, 1\}$  is  $f(x) y$

► Same sign means positive value. Different sign means negative

► Positive means correct. Negative means incorrect.

# Review: Loss Functions for Classification

## ► Notation

Outcome space  $\mathcal{Y} = \{-1, 1\}$

Action space  $\mathcal{A} = \{-1, 1\}$

## ► 0-1 Loss

$$\ell(f(x), y) = 1(f(x) \neq y)$$

Does not capture certainty about the classification

Functional Margin  
not Geometric Margin

## ► Notation

Output space  $\mathcal{Y} = \{-1, 1\}$

Action space  $\mathcal{A} = \mathbb{R}$

## ► Margin

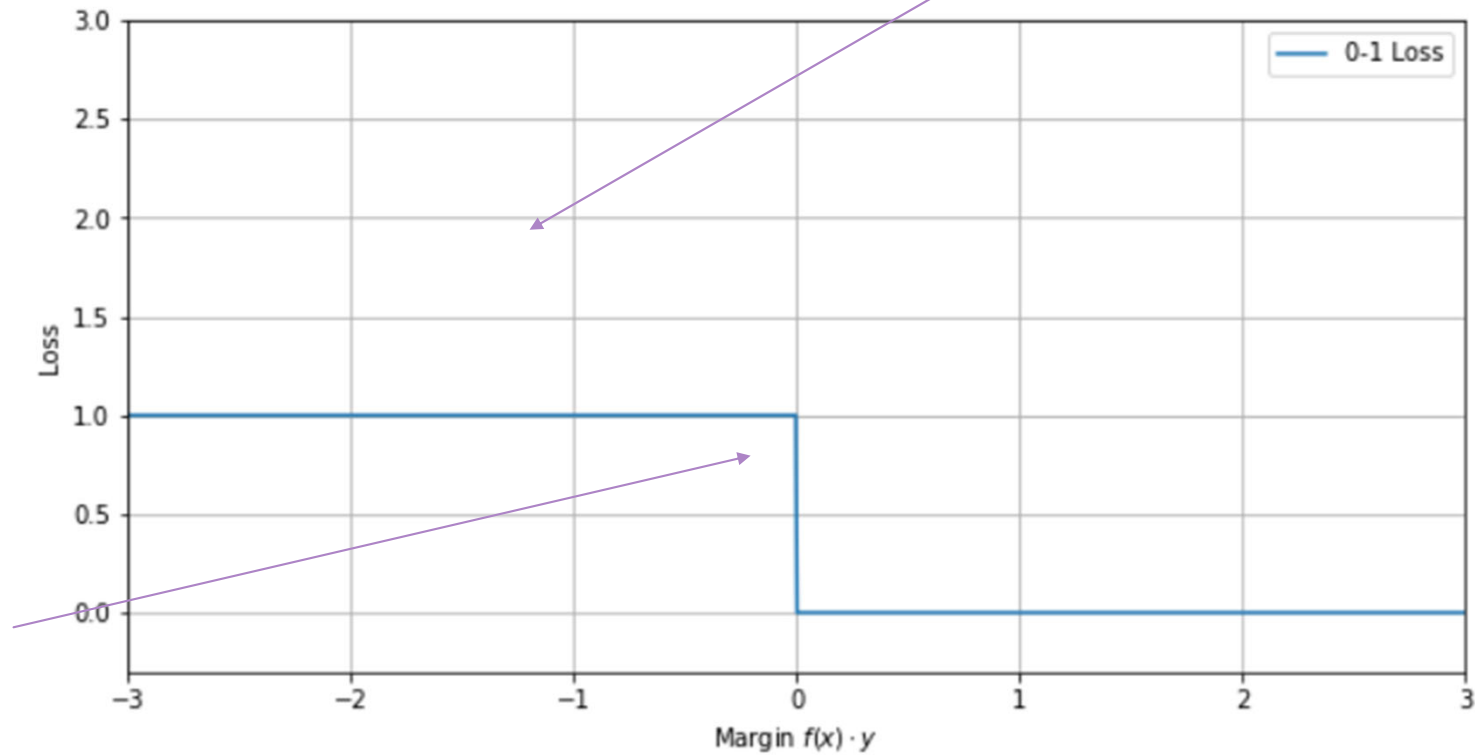
► For prediction  $f(x)$  and label  $y \in \{-1, 1\}$  is  $f(x) y$

► Same sign means positive value. Different sign means negative

► Positive means correct. Negative means incorrect.

## Review: Loss Functions for Classification

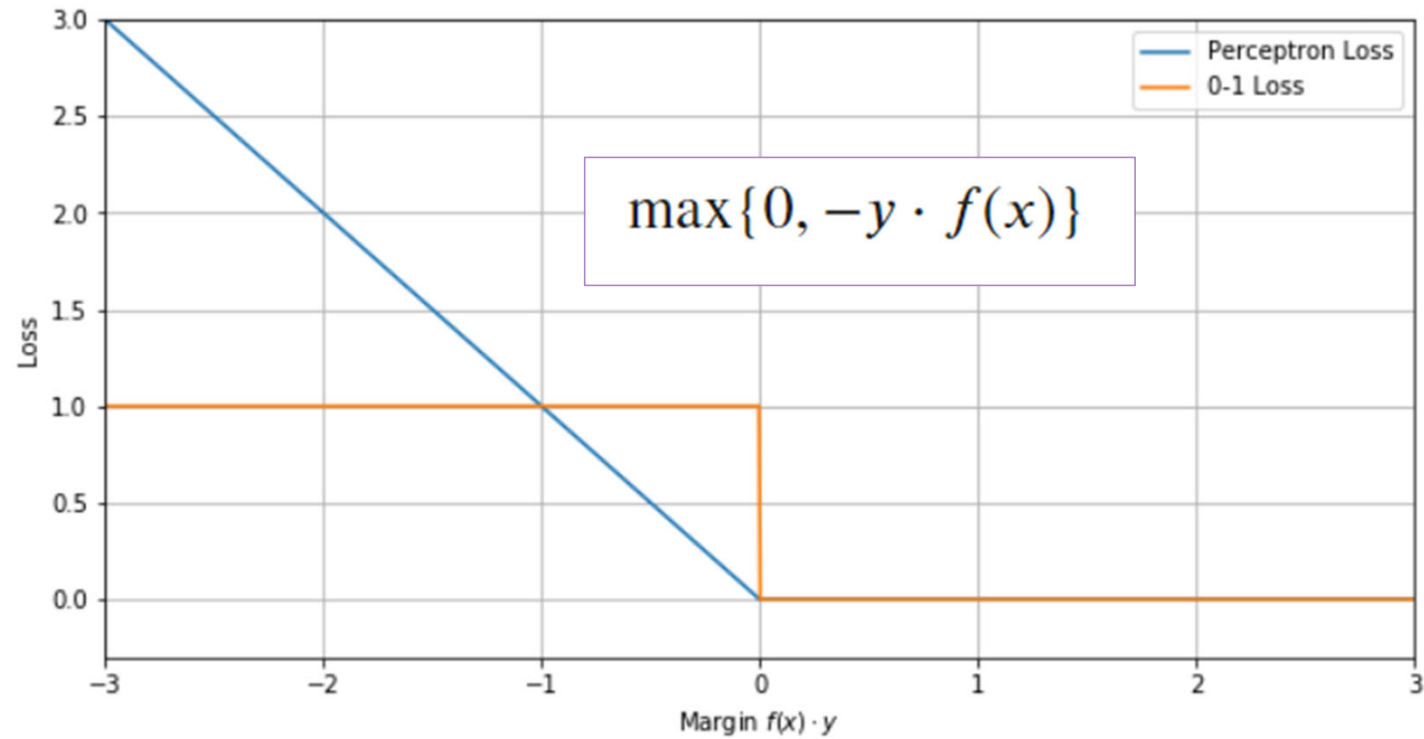
Not convex meaning  
no **subgradient** at  
decision boundary



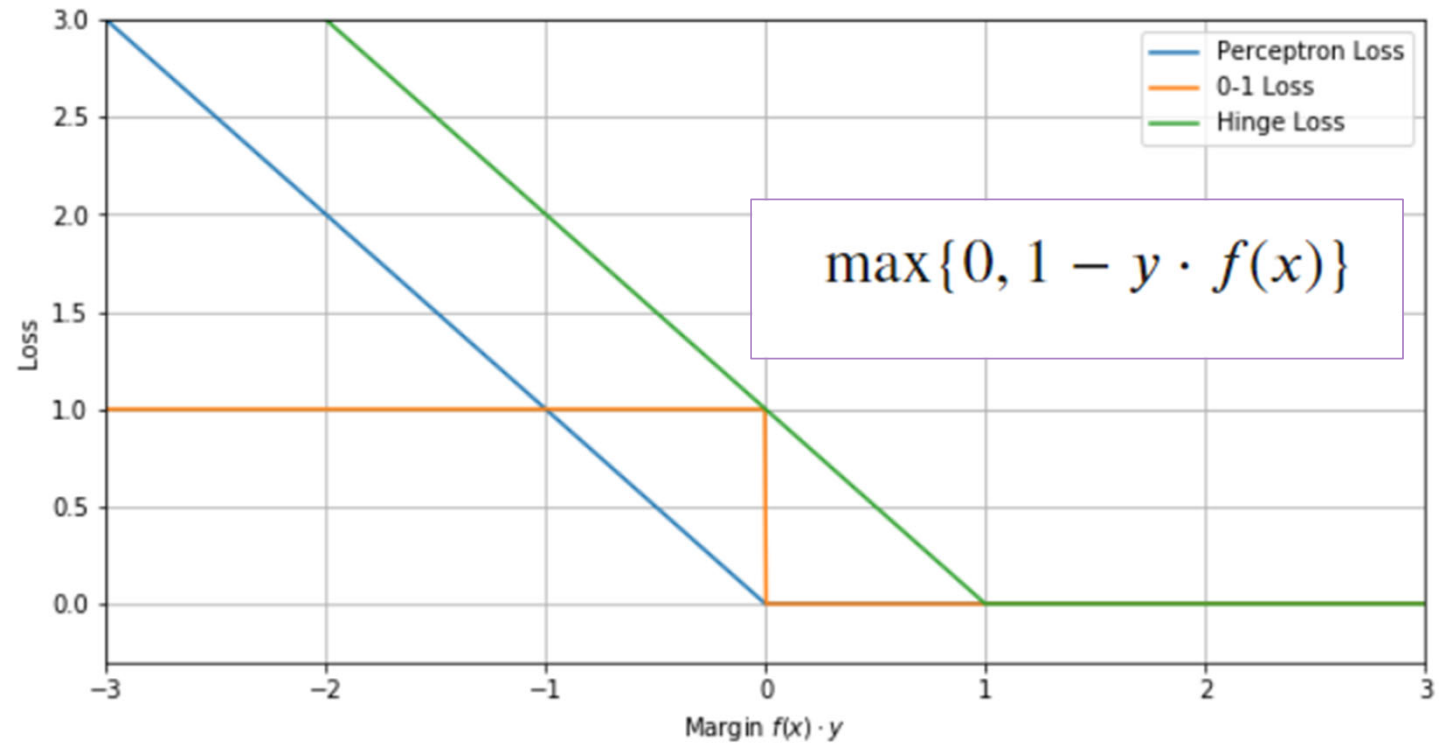
No gradient  
at decision  
boundary



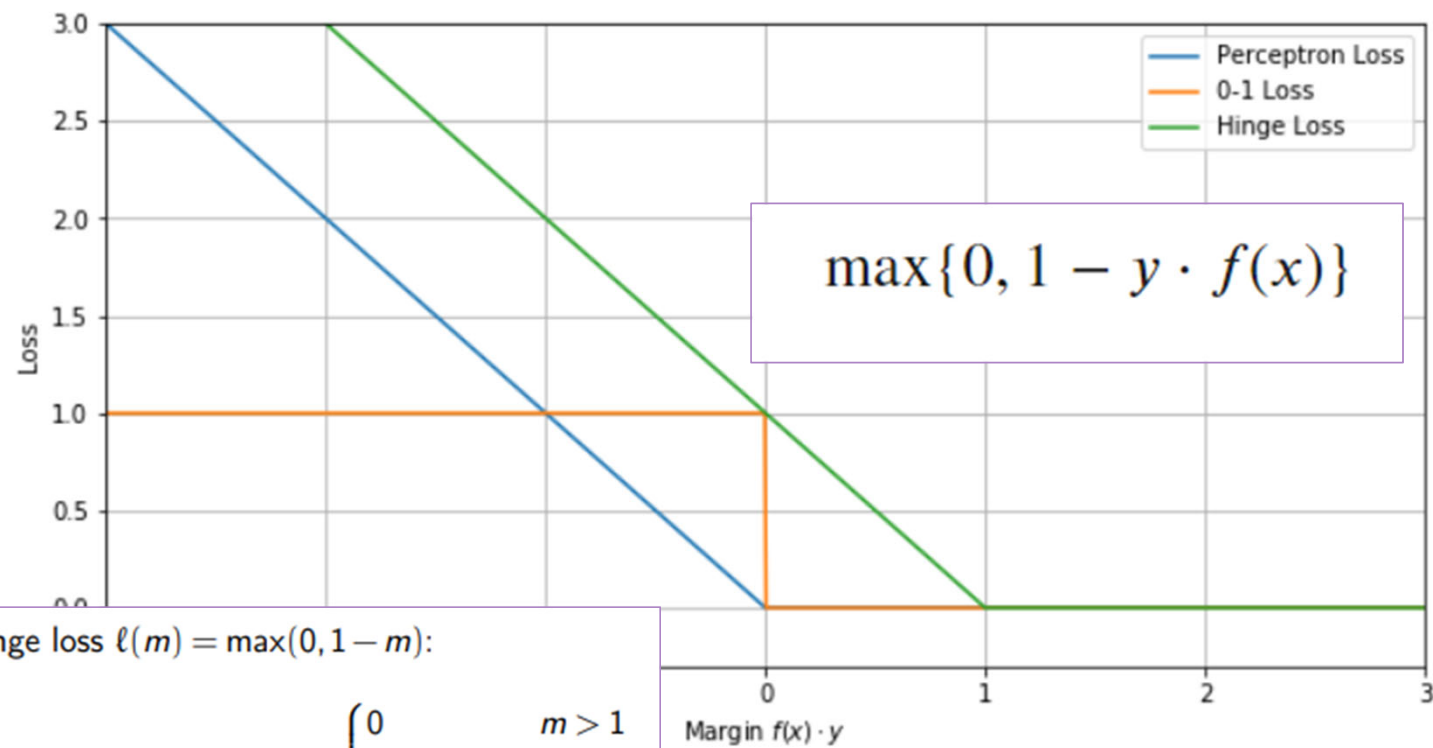
## Review: Loss Functions for Classification



# Review: Loss Functions for Classification



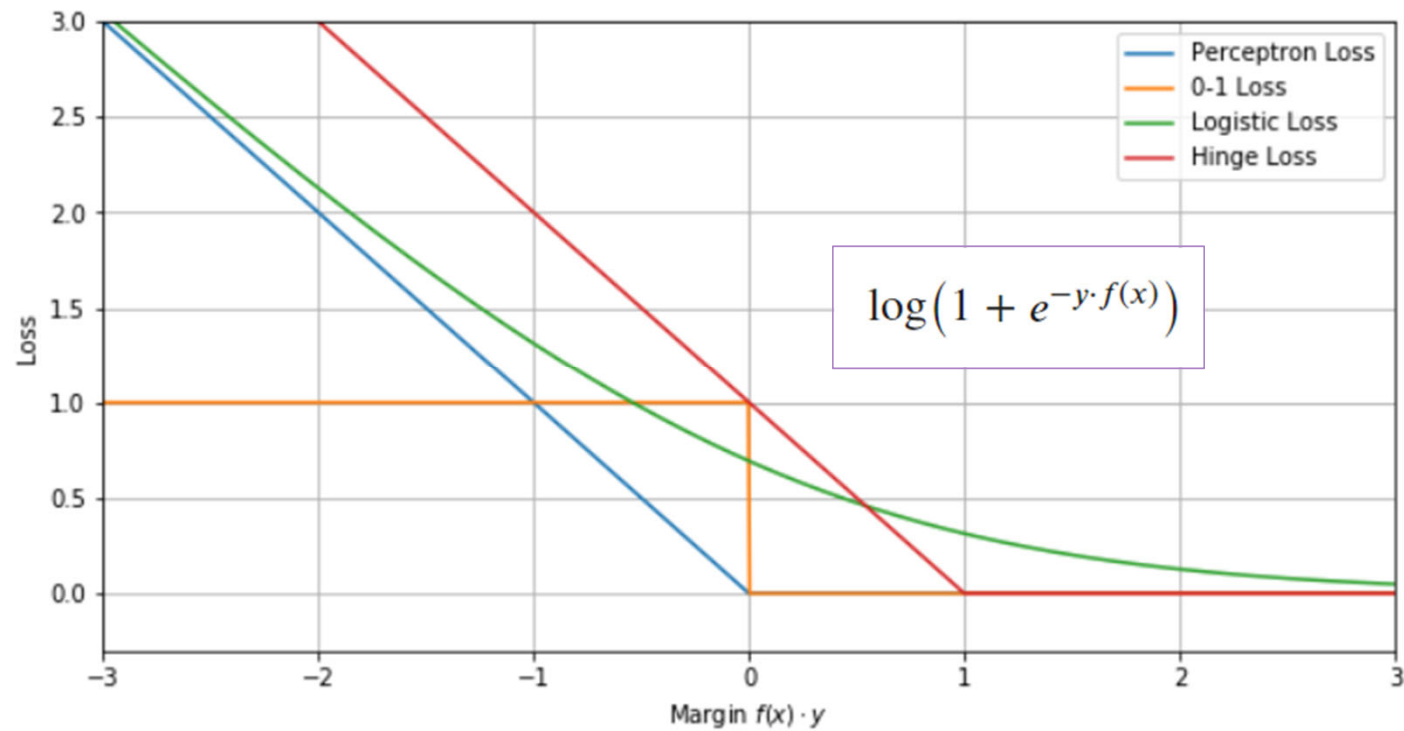
# Review: Loss Functions for Classification



Derivative of hinge loss  $\ell(m) = \max(0, 1 - m)$ :

$$\ell'(m) = \begin{cases} 0 & m > 1 \\ -1 & m < 1 \\ \text{undefined} & m = 1 \end{cases}$$

# Review: Loss Functions for Classification



# Agenda

- ▶ Lesson
  - ▶ Support Vector Machines
    - ▶ Hard Margin
    - ▶ Soft Margin
  - ▶ Convexity and Subgradients
  - ▶ Rearranging Optimization Problems
- ▶ Demo
  - ▶ Classifying Images with SVM

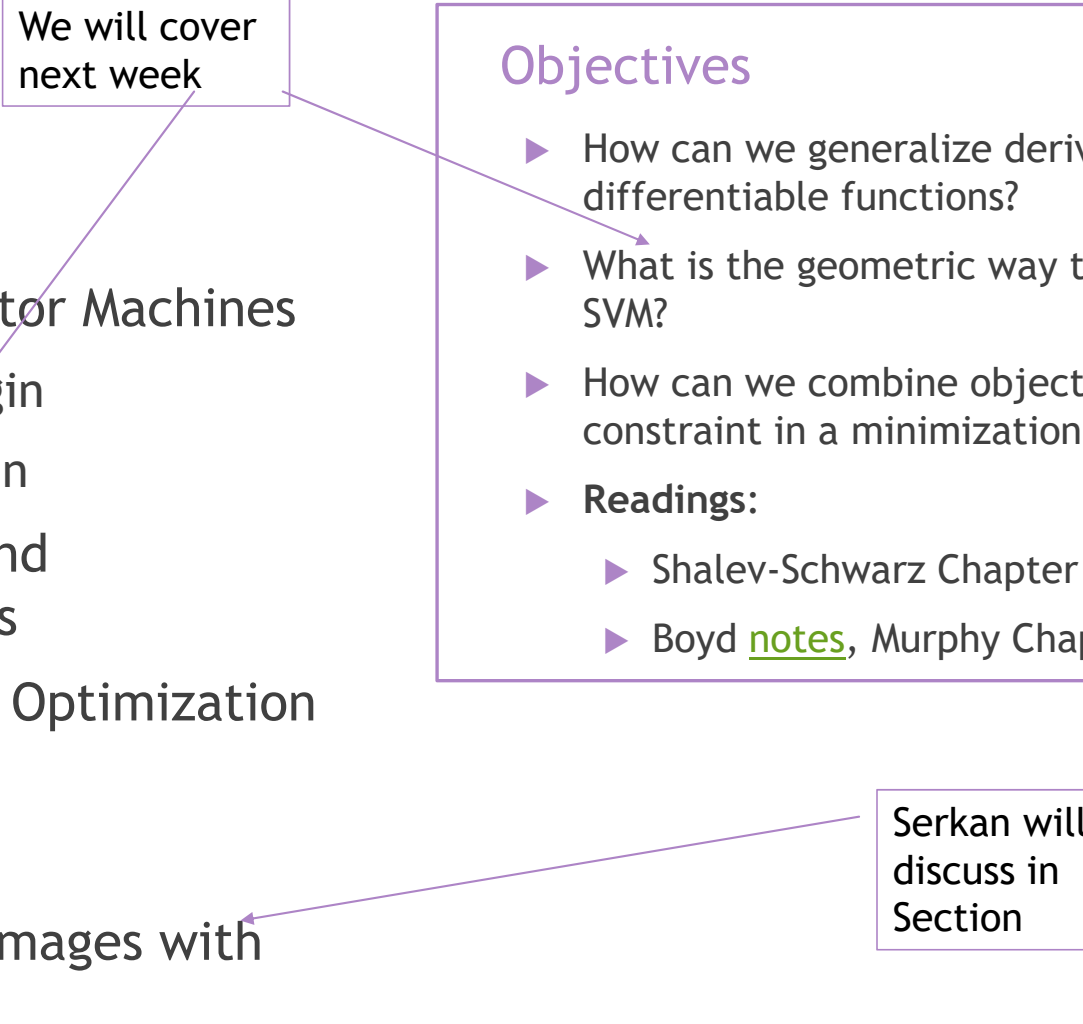
## Objectives

- ▶ How can we generalize derivatives to non-differentiable functions
- ▶ What is the geometric way to understand SVM?
- ▶ How can we combine objective and constraint in a minimization problem?
- ▶ **Readings:**
  - ▶ Shalev-Schwarz Chapter 9
  - ▶ Boyd [notes](#), Murphy Chapter 8.3

# Agenda

- ▶ Lesson
  - ▶ Support Vector Machines
    - ▶ Hard Margin
    - ▶ Soft Margin
  - ▶ Convexity and Subgradients
  - ▶ Rearranging Optimization Problems
- ▶ Demo
  - ▶ Classifying Images with SVM

We will cover next week



## Objectives

- ▶ How can we generalize derivatives to non-differentiable functions?
- ▶ What is the geometric way to understand SVM?
- ▶ How can we combine objective and constraint in a minimization problem?
- ▶ Readings:
  - ▶ Shalev-Schwarz Chapter 9
  - ▶ Boyd [notes](#), Murphy Chapter 8.3

Serkan will discuss in Section

# Agenda

## ▶ Lesson

### ▶ Support Vector Machines

▶ Hard Margin

▶ Soft Margin

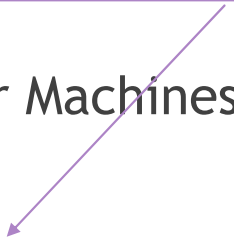
### ▶ Convexity and Subgradients

### ▶ Rearranging Optimization Problems

## ▶ Demo

### ▶ Classifying Images with SVM

Useful for  
classification and  
outlier detection



## Objectives

- ▶ How can we generalize derivatives to non-differentiable functions
- ▶ What is the geometric way to understand SVM?
- ▶ How can we combine objective and constraint in a minimization problem?
- ▶ **Readings:**
  - ▶ Shalev-Schwarz Chapter 9
  - ▶ Boyd [notes](#), Murphy Chapter 8.3

Serkan will  
discuss in  
Section



# Agenda

## ▶ Lesson

- ▶ Support Vector Machines
  - ▶ Hard Margin
  - ▶ Soft Margin
- ▶ Convexity and Subgradients
- ▶ Rearranging Optimization Problems


Useful for working with absolute value



## ▶ Demo

- ▶ Classifying Images with SVM

Serkan will discuss in Section



## Objectives

- ▶ How can we generalize derivatives to non-differentiable functions
- ▶ What is the geometric way to understand SVM?
- ▶ How can we combine objective and constraint in a minimization problem?
- ▶ **Readings:**
  - ▶ Shalev-Schwarz Chapter 9
  - ▶ Boyd [notes](#), Murphy Chapter 8.3




# Agenda

## ▶ Lesson

- ▶ Support Vector Machines
  - ▶ Hard Margin
  - ▶ Soft Margin
- ▶ Convexity and Subgradients
- ▶ Rearranging Optimization Problems


Useful for  
determining  
features



## ▶ Demo

- ▶ Classifying Images with SVM

Serkan will  
discuss in  
Section



## Objectives

- ▶ How can we generalize derivatives to non-differentiable functions
- ▶ What is the geometric way to understand SVM?
- ▶ How can we combine objective and constraint in a minimization problem?
- ▶ **Readings:**
  - ▶ Shalev-Schwarz Chapter 9
  - ▶ Boyd [notes](#), Murphy Chapter 8.3

# Support Vector Machines

Soft Margin

$$\min_{w \in \mathbf{R}^d, b \in \mathbf{R}} \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \max(0, 1 - y_i [w^T x_i + b]) .$$

# Support Vector Machines

Soft Margin

$$\min_{w \in \mathbf{R}^d, b \in \mathbf{R}} \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \max(0, 1 - y_i [w^T x_i + b]).$$

Penalization form not  
constraint form with l2  
regularization not l1  
regularization

# Support Vector Machines

Soft Margin

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \max(0, 1 - y_i [w^T x_i + b]).$$

c not lambda

Penalization form not  
constraint form with l2  
regularization not l1  
regularization

# Support Vector Machines

b is intercept term in line...for classification with lines b is threshold

*Soft Margin*

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \max(0, 1 - y_i [w^T x_i + b]).$$

c not lambda

Penalization form not  
constraint form with l2  
regularization not l1  
regularization

# Support Vector Machines

*Soft Margin*

b is intercept term in line...for classification with lines b is threshold

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \max(0, 1 - y_i [w^T x_i + b]).$$

c not lambda

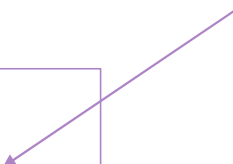
Penalization form not constraint form with l2 regularization not l1 regularization

While w and b are unconstrained, the objective is not differentiable...so use make sense of gradient or rearrange

# Support Vector Machines

Soft Margin

$$\min_{w \in \mathbf{R}^d, b \in \mathbf{R}} \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \max(0, 1 - y_i [w^T x_i + b]) .$$


$$\begin{array}{ll} \text{minimize} & \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \xi_i \\ \text{subject to} & \xi_i \geq \max(0, 1 - y_i [w^T x_i + b]) . \end{array}$$

# Support Vector Machines

Soft Margin

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \max(0, 1 - y_i [w^T x_i + b]).$$

minimize  $\frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \xi_i$

subject to  $\xi_i \geq \max(0, 1 - y_i [w^T x_i + b]).$

minimize  $\frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \xi_i$

subject to  $\xi_i \geq (1 - y_i [w^T x_i + b])$  for  $i = 1, \dots, n$   
 $\xi_i \geq 0$  for  $i = 1, \dots, n$



# Support Vector Machines

Soft Margin

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \max(0, 1 - y_i [w^T x_i + b]).$$

minimize  $\frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \xi_i$

subject to  $\xi_i \geq \max(0, 1 - y_i [w^T x_i + b]).$

Differentiable with  $n$   
+  $d + 1$  unknowns

minimize  $\frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \xi_i$

subject to  $\xi_i \geq (1 - y_i [w^T x_i + b])$  for  $i = 1, \dots, n$   
 $\xi_i \geq 0$  for  $i = 1, \dots, n$

# Support Vector Machines

Soft Margin

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \max(0, 1 - y_i [w^T x_i + b]).$$

minimize  $\frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \xi_i$   
subject to  $\xi_i \geq \max(0, 1 - y_i [w^T x_i + b]).$

Quadratic  
Programming  
Problem...could solve  
with [CVXOPT](#)

Differentiable with  $n$   
+  $d + 1$  unknowns

minimize  $\frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \xi_i$   
subject to  $\xi_i \geq (1 - y_i [w^T x_i + b])$  for  $i = 1, \dots, n$   
 $\xi_i \geq 0$  for  $i = 1, \dots, n$

# Support Vector Machines

Derivative of hinge loss  $\ell(m) = \max(0, 1 - m)$ :

$$\ell'(m) = \begin{cases} 0 & m > 1 \\ -1 & m < 1 \\ \text{undefined} & m = 1 \end{cases}$$

$$\begin{aligned} \nabla_w \ell(y_i w^T x_i) &= \ell'(y_i w^T x_i) y_i x_i \text{ (chain rule)} \\ &= \left( \begin{cases} 0 & y_i w^T x_i > 1 \\ -1 & y_i w^T x_i < 1 \\ \text{undefined} & y_i w^T x_i = 1 \end{cases} \right) y_i x_i \text{ (expanded } m \text{ in } \ell'(m)) \\ &= \begin{cases} 0 & y_i w^T x_i > 1 \\ -y_i x_i & y_i w^T x_i < 1 \\ \text{undefined} & y_i w^T x_i = 1 \end{cases} \end{aligned}$$

# Support Vector Machines

$$\nabla_w \ell(y_i w^T x_i) = \begin{cases} 0 & y_i w^T x_i > 1 \\ -y_i x_i & y_i w^T x_i < 1 \\ \text{undefined} & y_i w^T x_i = 1 \end{cases}$$

# Support Vector Machines

$$\nabla_w \ell(y_i w^T x_i) = \begin{cases} 0 & y_i w^T x_i > 1 \\ -y_i x_i & y_i w^T x_i < 1 \\ \text{undefined} & y_i w^T x_i = 1 \end{cases}$$

$$\begin{aligned} \nabla_w J(w) &= \nabla_w \left( \frac{1}{n} \sum_{i=1}^n \ell(y_i w^T x_i) + \lambda \|w\|^2 \right) \\ &= \frac{1}{n} \sum_{i=1}^n \nabla_w \ell(y_i w^T x_i) + 2\lambda w \\ &= \begin{cases} \frac{1}{n} \sum_{i: y_i w^T x_i < 1} (-y_i x_i) + 2\lambda w & \text{all } y_i w^T x_i \neq 1 \\ \text{undefined} & \text{otherwise} \end{cases} \end{aligned}$$

# Support Vector Machines

$$\nabla_w \ell(y_i w^T x_i) = \begin{cases} 0 & y_i w^T x_i > 1 \\ -y_i x_i & y_i w^T x_i < 1 \\ \text{undefined} & y_i w^T x_i = 1 \end{cases}$$

$$\begin{aligned} \nabla_w J(w) &= \nabla_w \left( \frac{1}{n} \sum_{i=1}^n \ell(y_i w^T x_i) + \lambda \|w\|^2 \right) \\ &= \frac{1}{n} \sum_{i=1}^n \nabla_w \ell(y_i w^T x_i) + 2\lambda w \\ &= \begin{cases} \frac{1}{n} \sum_{i: y_i w^T x_i < 1} (-y_i x_i) + 2\lambda w & \text{all } y_i w^T x_i \neq 1 \\ \text{undefined} & \text{otherwise} \end{cases} \end{aligned}$$

Does it make sense to check this on the computer...with floating point numbers

# Demo

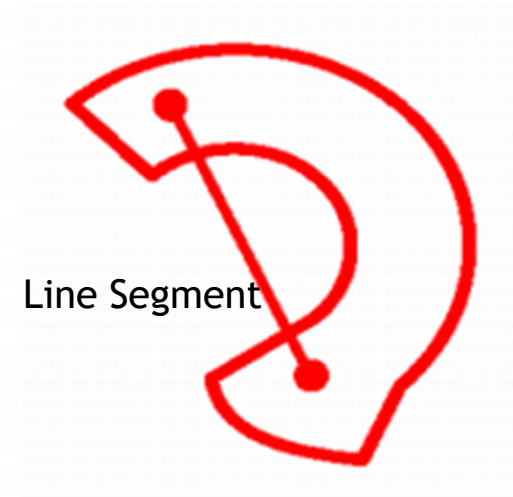
- ▶ Support Vector Machine
  - ▶ Iris Dataset
  - ▶ Features
    - ▶ Petal Width
    - ▶ Petal Length
  - ▶ Classification
    - ▶ Iris-Versicolor
    - ▶ Iris-Setosa

## Take-Aways

- ▶ Why is SVM affected by scaling?
- ▶ How can soft margin SVM be used to detect outliers?
- ▶ How does changing  $C$  affect the classification? What prevents against overfitting.
- ▶ How do we use SVM in sklearn?

# Convexity

Convex  
Sets

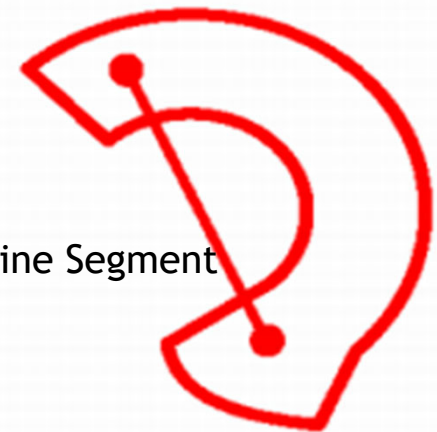


Line Segment



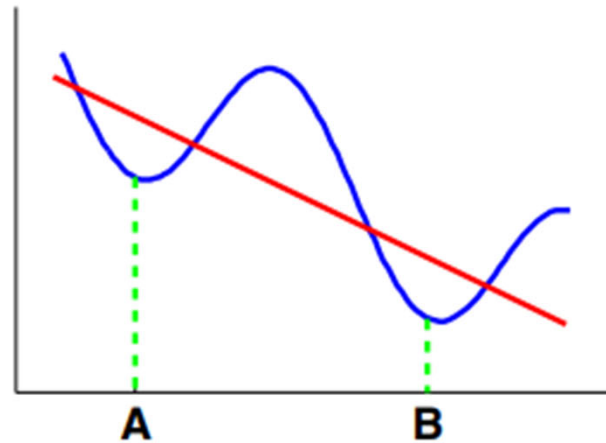
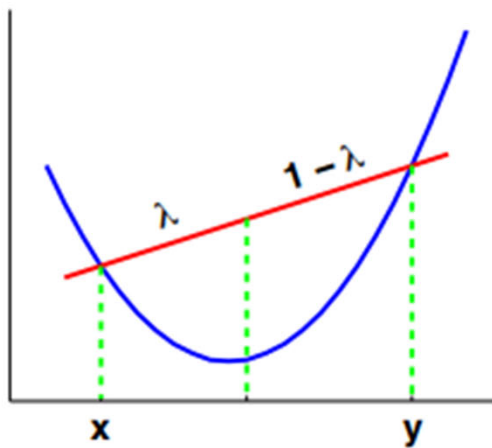
# Convexity

Convex  
Sets



Line Segment

Convex  
Function

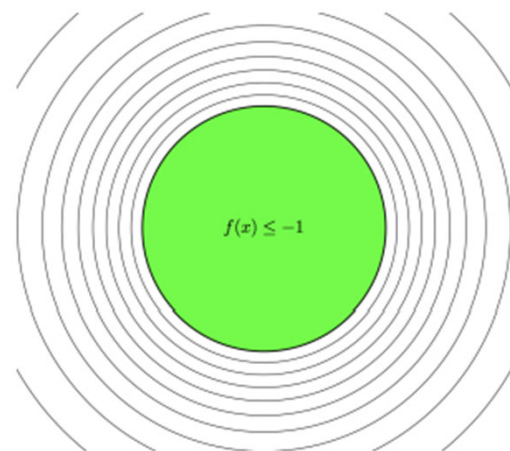


## Question

- What is a concave function?
- Can a function be both convex and concave?

# Convexity

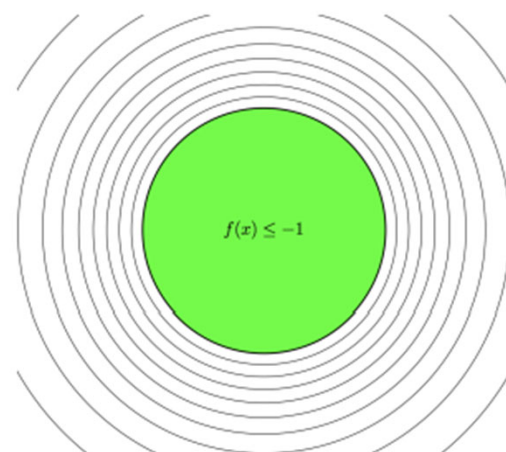
$f : \mathbf{R}^d \rightarrow \mathbf{R}$  be a function.



A level set or contour line for the value  $c$  is the set of points  $x \in \mathbf{R}^d$  for which  $f(x) = c$ .

# Convexity

$f : \mathbf{R}^d \rightarrow \mathbf{R}$  be a function.

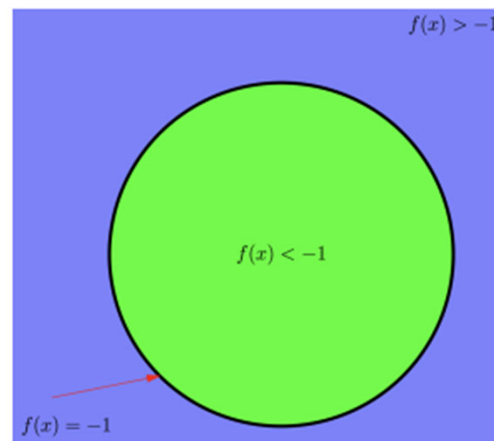


A **level set** or **contour line** for the value  $c$  is the set of points  $x \in \mathbf{R}^d$  for which  $f(x) = c$ .

A **sublevel set** for the value  $c$  is the set of points  $x \in \mathbf{R}^d$  for which  $f(x) \leq c$ .

# Convexity

$f : \mathbf{R}^d \rightarrow \mathbf{R}$  be a function.



A **level set** or **contour line** for the value  $c$  is the set of points  $x \in \mathbf{R}^d$  for which  $f(x) = c$ .

A **sublevel set** for the value  $c$  is the set of points  $x \in \mathbf{R}^d$  for which  $f(x) \leq c$ .

*If  $f : \mathbf{R}^d \rightarrow \mathbf{R}$  is convex, then the sublevel sets are convex.*

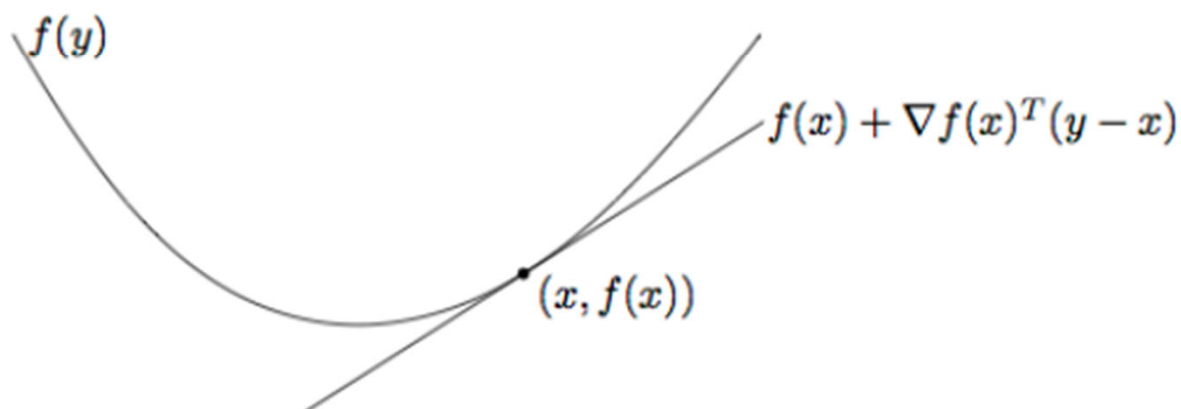
# Convexity

Suppose  $f : \mathbf{R}^d \rightarrow \mathbf{R}$  is differentiable.

Predict  $f(y)$  given  $f(x)$  and  $\nabla f(x)$ ?

Linear (i.e. “first order”) approximation:

$$f(y) \approx f(x) + \nabla f(x)^T (y - x)$$



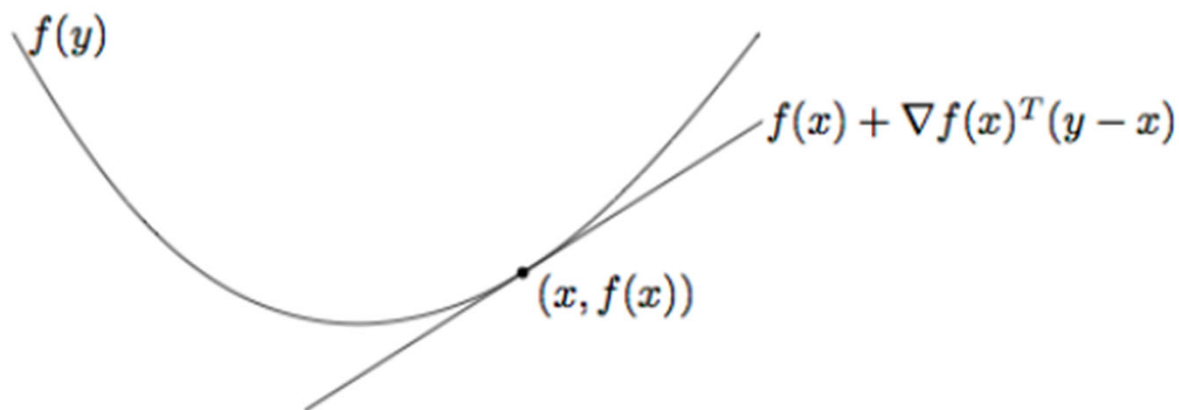
# Convexity

Suppose  $f : \mathbf{R}^d \rightarrow \mathbf{R}$  is **convex** and **differentiable**.

Then for any  $x, y \in \mathbf{R}^d$

$$f(y) \geq f(x) + \nabla f(x)^T (y - x)$$

The linear approximation to  $f$  at  $x$  is a **global underestimator** of  $f$ :



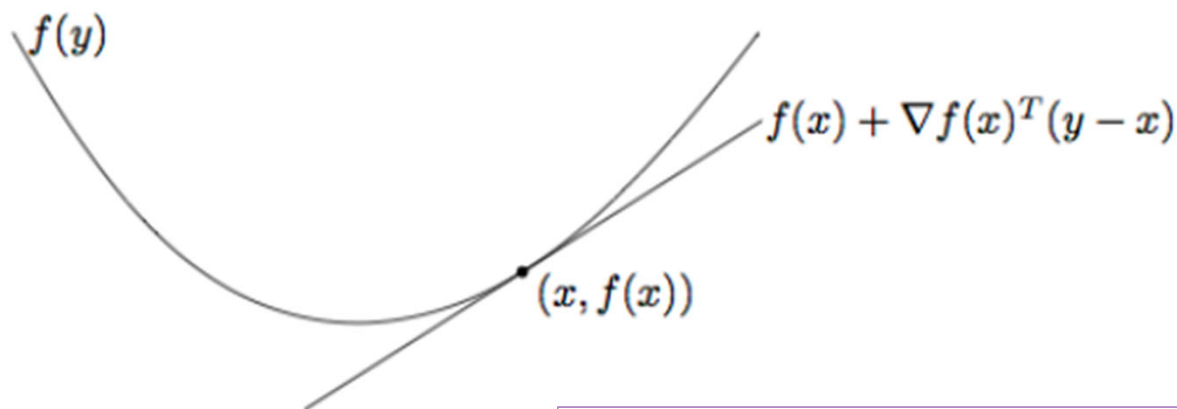
# Convexity

Suppose  $f : \mathbf{R}^d \rightarrow \mathbf{R}$  is **convex** and **differentiable**.

Then for any  $x, y \in \mathbf{R}^d$

$$f(y) \geq f(x) + \nabla f(x)^T (y - x)$$

The linear approximation to  $f$  at  $x$  is a **global underestimator** of  $f$ :

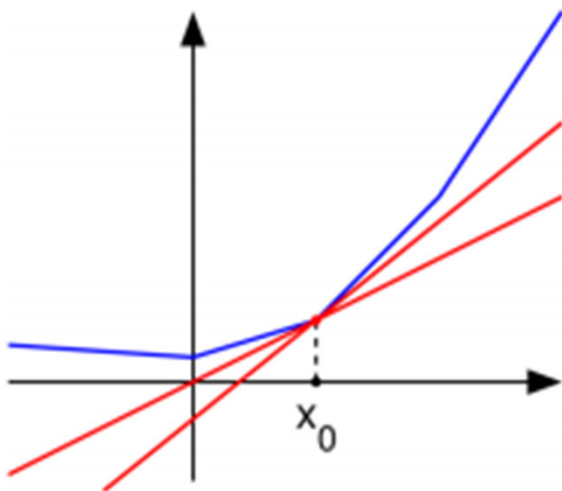


*If  $\nabla f(x) = 0$  then  $x$  is a global minimizer of  $f$ .*

# Subgradients

A vector  $g \in \mathbb{R}^d$  is a **subgradient** of  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  at  $x$  if for all  $z$ ,

$$f(z) \geq f(x) + g^T(z - x).$$

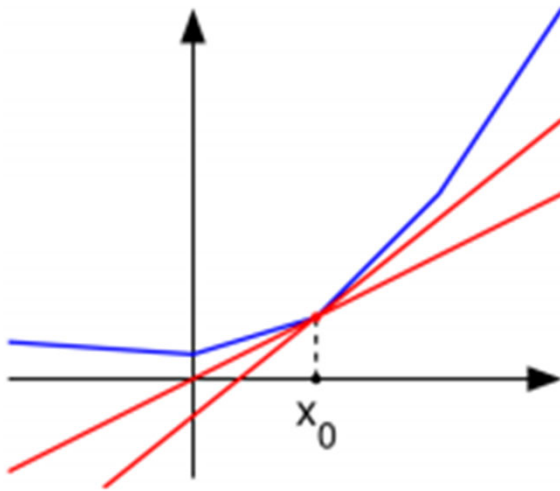




# Subgradients

A vector  $g \in \mathbb{R}^d$  is a **subgradient** of  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  at  $x$  if for all  $z$ ,

$$f(z) \geq f(x) + g^T(z - x).$$



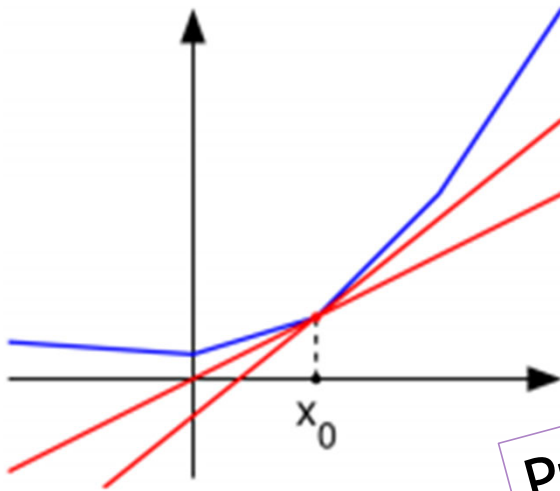
$f$  is **subdifferentiable** at  $x$  if  $\exists$  at least one subgradient at  $x$ .

The set of all subgradients at  $x$  is called the **subdifferential**:  $\partial f(x)$

# Subgradients

A vector  $g \in \mathbb{R}^d$  is a **subgradient** of  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  at  $x$  if for all  $z$ ,

$$f(z) \geq f(x) + g^T(z - x).$$



$f$  is **subdifferentiable** at  $x$  if  $\exists$  at least one subgradient at  $x$ .

The set of all subgradients at  $x$  is called the **subdifferential**:  $\partial f(x)$

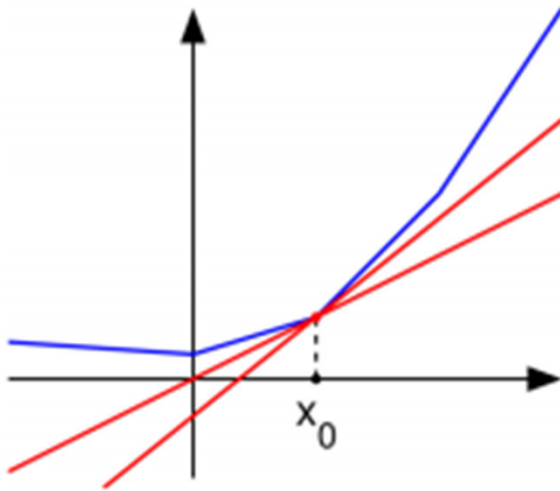
## Properties

- $f$  is convex and differentiable  $\implies \partial f(x) = \{\nabla f(x)\}$ .
- Any point  $x$ , there can be 0, 1, or infinitely many subgradients.
- $\partial f(x) = \emptyset \implies f$  is not convex.

# Subgradients

A vector  $g \in \mathbb{R}^d$  is a **subgradient** of  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  at  $x$  if for all  $z$ ,

$$f(z) \geq f(x) + g^T(z - x).$$



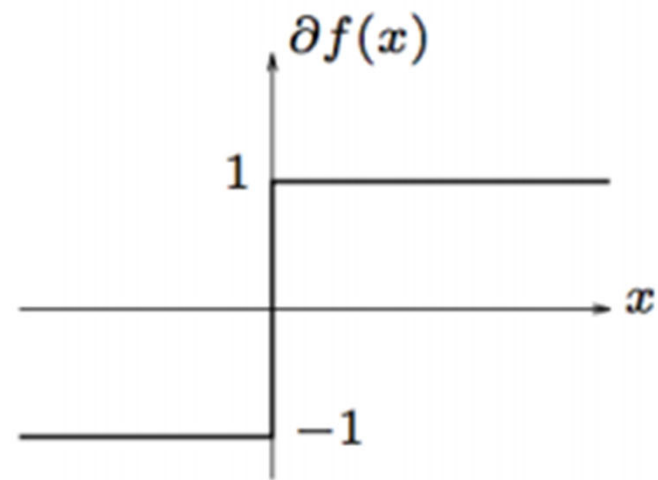
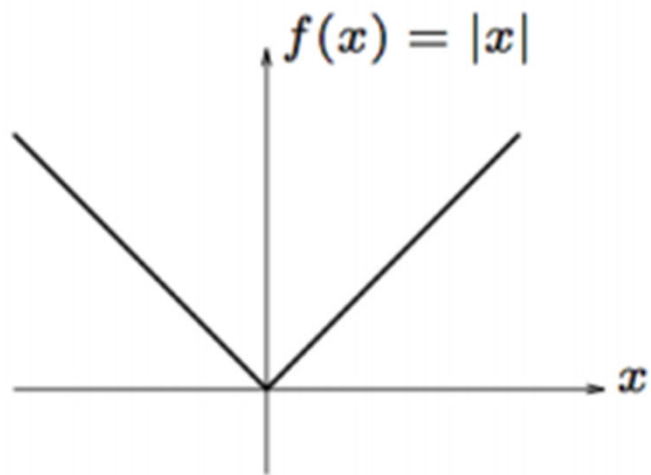
$f$  is **subdifferentiable** at  $x$  if  $\exists$  at least one subgradient at  $x$ .  
The set of all subgradients at  $x$  is called the **subdifferential**:  $\partial f(x)$

What if

$$0 \in \partial f(x)$$

- $f$  is convex and differentiable  $\implies \partial f(x) = \{\nabla f(x)\}$ .
- Any point  $x$ , there can be 0, 1, or infinitely many subgradients.
- $\partial f(x) = \emptyset \implies f$  is not convex.

# Subgradients

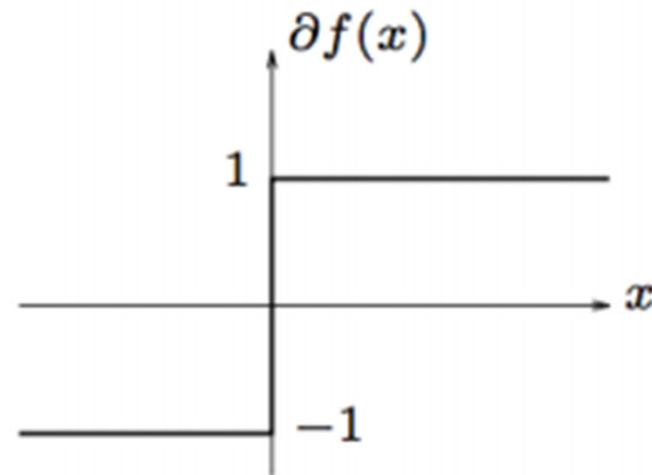
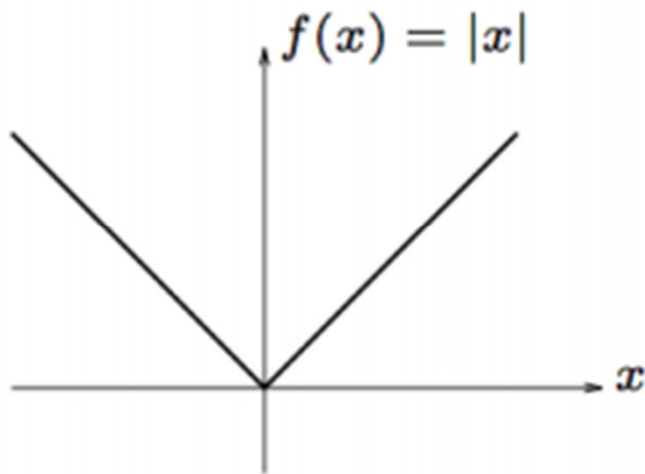


# Subgradients

Question

Let  $\mathcal{X} = \{1, \dots, 10\}$ , let  $\mathcal{Y} = \{1, \dots, 10\}$ , and let  $A = \mathcal{Y}$ . Suppose the data generating distribution,  $P$ , has marginal  $X \sim \text{Unif}\{1, \dots, 10\}$  and conditional distribution  $Y|X = x \sim \text{Unif}\{1, \dots, x\}$ . For each loss function below give a target function

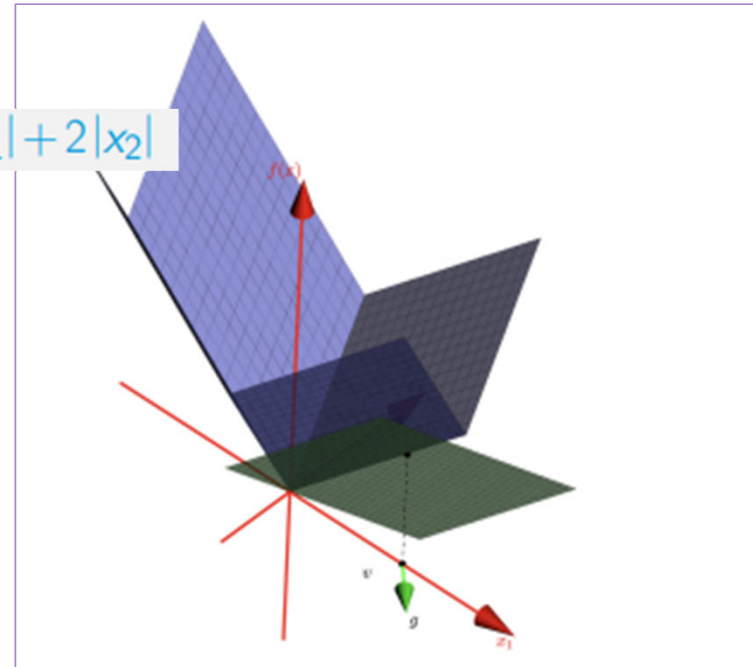
- (a)  $\ell(a, y) = (a - y)^2$ ,
- (b)  $\ell(a, y) = |a - y|$ ,
- (c)  $\ell(a, y) = 1(a \neq y)$ .



# Subgradients

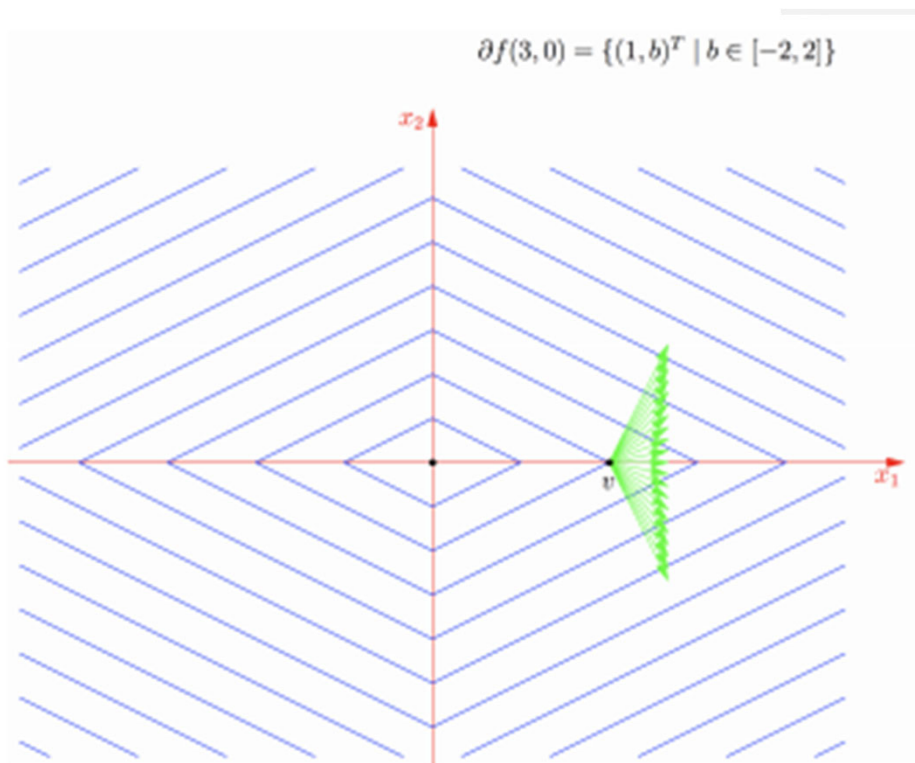
Example

$$f(x_1, x_2) = |x_1| + 2|x_2|$$

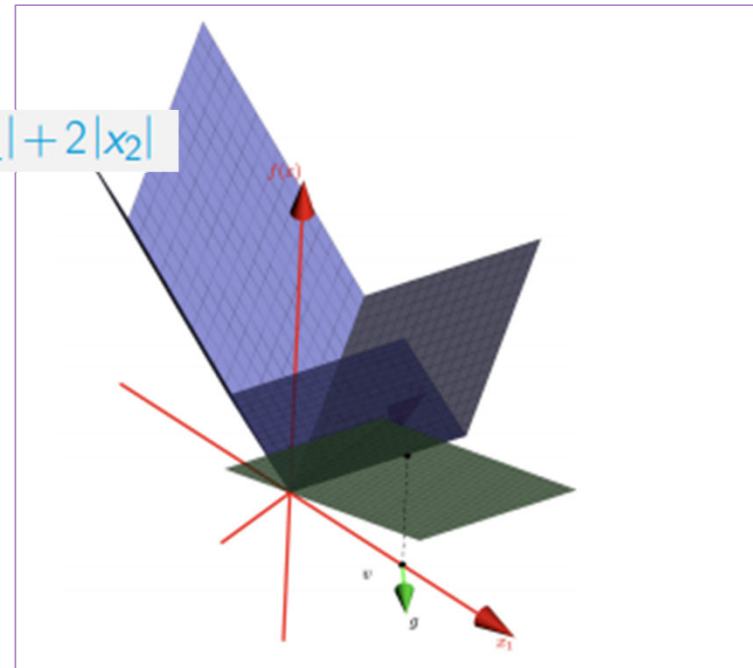


# Subgradients

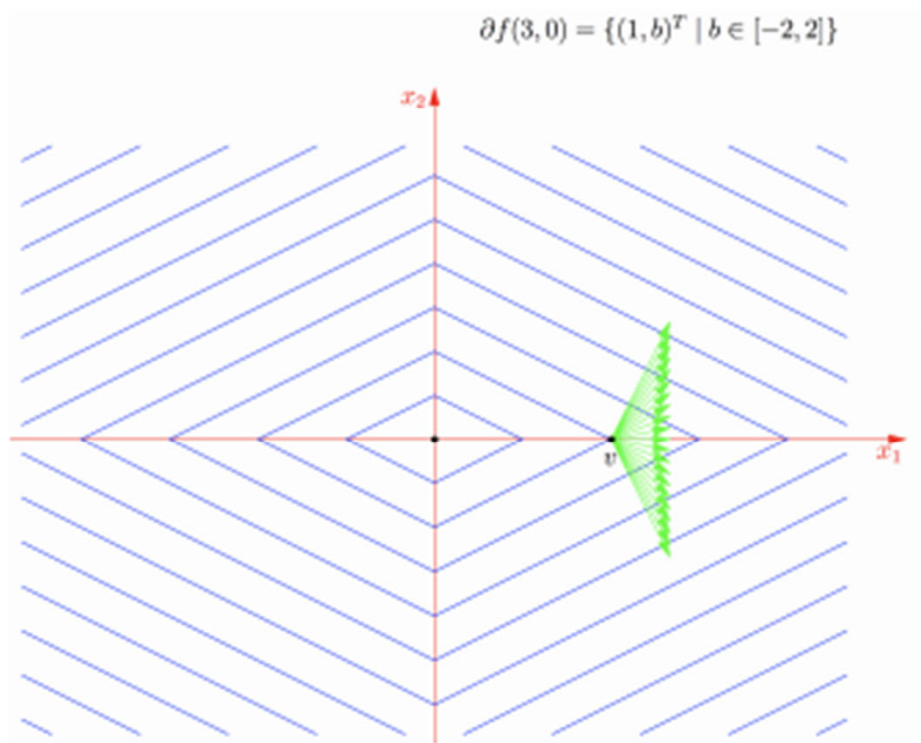
Example



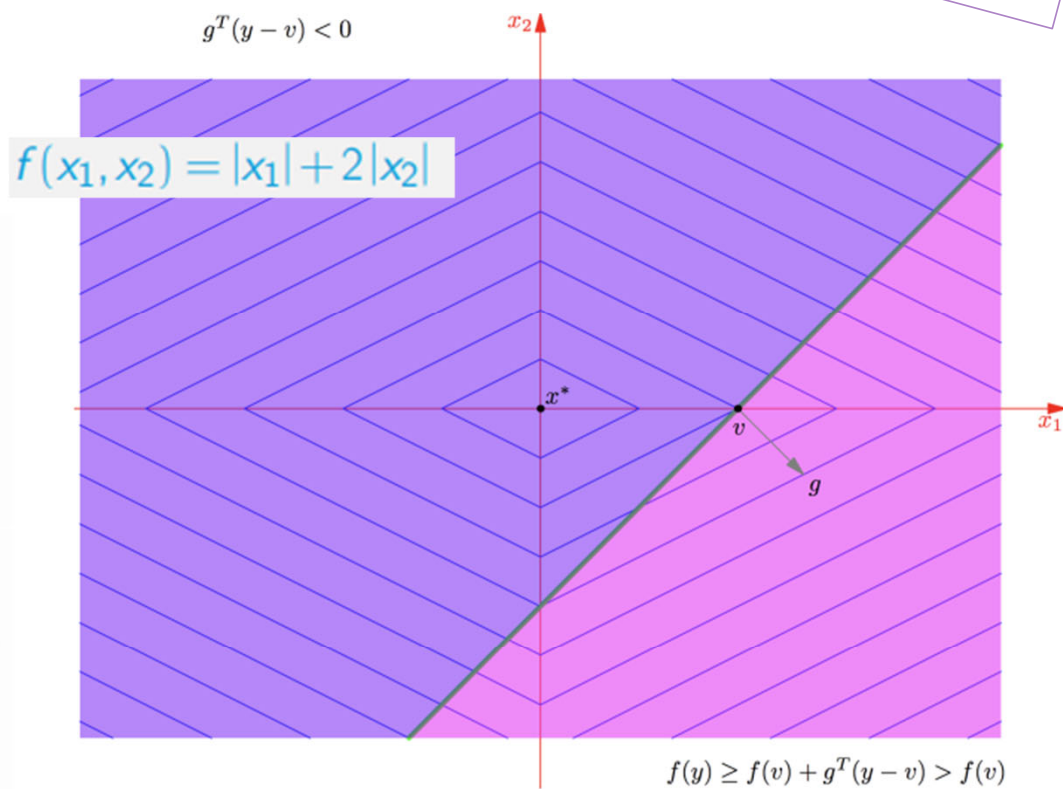
$$= |x_1| + 2|x_2|$$



# Subgradients



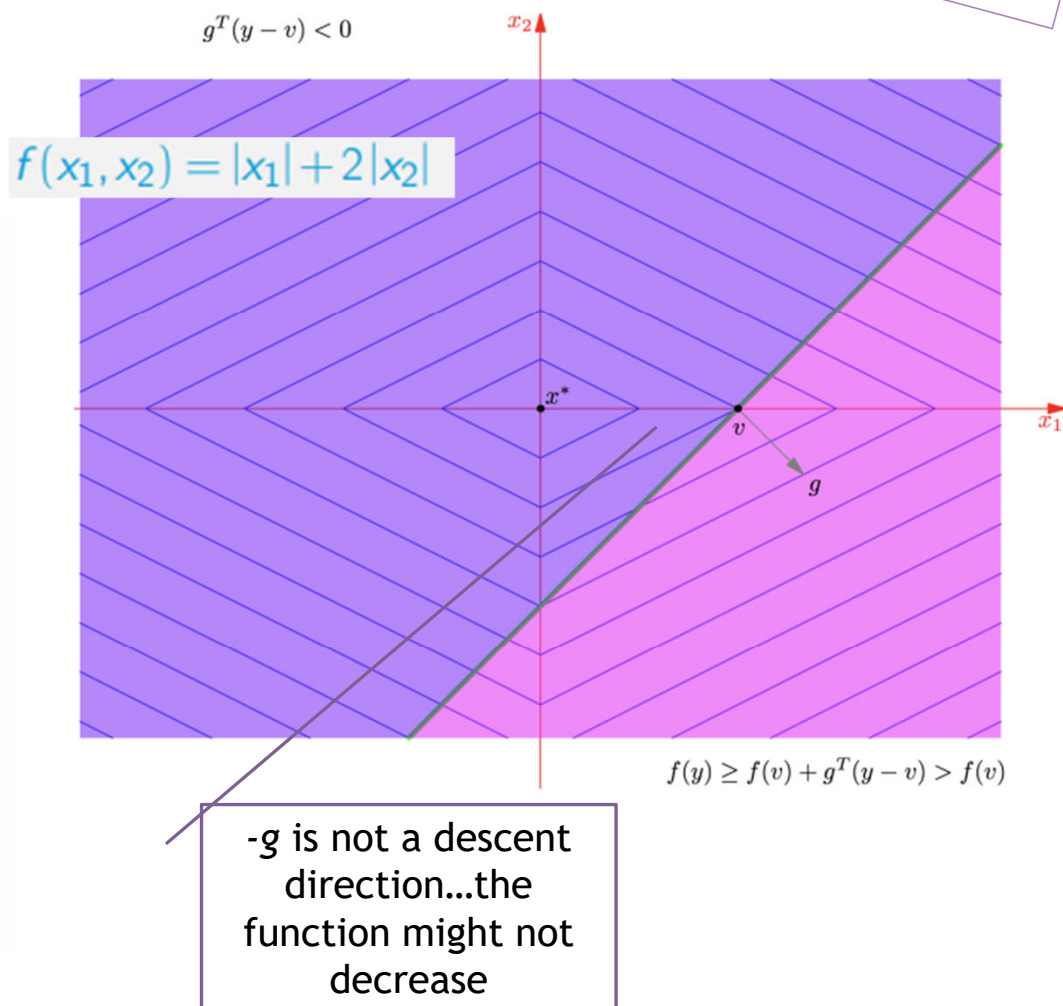
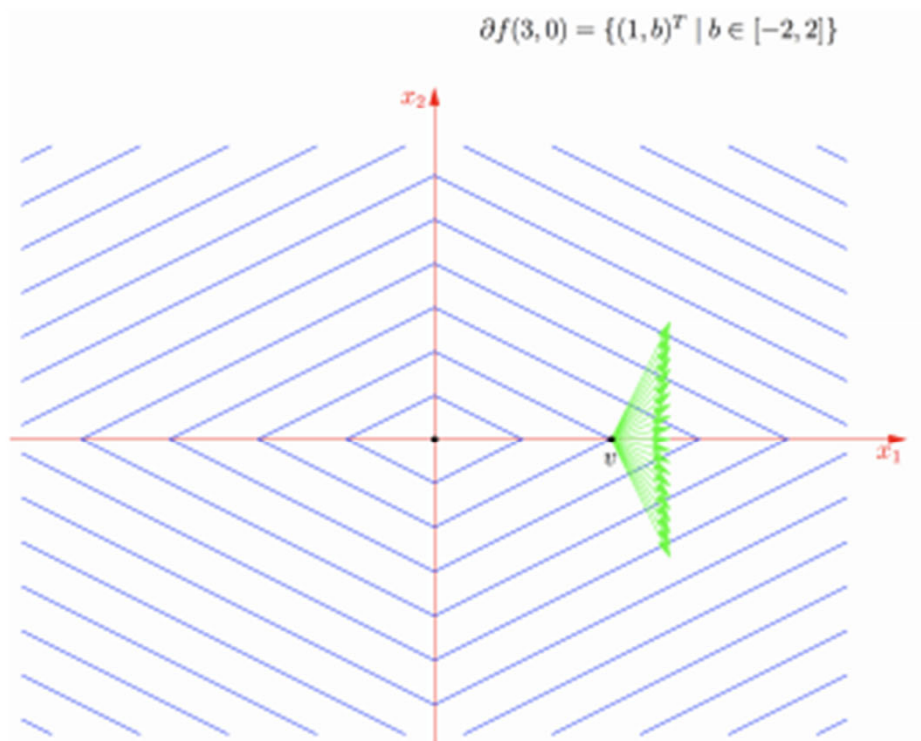
Example





# Subgradients

Example



# Subgradient Descent

*Suppose  $f$  is convex.*

- *Let  $x = x_0 - tg$ , for  $g \in \partial f(x_0)$ .*
- *Let  $z$  be any point for which  $f(z) < f(x_0)$ .*
- *Then for small enough  $t > 0$ ,*

$$\|x - z\|_2 < \|x_0 - z\|_2.$$

# Subgradient Descent

- Let  $x = x_0 - tg$ , for  $g \in \partial f(x_0)$  and  $t > 0$ .
- Let  $z$  be any point for which  $f(z) < f(x_0)$ .

*Suppose  $f$  is convex.*

- *Let  $x = x_0 - tg$ , for  $g \in \partial f(x_0)$ .*
- *Let  $z$  be any point for which  $f(z) < f(x_0)$ .*
- *Then for small enough  $t > 0$ ,*

$$\|x - z\|_2 < \|x_0 - z\|_2.$$

# Subgradient Descent

- Let  $x = x_0 - tg$ , for  $g \in \partial f(x_0)$  and  $t > 0$ .
- Let  $z$  be any point for which  $f(z) < f(x_0)$ .
- Then

*Suppose  $f$  is convex.*

- *Let  $x = x_0 - tg$ , for  $g \in \partial f(x_0)$ .*
- *Let  $z$  be any point for which  $f(z) < f(x_0)$ .*
- *Then for small enough  $t > 0$ ,*

$$\|x - z\|_2 < \|x_0 - z\|_2.$$

$$\begin{aligned}\|x - z\|_2^2 &= \|x_0 - tg - z\|_2^2 \\ &= \|x_0 - z\|_2^2 - 2tg^T(x_0 - z) + t^2\|g\|_2^2 \\ &\leq \|x_0 - z\|_2^2 - 2t[f(x_0) - f(z)] + t^2\|g\|_2^2\end{aligned}$$

When are these terms negative?

# Subgradient Descent

- Let  $x = x_0 - tg$ , for  $g \in \partial f(x_0)$  and  $t > 0$ .
- Let  $z$  be any point for which  $f(z) < f(x_0)$ .
- Then

Suppose  $f$  is convex.

- Let  $x = x_0 - tg$ , for  $g \in \partial f(x_0)$ .
- Let  $z$  be any point for which  $f(z) < f(x_0)$ .
- Then for small enough  $t > 0$ ,

$$\|x - z\|_2 < \|x_0 - z\|_2.$$

$$\begin{aligned}\|x - z\|_2^2 &= \|x_0 - tg - z\|_2^2 \\ &= \|x_0 - z\|_2^2 - 2tg^T(x_0 - z) + t^2\|g\|_2^2 \\ &\leq \|x_0 - z\|_2^2 - 2t[f(x_0) - f(z)] + t^2\|g\|_2^2\end{aligned}$$

- Consider  $-2t[f(x_0) - f(z)] + t^2\|g\|_2^2$ .
  - It's a convex quadratic (facing upwards).
  - Has zeros at  $t = 0$  and  $t = 2(f(x_0) - f(z)) / \|g\|_2^2 > 0$ .
  - Therefore, it's negative for any

$$t \in \left(0, \frac{2(f(x_0) - f(z))}{\|g\|_2^2}\right).$$

When are these terms negative?

# Rearranging Optimization Problems

Example

- How to solve the Lasso?

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|_1$$

- $\|w\|_1 = |w_1| + |w_2|$  is not differentiable!

# Rearranging Optimization Problems

Example

- How to solve the Lasso?

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|_1$$

- Replace each  $w_i$  by  $w_i^+ - w_i^-$ .
- Write  $w^+ = (w_1^+, \dots, w_d^+)$  and  $w^- = (w_1^-, \dots, w_d^-)$ .

- $\|w\|_1 = |w_1| + |w_2|$  is not differentiable!

# Rearranging Optimization Problems

Example

- How to solve the Lasso?

$$\min_{w \in \mathbf{R}^d} \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|_1$$

- Replace each  $w_i$  by  $w_i^+ - w_i^-$ .
- Write  $w^+ = (w_1^+, \dots, w_d^+)$  and  $w^- = (w_1^-, \dots, w_d^-)$ .

- $\|w\|_1 = |w_1| + |w_2|$  is not differentiable!

$$\begin{aligned} & \min_{w^+, w^- \in \mathbf{R}^d} \sum_{i=1}^n \left( (w^+ - w^-)^T x_i - y_i \right)^2 + \lambda \mathbf{1}^T (w^+ + w^-) \\ & \text{subject to } w_i^+ \geq 0 \text{ for all } i \\ & \quad \quad \quad w_i^- \geq 0 \text{ for all } i \end{aligned}$$



# Rearranging Optimization Problems

Example

Switching the order is helpful operation

- How to solve the Lasso?

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|_1$$

- $\|w\|_1 = |w_1| + |w_2|$  is not differentiable!

- Replace each  $w_i$  by  $w_i^+ - w_i^-$ .
- Write  $w^+ = (w_1^+, \dots, w_d^+)$  and  $w^- = (w_1^-, \dots, w_d^-)$ .

$$\begin{aligned} \min_{w^+, w^- \in \mathbb{R}^d} \sum_{i=1}^n \left( (w^+ - w^-)^T x_i - y_i \right)^2 + \lambda 1^T (w^+ + w^-) \\ \text{subject to } w_i^+ \geq 0 \text{ for all } i \\ w_i^- \geq 0 \text{ for all } i \end{aligned}$$

# Rearranging Optimization Problems

Example

- Suppose you want to minimize penalization form

$$\min_x f(x) + \lambda g(x)$$

Example

## Rearranging Optimization Problems

- Suppose you want to minimize penalization form

$$\min_x f(x) + \lambda g(x)$$

- Question: Is it equivalent to constraint form?

$$\min_x f(x)$$

$$\text{s.t. } g(x) \leq R$$

# Rearranging Optimization Problems

Example

- Suppose you want to minimize penalization form

$$\min_x f(x) + \lambda g(x)$$

- Question: Is it equivalent to constraint form?

$$\min_x f(x)$$

$$\text{s.t. } g(x) \leq R$$

- Suppose minimizer for penalization form is not minimizer for constraint form

$$\text{Set } R = g(x^*).$$

# Rearranging Optimization Problems

Example

- Suppose you want to minimize penalization form

$$\min_x f(x) + \lambda g(x)$$

- Question: Is it equivalent to constraint form?

$$\min_x f(x)$$

$$\text{s.t. } g(x) \leq R$$

- Suppose minimizer for penalization form is not minimizer for constraint form

$$\text{Set } R = g(x^*).$$

- Note

$$f(x') < f(x^*)$$

# Rearranging Optimization Problems

Example

- Suppose you want to minimize penalization form

$$\min_x f(x) + \lambda g(x)$$

- Question: Is it equivalent to constraint form?

$$\min_x f(x)$$

$$\text{s.t. } g(x) \leq R$$

- Suppose minimizer for penalization form is not minimizer for constraint form

$$\text{Set } R = g(x^*).$$

- Note

$$f(x') < f(x^*)$$

- Therefore

$$\underline{f(x') + \lambda g(x') < f(x^*) + \lambda g(x^*)}$$

# Rearranging Optimization Problems

Example

- Suppose you want to minimize penalization form

$$\min_x f(x) + \lambda g(x)$$

- Question: Is it equivalent to constraint form?

$$\min_x f(x)$$

$$\text{s.t. } g(x) \leq R$$

Combining objective function and constraint is helpful

- Suppose minimizer for penalization form is not minimizer for constraint form

$$\text{Set } R = g(x^*).$$

- Note

$$f(x') < f(x^*)$$

- Therefore

$$\underline{f(x') + \lambda g(x') < f(x^*) + \lambda g(x^*)}$$

## Rearranging Optimization Problems

Suppose we have two functions  $f : \mathbf{R}^d \rightarrow \mathbf{R}$  and  $g : \mathbf{R}^d \rightarrow \mathbf{R}$ . Now consider the following optimization problem:

$$\min_{x \in \mathbf{R}^d} f(x) + g(x).$$



## Rearranging Optimization Problems

Suppose we have two functions  $f : \mathbf{R}^d \rightarrow \mathbf{R}$  and  $g : \mathbf{R}^d \rightarrow \mathbf{R}$ . Now consider the following optimization problem:

$$\min_{x \in \mathbf{R}^d} f(x) + g(x).$$

This is an unconstrained optimization problem. Let's also consider the following constrained optimization problem:

$$\begin{array}{ll} \text{minimize} & f(x) + \xi \\ \text{subject to} & \xi \geq g(x). \end{array}$$

# Rearranging Optimization Problems

Suppose we have two functions  $f : \mathbf{R}^d \rightarrow \mathbf{R}$  and  $g : \mathbf{R}^d \rightarrow \mathbf{R}$ . Now consider the following optimization problem:

$$\min_{x \in \mathbf{R}^d} f(x) + g(x).$$

This is an unconstrained optimization problem. Let's also consider the following constrained optimization problem:

$$\begin{array}{ll} \text{minimize} & f(x) + \xi \\ \text{subject to} & \xi \geq g(x). \end{array}$$

Need to go both ways to have equivalent problem...there cannot be a gap

## Rearranging Optimization Problems

Example

$$A = \begin{bmatrix} 5 & 5 & 5 & 5 & 5 \\ 8 & 8 & 1 & 8 & 8 \\ +\infty & +\infty & +\infty & 0 & +\infty \end{bmatrix}$$

# Rearranging Optimization Problems

Example

$$A = \begin{bmatrix} 5 & 5 & 5 & 5 & 5 \\ 8 & 8 & 1 & 8 & 8 \\ +\infty & +\infty & +\infty & 0 & +\infty \end{bmatrix}$$

► We always have

$$\max_j \min_i a_{ij} = d^* \leq p^* = \min_i \max_j a_{ij}.$$

$$p^* = \min_i \max_j a_{ij}$$

$$d^* = \max_j \min_i a_{ij}$$

Example

## Rearranging Optimization Problems

$$A = \begin{bmatrix} 5 & 5 & 5 & 5 & 5 \\ 8 & 8 & 1 & 8 & 8 \\ +\infty & +\infty & +\infty & 0 & +\infty \end{bmatrix}$$

► We always have

$$\max_j \min_i a_{ij} = d^* \leq p^* = \min_i \max_j a_{ij}.$$

$$p^* = \min_i \max_j a_{ij}$$

$$d^* = \max_j \min_i a_{ij}$$

because

$$d^* = a_{i_d j_d} \leq a_{i_p j_d} \leq a_{i_p j_p} = p^*.$$

# Rearranging Optimization Problems

Example

$$A = \begin{bmatrix} 5 & 5 & 5 & 5 & 5 \\ 8 & 8 & 1 & 8 & 8 \\ +\infty & +\infty & +\infty & 0 & +\infty \end{bmatrix}$$

► We always have

$$\max_j \min_i a_{ij} = d^* \leq p^* = \min_i \max_j a_{ij}.$$

because

$$d^* = a_{i_d j_d} \leq a_{i_p j_d} \leq a_{i_p j_p} = p^*.$$

Primal Problem and Dual Problem may not be equal meaning you cannot switch max and min

$$p^* = \min_i \max_j a_{ij}$$

$$d^* = \max_j \min_i a_{ij}$$

# Rearranging Optimization Problems

Example

- ▶ Minimize  $x + y$  subject to constraint  $x^2 + y^2 = 1$

## Rearranging Optimization Problems

- ▶ Minimize  $x + y$  subject to constraint  $x^2 + y^2 = 1$
- ▶ We can combine the objective and constraint into a single function called the Lagrangian

$$L(x, y, \lambda) = x + y + \lambda(x^2 + y^2 - 1)$$



## Rearranging Optimization Problems

- ▶ Minimize  $x + y$  subject to constraint  $x^2 + y^2 = 1$
- ▶ We can combine the objective and constraint into a single function called the Lagrangian

$$L(x, y, \lambda) = x + y + \lambda(x^2 + y^2 - 1)$$

- ▶ Take derivative to find minimum

$$\nabla L = \begin{pmatrix} \frac{\partial L}{\partial x} \\ \frac{\partial L}{\partial y} \\ \frac{\partial L}{\partial \lambda} \end{pmatrix} = \begin{pmatrix} 1 + 2\lambda x \\ 1 + 2\lambda y \\ x^2 + y^2 - 1 \end{pmatrix}$$

## Rearranging Optimization Problems

- ▶ Minimize  $x + y$  subject to constraint  $x^2 + y^2 = 1$
- ▶ We can combine the objective and constraint into a single function called the Lagrangian

$$L(x, y, \lambda) = x + y + \lambda(x^2 + y^2 - 1)$$

- ▶ Take derivative to find minimum

$$\nabla L = \begin{pmatrix} \frac{\partial L}{\partial x} \\ \frac{\partial L}{\partial y} \\ \frac{\partial L}{\partial \lambda} \end{pmatrix} = \begin{pmatrix} 1 + 2\lambda x \\ 1 + 2\lambda y \\ x^2 + y^2 - 1 \end{pmatrix}$$

- ▶ Solutions at

$$\left(\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2}\right) \text{ and } \left(-\frac{\sqrt{2}}{2}, -\frac{\sqrt{2}}{2}\right).$$

# Rearranging Optimization Problems

Example

- ▶ Minimize  $x + y$  subject to constraint  $x^2 + y^2 \leq 1$

Example

## Rearranging Optimization Problems

- ▶ Minimize  $x + y$  subject to constraint  $x^2 + y^2 \leq 1$
- ▶ We can combine the objective and constraint into a single function called the Lagrangian

$$L(x, y, \lambda) = x + y + \lambda (1 - (x^2 + y^2))$$

- ▶ Here  $\lambda > 0$ .

## Example

# Rearranging Optimization Problems

- ▶ Minimize  $x + y$  subject to constraint  $x^2 + y^2 \leq 1$
- ▶ We can combine the objective and constraint into a single function called the Lagrangian

$$L(x, y, \lambda) = x + y + \lambda (1 - (x^2 + y^2))$$

- ▶ Here  $\lambda > 0$

Take max over the dual  
variables and min over the  
primal variables

## Example

# Rearranging Optimization Problems

- ▶ Minimize  $x + y$  subject to constraint  $x^2 + y^2 \leq 1$
- ▶ We can combine the objective and constraint into a single function called the Lagrangian

$$L(x, y, \lambda) = x + y + \lambda (1 - (x^2 + y^2))$$

- ▶ Here  $\lambda > 0$ .

So penalization form and  
constraint form are definitely  
the same!

Example

## Rearranging Optimization Problems

- ▶ Minimize  $x + y$  subject to constraint  $x^2 + y^2 \leq 1$
- ▶ We can combine the objective and constraint into a single function called the Lagrangian

$$L(x, y, \lambda) = x + y + \lambda (1 - (x^2 + y^2))$$

- ▶ Here  $\lambda > 0$ . Take derivative to find minimum

$$\nabla L = \begin{pmatrix} \frac{\partial L}{\partial x} \\ \frac{\partial L}{\partial y} \end{pmatrix} = \begin{pmatrix} 1 + 2\lambda x \\ 1 + 2\lambda y \end{pmatrix}$$

## Rearranging Optimization Problems

- ▶ Minimize  $x + y$  subject to constraint  $x^2 + y^2 \leq 1$
- ▶ We can combine the objective and constraint into a single function called the Lagrangian

$$L(x, y, \lambda) = x + y + \lambda (1 - (x^2 + y^2))$$

- ▶ Here  $\lambda > 0$ . Take derivative to find minimum

$$\nabla L = \begin{pmatrix} \frac{\partial L}{\partial x} \\ \frac{\partial L}{\partial y} \end{pmatrix} = \begin{pmatrix} 1 + 2\lambda x \\ 1 + 2\lambda y \end{pmatrix}$$

- ▶ Solutions at

$$\left(\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2}\right) \text{ and } \left(-\frac{\sqrt{2}}{2}, -\frac{\sqrt{2}}{2}\right).$$



## Example

# Rearranging Optimization Problems

- ▶ Minimize  $x + y$  subject to constraint  $x^2 + y^2 \leq 1$
- ▶ We can combine the objective and constraint into a single function called the Lagrangian

$$L(x, y, \lambda) = x + y + \lambda (1 - (x^2 + y^2))$$

- ▶ Here  $\lambda > 0$ . Take derivative to find minimum

$$\nabla L = \begin{pmatrix} \frac{\partial L}{\partial x} \\ \frac{\partial L}{\partial y} \end{pmatrix} = \begin{pmatrix} 1 + 2\lambda x \\ 1 + 2\lambda y \end{pmatrix}$$

- ▶ Solutions at

$$\left(\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2}\right) \text{ and } \left(-\frac{\sqrt{2}}{2}, -\frac{\sqrt{2}}{2}\right).$$

Note that solutions not unique. Is the objective convex? Is the objective concave?

## Example

# Rearranging Optimization Problems

- ▶ Minimize  $x + y$  subject to constraint  $x^2 + y^2 \leq 1$
- ▶ We can combine the objective and constraint into a single function called the Lagrangian

$$L(x, y, \lambda) = x + y + \lambda (1 - (x^2 + y^2))$$

- ▶ Here  $\lambda > 0$ . Take derivative to find minimum

$$\nabla L = \begin{pmatrix} \frac{\partial L}{\partial x} \\ \frac{\partial L}{\partial y} \end{pmatrix} = \begin{pmatrix} 1 - 2\lambda x \\ 1 - 2\lambda y \end{pmatrix}$$

- ▶ Solutions at

$$\left(\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2}\right) \text{ and } \left(-\frac{\sqrt{2}}{2}, -\frac{\sqrt{2}}{2}\right).$$

At the minimizer the constraint is satisfied...this is example of complementary slackness

# Summary

- ▶ Support Vector Machines
  - ▶ Hard Margin: Only applies to linearly separable data
  - ▶ Soft Margin: Allows for slack variables. Useful for outlier detection
- ▶ Subgradients
  - ▶ Useful for convex functions. Takes any vector with properties of gradient.
  - ▶ Subgradient Descent variant of Gradient Descent
- ▶ Rearranging Optimization Problems
  - ▶ Combine objective and constraint
  - ▶ Switch order of minimization / maximization
  - ▶ Lagrangians, First Order Conditions and Complementary Slackness