

# Machine Learning

A. G. Schwing

University of Illinois at Urbana-Champaign, 2018

## L5: Optimization Dual

## **Goals of this lecture**

## **Goals of this lecture**

- Constrained optimization

## **Goals of this lecture**

- Constrained optimization
- Understanding duality for optimization

## **Goals of this lecture**

- Constrained optimization
- Understanding duality for optimization

## **Reading Material**

## **Goals of this lecture**

- Constrained optimization
- Understanding duality for optimization

## **Reading Material**

- S. Boyd and L. Vandenberghe; Convex Optimization; Chapter 5

Optimization problems that we have seen so far:



Optimization problems that we have seen so far:

- Linear Regression

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left( y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

Optimization problems that we have seen so far:

- Linear Regression

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left( y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

- Logistic Regression

$$\min_{\mathbf{w}} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \log \left( 1 + \exp(-y^{(i)} \mathbf{w}^\top \phi(x^{(i)})) \right)$$

Optimization problems that we have seen so far:

- Linear Regression

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left( y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

- Logistic Regression

$$\min_{\mathbf{w}} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \log \left( 1 + \exp(-y^{(i)} \mathbf{w}^\top \phi(x^{(i)})) \right)$$

Finding optimum:

Optimization problems that we have seen so far:

- Linear Regression

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left( y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

- Logistic Regression

$$\min_{\mathbf{w}} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \log \left( 1 + \exp(-y^{(i)} \mathbf{w}^\top \phi(x^{(i)})) \right)$$

Finding optimum:

Analytically computable optimum vs. gradient descent

## **The Problem more generally:**

## The Problem more generally:

$$\begin{array}{ll} \min_{\mathbf{w}} & f_0(\mathbf{w}) \\ \text{s.t.} & f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C\} \end{array}$$

## The Problem more generally:

$$\begin{array}{ll} \min_{\mathbf{w}} & f_0(\mathbf{w}) \\ \text{s.t.} & f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C\} \end{array}$$

**Solution:**

## The Problem more generally:

$$\begin{array}{ll} \min_{\mathbf{w}} & f_0(\mathbf{w}) \\ \text{s.t.} & f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C\} \end{array}$$

### Solution:

Solution  $\mathbf{w}^*$  has smallest value  $f_0(\mathbf{w}^*)$  among all values that satisfy constraints

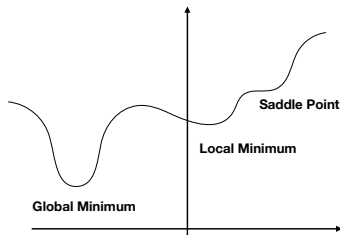


## The Problem more generally:

$$\begin{array}{ll}\min_{\mathbf{w}} & f_0(\mathbf{w}) \\ \text{s.t.} & f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C\}\end{array}$$

### Solution:

Solution  $\mathbf{w}^*$  has smallest value  $f_0(\mathbf{w}^*)$  among all values that satisfy constraints



## **Original/Primal Problem:**

## Original/Primal Problem:

$$\begin{array}{ll}\min_{\mathbf{w}} & f_0(\mathbf{w}) \\ \text{s.t.} & f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C_1\} \\ & h_i(\mathbf{w}) = 0 \quad \forall i \in \{1, \dots, C_2\}\end{array}$$

## Original/Primal Problem:

$$\begin{array}{ll}\min_{\mathbf{w}} & f_0(\mathbf{w}) \\ \text{s.t.} & f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C_1\} \\ & h_i(\mathbf{w}) = 0 \quad \forall i \in \{1, \dots, C_2\}\end{array}$$

How to optimize this?

When can we find the optimum?

When can we find the optimum?

- Least squares, linear and convex programs can be solved efficiently and reliably

When can we find the optimum?

- Least squares, linear and convex programs can be solved efficiently and reliably
- General optimization problems are very difficult to solve

When can we find the optimum?

- Least squares, linear and convex programs can be solved efficiently and reliably
- General optimization problems are very difficult to solve
- Often compromise between accuracy and computation time



What's the form of 'least squares,' 'linear,' and 'convex' programs?

- Least squares program

What's the form of 'least squares,' 'linear,' and 'convex' programs?

- Least squares program

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left( y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

What's the form of 'least squares,' 'linear,' and 'convex' programs?

- Least squares program

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left( y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

- Linear program

What's the form of 'least squares,' 'linear,' and 'convex' programs?

- Least squares program

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left( y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

- Linear program

$$\min_{\mathbf{w}} \mathbf{c}^\top \mathbf{w} \quad \text{s.t.} \quad \mathbf{A}\mathbf{w} \leq \mathbf{b}$$

What's the form of 'least squares,' 'linear,' and 'convex' programs?

- Least squares program

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left( y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

- Linear program

$$\min_{\mathbf{w}} \mathbf{c}^\top \mathbf{w} \quad \text{s.t.} \quad \mathbf{A}\mathbf{w} \leq \mathbf{b}$$

- Convex program

What's the form of 'least squares,' 'linear,' and 'convex' programs?

- Least squares program

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left( y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

- Linear program

$$\min_{\mathbf{w}} \mathbf{c}^\top \mathbf{w} \quad \text{s.t.} \quad \mathbf{A}\mathbf{w} \leq \mathbf{b}$$

- Convex program when all  $f_i$  **convex** (generalizes the above)

$$\min_{\mathbf{w}} f_0(\mathbf{w}) \quad \text{s.t.} \quad f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C\}$$

# Optimality of convex optimization

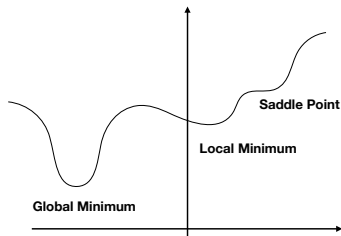
## Optimality of convex optimization

- A point  $\mathbf{w}^*$  is locally optimal if  $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$  in a neighborhood of  $\mathbf{w}^*$ ; globally optimal if  $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$



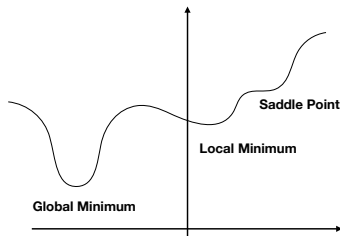
## Optimality of convex optimization

- A point  $\mathbf{w}^*$  is locally optimal if  $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$  in a neighborhood of  $\mathbf{w}^*$ ; globally optimal if  $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$



## Optimality of convex optimization

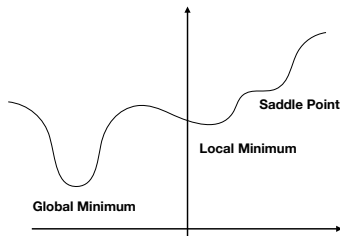
- A point  $\mathbf{w}^*$  is locally optimal if  $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$  in a neighborhood of  $\mathbf{w}^*$ ; globally optimal if  $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$



**For convex problems global optimality follows directly from local optimality.**

## Optimality of convex optimization

- A point  $\mathbf{w}^*$  is locally optimal if  $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$  in a neighborhood of  $\mathbf{w}^*$ ; globally optimal if  $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$

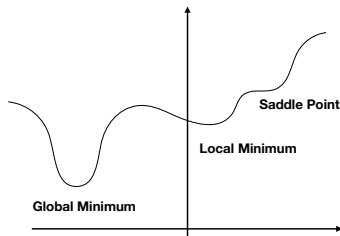


**For convex problems global optimality follows directly from local optimality.**

- For a local minimum of  $f$ ,  $\nabla f(\mathbf{w}^*) = 0$

## Optimality of convex optimization

- A point  $\mathbf{w}^*$  is locally optimal if  $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$  in a neighborhood of  $\mathbf{w}^*$ ; globally optimal if  $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$

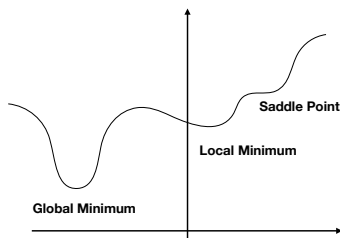


**For convex problems global optimality follows directly from local optimality.**

- For a local minimum of  $f$ ,  $\nabla f(\mathbf{w}^*) = 0$
- If  $f$  convex, then  $\nabla f(\mathbf{w}^*) = 0$  sufficient for global optimality

## Optimality of convex optimization

- A point  $\mathbf{w}^*$  is locally optimal if  $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$  in a neighborhood of  $\mathbf{w}^*$ ; globally optimal if  $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$



**For convex problems global optimality follows directly from local optimality.**

- For a local minimum of  $f$ ,  $\nabla f(\mathbf{w}^*) = 0$
- If  $f$  convex, then  $\nabla f(\mathbf{w}^*) = 0$  sufficient for global optimality

This makes convex optimization special!

Algorithms to search for the optimum?

## Descent methods

$$\min_{\mathbf{w}} f(\mathbf{w})$$

Intuition

## Descent methods

$$\min_{\mathbf{w}} f(\mathbf{w})$$

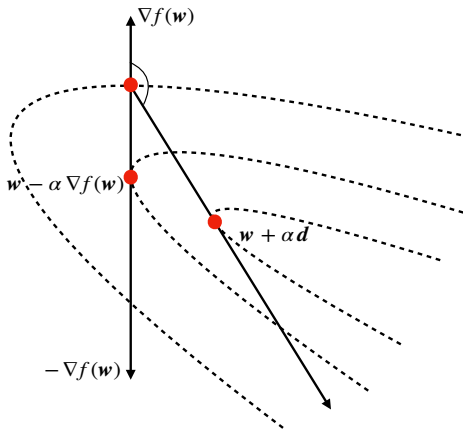
Intuition (find a stationary point with  $\nabla f(\mathbf{w}) = 0$ )



## Descent methods

$$\min_{\mathbf{w}} f(\mathbf{w})$$

Intuition (find a stationary point with  $\nabla f(\mathbf{w}) = 0$ )



## Iterative algorithm

## Iterative algorithm

- Start with some guess  $\mathbf{w}$

## Iterative algorithm

- Start with some guess  $\mathbf{w}$
- Iterate  $k = 1, 2, 3, \dots$

## Iterative algorithm

- Start with some guess  $\mathbf{w}$
- Iterate  $k = 1, 2, 3, \dots$ 
  - ▶ Select direction  $\mathbf{d}_k$  and stepsize  $\alpha_k$

## Iterative algorithm

- Start with some guess  $\mathbf{w}$
- Iterate  $k = 1, 2, 3, \dots$ 
  - ▶ Select direction  $\mathbf{d}_k$  and stepsize  $\alpha_k$
  - ▶  $\mathbf{w} \leftarrow \mathbf{w} + \alpha_k \mathbf{d}_k$

## Iterative algorithm

- Start with some guess  $\mathbf{w}$
- Iterate  $k = 1, 2, 3, \dots$ 
  - ▶ Select direction  $\mathbf{d}_k$  and stepsize  $\alpha_k$
  - ▶  $\mathbf{w} \leftarrow \mathbf{w} + \alpha_k \mathbf{d}_k$
  - ▶ Check whether we should stop (e.g., if  $\nabla f(\mathbf{w}) \approx 0$ )

## Iterative algorithm

- Start with some guess  $\mathbf{w}$
- Iterate  $k = 1, 2, 3, \dots$ 
  - ▶ Select direction  $\mathbf{d}_k$  and stepsize  $\alpha_k$
  - ▶  $\mathbf{w} \leftarrow \mathbf{w} + \alpha_k \mathbf{d}_k$
  - ▶ Check whether we should stop (e.g., if  $\nabla f(\mathbf{w}) \approx 0$ )

Descent direction  $\mathbf{d}_k$  satisfies



## Iterative algorithm

- Start with some guess  $\mathbf{w}$
- Iterate  $k = 1, 2, 3, \dots$ 
  - ▶ Select direction  $\mathbf{d}_k$  and stepsize  $\alpha_k$
  - ▶  $\mathbf{w} \leftarrow \mathbf{w} + \alpha_k \mathbf{d}_k$
  - ▶ Check whether we should stop (e.g., if  $\nabla f(\mathbf{w}) \approx 0$ )

Descent direction  $\mathbf{d}_k$  satisfies  $\nabla f(\mathbf{w})^\top \mathbf{d}_k < 0$

How to select direction:

How to select direction:

- Steepest descent:  $\mathbf{d}_k = -\nabla f(\mathbf{w}_k)$

How to select direction:

- Steepest descent:  $\mathbf{d}_k = -\nabla f(\mathbf{w}_k)$
- Scaled gradient:  $\mathbf{d}_k = -\mathbf{D}_k \nabla f(\mathbf{w}_k)$  for  $\mathbf{D}_k \succ 0$

How to select direction:

- Steepest descent:  $\mathbf{d}_k = -\nabla f(\mathbf{w}_k)$
- Scaled gradient:  $\mathbf{d}_k = -\mathbf{D}_k \nabla f(\mathbf{w}_k)$  for  $\mathbf{D}_k \succ 0$ 
  - ▶ E.g., Newton's method:  $\mathbf{D}_k = [\nabla^2 f(\mathbf{w}_k)]^{-1}$

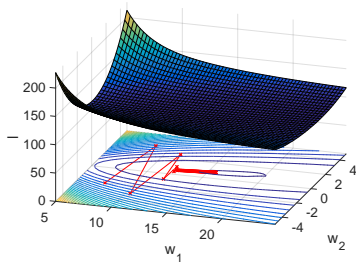
How to select direction:

- Steepest descent:  $\mathbf{d}_k = -\nabla f(\mathbf{w}_k)$
- Scaled gradient:  $\mathbf{d}_k = -\mathbf{D}_k \nabla f(\mathbf{w}_k)$  for  $\mathbf{D}_k \succ 0$ 
  - ▶ E.g., Newton's method:  $\mathbf{D}_k = [\nabla^2 f(\mathbf{w}_k)]^{-1}$
- Gradient with momentum

# Gradient with momentum

# Gradient with momentum

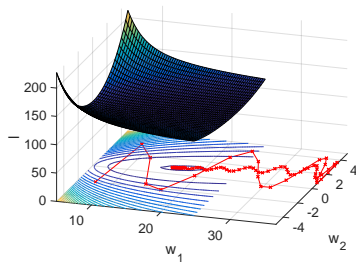
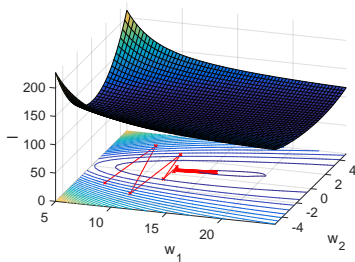
Intuition:





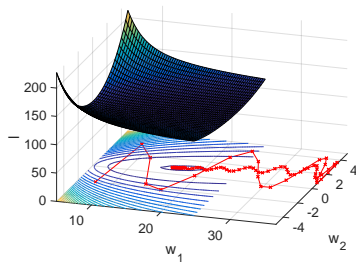
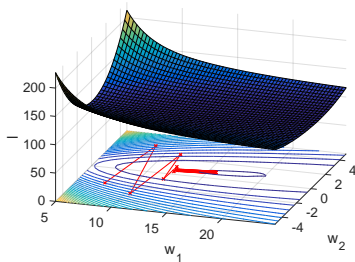
# Gradient with momentum

Intuition:



# Gradient with momentum

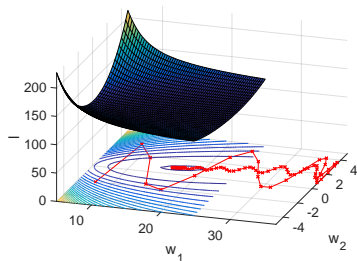
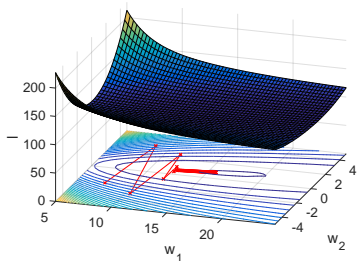
Intuition:



Video

# Gradient with momentum

Intuition:

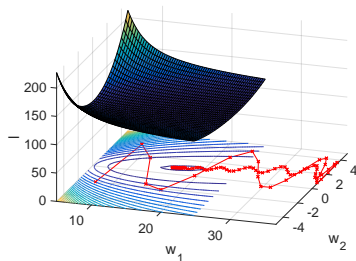
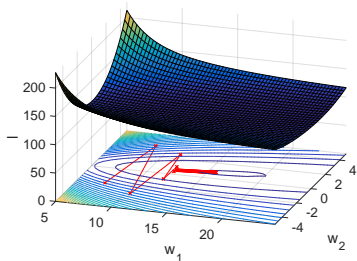


Video

- Polyak's method (aka heavy-ball)

# Gradient with momentum

Intuition:



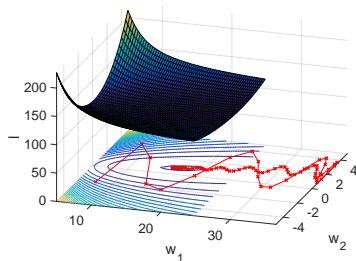
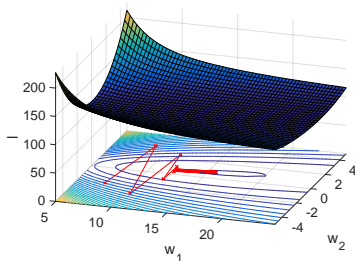
Video

- Polyak's method (aka heavy-ball)

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k) + \beta_k (\mathbf{w}_k - \mathbf{w}_{k-1})$$

# Gradient with momentum

Intuition:



Video

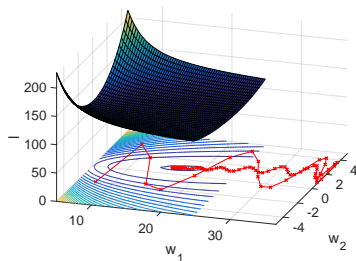
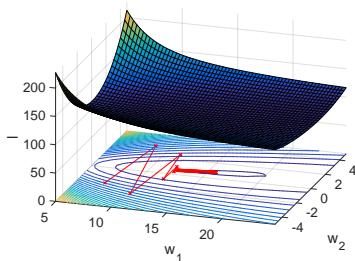
- Polyak's method (aka heavy-ball)

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k) + \beta_k (\mathbf{w}_k - \mathbf{w}_{k-1})$$

- Momentum method in deep learning

# Gradient with momentum

Intuition:



Video

- Polyak's method (aka heavy-ball)

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k) + \beta_k (\mathbf{w}_k - \mathbf{w}_{k-1})$$

- Momentum method in deep learning

$$\mathbf{v}_{k+1} = \beta \mathbf{v}_k + \nabla f(\mathbf{w}_k)$$

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \mathbf{v}_{k+1}$$

How to select stepsize:

How to select stepsize:

- Exact:  $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{w}_k + \alpha \mathbf{d}_k)$



How to select stepsize:

- Exact:  $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{w}_k + \alpha \mathbf{d}_k)$
- Constant:  $\alpha_k = 1/L$  (for suitable  $L$ )

How to select stepsize:

- Exact:  $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{w}_k + \alpha \mathbf{d}_k)$
- Constant:  $\alpha_k = 1/L$  (for suitable  $L$ )
- Diminishing:  $\alpha_k \rightarrow 0$  but  $\sum_k \alpha_k = \infty$  (e.g.,  $\alpha_k = 1/k$ )

How to select stepsize:

- Exact:  $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{w}_k + \alpha \mathbf{d}_k)$
- Constant:  $\alpha_k = 1/L$  (for suitable  $L$ )
- Diminishing:  $\alpha_k \rightarrow 0$  but  $\sum_k \alpha_k = \infty$  (e.g.,  $\alpha_k = 1/k$ )
- Armijo Rule

Recall the structure of our optimization problems:

Recall the structure of our optimization problems:

$$\min_{\mathbf{w}} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \ell(y_i, F(x^{(i)}, \mathbf{w}))$$

Recall the structure of our optimization problems:

$$\min_{\mathbf{w}} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \ell(y_i, F(x^{(i)}, \mathbf{w}))$$

- So far we didn't consider the time for computing the gradient

Recall the structure of our optimization problems:

$$\min_{\mathbf{w}} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \ell(y_i, F(x^{(i)}, \mathbf{w}))$$

- So far we didn't consider the time for computing the gradient
- Iteration complexity is linear in the number of samples  $|\mathcal{D}|$

Recall the structure of our optimization problems:

$$\min_{\mathbf{w}} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \ell(y_i, F(x^{(i)}, \mathbf{w}))$$

- So far we didn't consider the time for computing the gradient
- Iteration complexity is linear in the number of samples  $|\mathcal{D}|$
- A large dataset makes gradient computation slow



Recall the structure of our optimization problems:

$$\min_{\mathbf{w}} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \ell(y_i, F(x^{(i)}, \mathbf{w}))$$

- So far we didn't consider the time for computing the gradient
- Iteration complexity is linear in the number of samples  $|\mathcal{D}|$
- A large dataset makes gradient computation slow

How to deal with this?

# Stochastic gradient descent

## **Stochastic gradient descent**

Consider a subset of samples and approximate the gradient based on this batch of data.

## Stochastic gradient descent

Consider a subset of samples and approximate the gradient based on this batch of data.

- Select a subset of samples  $\mathcal{B}_k$

## Stochastic gradient descent

Consider a subset of samples and approximate the gradient based on this batch of data.

- Select a subset of samples  $\mathcal{B}_k$
- Gradient update using approximation

$$\nabla f(\mathbf{w}) \approx \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{B}_k} \nabla \ell(y^{(i)}, F(x^{(i)}, \mathbf{w}))$$

How about constraints?

## Original/Primal Problem:

$$\begin{array}{ll}\min_{\mathbf{w}} & f_0(\mathbf{w}) \\ \text{s.t.} & f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C_1\} \\ & h_i(\mathbf{w}) = 0 \quad \forall i \in \{1, \dots, C_2\}\end{array}$$

## Original/Primal Problem:

$$\begin{array}{ll}\min_{\mathbf{w}} & f_0(\mathbf{w}) \\ \text{s.t.} & f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C_1\} \\ & h_i(\mathbf{w}) = 0 \quad \forall i \in \{1, \dots, C_2\}\end{array}$$

## Lagrangian



## Original/Primal Problem:

$$\begin{aligned} \min_{\mathbf{w}} \quad & f_0(\mathbf{w}) \\ \text{s.t.} \quad & f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C_1\} \\ & h_i(\mathbf{w}) = 0 \quad \forall i \in \{1, \dots, C_2\} \end{aligned}$$

## Lagrangian

$$L(\mathbf{w}, \lambda, \nu) = f_0(\mathbf{w}) + \sum_{i=1}^{C_1} \lambda_i f_i(\mathbf{w}) + \sum_{i=1}^{C_2} \nu_i h_i(\mathbf{w})$$

## Original/Primal Problem:

$$\begin{aligned} \min_{\mathbf{w}} \quad & f_0(\mathbf{w}) \\ \text{s.t.} \quad & f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C_1\} \\ & h_i(\mathbf{w}) = 0 \quad \forall i \in \{1, \dots, C_2\} \end{aligned}$$

## Lagrangian

$$L(\mathbf{w}, \lambda, \nu) = f_0(\mathbf{w}) + \sum_{i=1}^{C_1} \lambda_i f_i(\mathbf{w}) + \sum_{i=1}^{C_2} \nu_i h_i(\mathbf{w})$$

- $\lambda_i$  are Lagrange multiplier associated with inequality constraints

## Original/Primal Problem:

$$\begin{array}{ll}\min_{\mathbf{w}} & f_0(\mathbf{w}) \\ \text{s.t.} & f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C_1\} \\ & h_i(\mathbf{w}) = 0 \quad \forall i \in \{1, \dots, C_2\}\end{array}$$

## Lagrangian

$$L(\mathbf{w}, \lambda, \nu) = f_0(\mathbf{w}) + \sum_{i=1}^{C_1} \lambda_i f_i(\mathbf{w}) + \sum_{i=1}^{C_2} \nu_i h_i(\mathbf{w})$$

- $\lambda_i$  are Lagrange multiplier associated with inequality constraints
- $\nu_i$  are Lagrange multiplier associated with equality constraints

## Properties of Lagrangian:

$$L(\mathbf{w}, \lambda, \nu) = f_0(\mathbf{w}) + \sum_{i=1}^{C_1} \lambda_i f_i(\mathbf{w}) + \sum_{i=1}^{C_2} \nu_i h_i(\mathbf{w})$$

## Properties of Lagrangian:

$$L(\mathbf{w}, \lambda, \nu) = f_0(\mathbf{w}) + \sum_{i=1}^{C_1} \lambda_i f_i(\mathbf{w}) + \sum_{i=1}^{C_2} \nu_i h_i(\mathbf{w})$$

If  $\hat{\mathbf{w}}$  feasible and  $\lambda_i \geq 0 \ \forall i$  then

## Properties of Lagrangian:

$$L(\mathbf{w}, \lambda, \nu) = f_0(\mathbf{w}) + \sum_{i=1}^{C_1} \lambda_i f_i(\mathbf{w}) + \sum_{i=1}^{C_2} \nu_i h_i(\mathbf{w})$$

If  $\hat{\mathbf{w}}$  feasible and  $\lambda_i \geq 0 \ \forall i$  then

$$f_0(\hat{\mathbf{w}}) \geq L(\hat{\mathbf{w}}, \lambda, \nu) \geq \min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}, \lambda, \nu) = g(\lambda, \nu) \quad \forall \lambda \geq 0, \nu$$

$$f_0(\mathbf{w}^*) \geq g(\lambda, \nu) \quad \forall \lambda \geq 0, \nu$$

$\mathcal{W}$  denotes all the constraints that are not part of the Lagrangian  
(larger than feasible set)

## Properties of Lagrangian:

$$L(\mathbf{w}, \lambda, \nu) = f_0(\mathbf{w}) + \sum_{i=1}^{C_1} \lambda_i f_i(\mathbf{w}) + \sum_{i=1}^{C_2} \nu_i h_i(\mathbf{w})$$

If  $\hat{\mathbf{w}}$  feasible and  $\lambda_i \geq 0 \ \forall i$  then

$$f_0(\hat{\mathbf{w}}) \geq L(\hat{\mathbf{w}}, \lambda, \nu) \geq \min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}, \lambda, \nu) = g(\lambda, \nu) \quad \forall \lambda \geq 0, \nu$$

$$f_0(\mathbf{w}^*) \geq g(\lambda, \nu) \quad \forall \lambda \geq 0, \nu$$

$\mathcal{W}$  denotes all the constraints that are not part of the Lagrangian  
(larger than feasible set)

Dual Program:

$$\max_{\lambda, \nu} g(\lambda, \nu) \quad \text{s.t.} \quad \lambda_i \geq 0 \quad \forall i$$

## **Recipe for computing dual program:**



## **Recipe for computing dual program:**

- Bring primal program into standard form

## **Recipe for computing dual program:**

- Bring primal program into standard form
- Assign Lagrange multipliers to a suitable set of constraints

## **Recipe for computing dual program:**

- Bring primal program into standard form
- Assign Lagrange multipliers to a suitable set of constraints
- Subsume all other constraints in  $\mathcal{W}$

## Recipe for computing dual program:

- Bring primal program into standard form
- Assign Lagrange multipliers to a suitable set of constraints
- Subsume all other constraints in  $\mathcal{W}$
- Write down the Lagrangian  $L$

## Recipe for computing dual program:

- Bring primal program into standard form
- Assign Lagrange multipliers to a suitable set of constraints
- Subsume all other constraints in  $\mathcal{W}$
- Write down the Lagrangian  $L$
- Minimize Lagrangian w.r.t. primal variables s.t.  $\mathbf{w} \in \mathcal{W}$

## **Examples:** Linear Program

## Examples: Linear Program

$$\min_{\mathbf{w}} \mathbf{c}^T \mathbf{w} \quad \text{s.t.} \quad \mathbf{A}\mathbf{w} \leq \mathbf{b}$$

## Examples: Linear Program

$$\min_{\mathbf{w}} \mathbf{c}^T \mathbf{w} \quad \text{s.t.} \quad \mathbf{A}\mathbf{w} \leq \mathbf{b}$$

Lagrangian:



## Examples: Linear Program

$$\min_{\mathbf{w}} \mathbf{c}^\top \mathbf{w} \quad \text{s.t.} \quad \mathbf{A}\mathbf{w} \leq \mathbf{b}$$

Lagrangian: ( $\lambda \geq 0$ )

$$L() = \mathbf{c}^\top \mathbf{w} + \lambda^\top (\mathbf{A}\mathbf{w} - \mathbf{b}) = (\mathbf{c} + \mathbf{A}^\top \lambda)^\top \mathbf{w} - \mathbf{b}^\top \lambda$$

## Examples: Linear Program

$$\min_{\mathbf{w}} \mathbf{c}^\top \mathbf{w} \quad \text{s.t.} \quad \mathbf{A}\mathbf{w} \leq \mathbf{b}$$

Lagrangian: ( $\lambda \geq 0$ )

$$L() = \mathbf{c}^\top \mathbf{w} + \lambda^\top (\mathbf{A}\mathbf{w} - \mathbf{b}) = (\mathbf{c} + \mathbf{A}^\top \lambda)^\top \mathbf{w} - \mathbf{b}^\top \lambda$$

Minimizing Lagrangian w.r.t. primal variables:

## Examples: Linear Program

$$\min_{\mathbf{w}} \mathbf{c}^\top \mathbf{w} \quad \text{s.t.} \quad \mathbf{A}\mathbf{w} \leq \mathbf{b}$$

Lagrangian: ( $\lambda \geq 0$ )

$$L() = \mathbf{c}^\top \mathbf{w} + \lambda^\top (\mathbf{A}\mathbf{w} - \mathbf{b}) = (\mathbf{c} + \mathbf{A}^\top \lambda)^\top \mathbf{w} - \mathbf{b}^\top \lambda$$

Minimizing Lagrangian w.r.t. primal variables:

$$\min_{\mathbf{w}} L() = \begin{cases} -\mathbf{b}^\top \lambda & \mathbf{A}^\top \lambda + \mathbf{c} = 0 \\ -\infty & \text{otherwise} \end{cases}$$

## Examples: Linear Program

$$\min_{\mathbf{w}} \mathbf{c}^\top \mathbf{w} \quad \text{s.t.} \quad \mathbf{A}\mathbf{w} \leq \mathbf{b}$$

Lagrangian: ( $\lambda \geq 0$ )

$$L() = \mathbf{c}^\top \mathbf{w} + \lambda^\top (\mathbf{A}\mathbf{w} - \mathbf{b}) = (\mathbf{c} + \mathbf{A}^\top \lambda)^\top \mathbf{w} - \mathbf{b}^\top \lambda$$

Minimizing Lagrangian w.r.t. primal variables:

$$\min_{\mathbf{w}} L() = \begin{cases} -\mathbf{b}^\top \lambda & \mathbf{A}^\top \lambda + \mathbf{c} = 0 \\ -\infty & \text{otherwise} \end{cases}$$

Dual Program:

## Examples: Linear Program

$$\min_{\mathbf{w}} \mathbf{c}^\top \mathbf{w} \quad \text{s.t.} \quad \mathbf{A}\mathbf{w} \leq \mathbf{b}$$

Lagrangian: ( $\lambda \geq 0$ )

$$L() = \mathbf{c}^\top \mathbf{w} + \lambda^\top (\mathbf{A}\mathbf{w} - \mathbf{b}) = (\mathbf{c} + \mathbf{A}^\top \lambda)^\top \mathbf{w} - \mathbf{b}^\top \lambda$$

Minimizing Lagrangian w.r.t. primal variables:

$$\min_{\mathbf{w}} L() = \begin{cases} -\mathbf{b}^\top \lambda & \mathbf{A}^\top \lambda + \mathbf{c} = 0 \\ -\infty & \text{otherwise} \end{cases}$$

Dual Program:

$$\max_{\lambda \geq 0} -\mathbf{b}^\top \lambda \quad \text{s.t.} \quad \mathbf{A}^\top \lambda + \mathbf{c} = 0,$$

## **Examples:** Logistic Regression

## Examples: Logistic Regression

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \log(1 + \exp(-y^{(i)} \mathbf{w}^\top \phi(x^{(i)})))$$

## Examples: Logistic Regression

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \log(1 + \exp(-y^{(i)} \mathbf{w}^\top \phi(x^{(i)})))$$

Reformulate:



## Examples: Logistic Regression

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \log(1 + \exp(-y^{(i)} \mathbf{w}^\top \phi(x^{(i)})))$$

Reformulate:

$$\min_{\mathbf{w}, z^{(i)}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \log(1 + \exp(-z^{(i)})) \quad \text{s.t.} \quad z^{(i)} = y^{(i)} \mathbf{w}^\top \phi(x^{(i)})$$

## Examples: Logistic Regression

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \log(1 + \exp(-y^{(i)} \mathbf{w}^\top \phi(x^{(i)})))$$

Reformulate:

$$\min_{\mathbf{w}, z^{(i)}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \log(1 + \exp(-z^{(i)})) \quad \text{s.t.} \quad z^{(i)} = y^{(i)} \mathbf{w}^\top \phi(x^{(i)})$$

Lagrangian:

## Examples: Logistic Regression

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \log(1 + \exp(-y^{(i)} \mathbf{w}^\top \phi(x^{(i)})))$$

Reformulate:

$$\min_{\mathbf{w}, z^{(i)}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \log(1 + \exp(-z^{(i)})) \quad \text{s.t.} \quad z^{(i)} = y^{(i)} \mathbf{w}^\top \phi(x^{(i)})$$

Lagrangian:

$$\begin{aligned} L() &= \frac{C}{2} \|\mathbf{w}\|_2^2 - \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \lambda^{(i)} y^{(i)} \mathbf{w}^\top \phi(x^{(i)}) \\ &\quad + \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left[ \log(1 + \exp(-z^{(i)})) + \lambda^{(i)} z^{(i)} \right] \end{aligned}$$

Minimize Lagrangian w.r.t. primal variables ( $\min_{\mathbf{w}, \mathbf{z}} L()$ ):

Minimize Lagrangian w.r.t. primal variables ( $\min_{\mathbf{w}, z} L()$ ):

$$\frac{\partial L}{\partial \mathbf{w}} :$$

$$\frac{\partial L}{\partial z^{(i)}} :$$

Minimize Lagrangian w.r.t. primal variables ( $\min_{\mathbf{w}, \mathbf{z}} L()$ ):

$$\frac{\partial L}{\partial \mathbf{w}} : \quad C\mathbf{w} = \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \lambda^{(i)} y^{(i)} \phi(x^{(i)})$$

$$\frac{\partial L}{\partial z^{(i)}} :$$

Minimize Lagrangian w.r.t. primal variables ( $\min_{\mathbf{w}, z} L()$ ):

$$\frac{\partial L}{\partial \mathbf{w}} : \quad C\mathbf{w} = \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \lambda^{(i)} y^{(i)} \phi(x^{(i)})$$

$$\frac{\partial L}{\partial z^{(i)}} : \quad \lambda^{(i)} = \frac{\exp(-z^{(i)})}{1 + \exp(-z^{(i)})} \quad \implies \lambda^{(i)} \geq 0$$

Minimize Lagrangian w.r.t. primal variables ( $\min_{\mathbf{w}, z} L()$ ):

$$\frac{\partial L}{\partial \mathbf{w}} : \quad C\mathbf{w} = \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \lambda^{(i)} y^{(i)} \phi(x^{(i)})$$

$$\begin{aligned} \frac{\partial L}{\partial z^{(i)}} : \quad \lambda^{(i)} &= \frac{\exp(-z^{(i)})}{1 + \exp(-z^{(i)})} \quad \implies \lambda^{(i)} \geq 0 \\ &\implies z^{(i)} = \log \frac{1 - \lambda^{(i)}}{\lambda^{(i)}} \quad \implies \lambda^{(i)} \leq 1 \end{aligned}$$



Minimize Lagrangian w.r.t. primal variables ( $\min_{\mathbf{w}, z} L()$ ):

$$\frac{\partial L}{\partial \mathbf{w}} : \quad C\mathbf{w} = \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \lambda^{(i)} y^{(i)} \phi(x^{(i)})$$

$$\begin{aligned} \frac{\partial L}{\partial z^{(i)}} : \quad \lambda^{(i)} &= \frac{\exp(-z^{(i)})}{1 + \exp(-z^{(i)})} \quad \implies \lambda^{(i)} \geq 0 \\ &\implies z^{(i)} = \log \frac{1 - \lambda^{(i)}}{\lambda^{(i)}} \quad \implies \lambda^{(i)} \leq 1 \end{aligned}$$

Dual function:

Minimize Lagrangian w.r.t. primal variables ( $\min_{\mathbf{w}, z} L()$ ):

$$\frac{\partial L}{\partial \mathbf{w}} : \quad C\mathbf{w} = \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} \lambda^{(i)} y^{(i)} \phi(\mathbf{x}^{(i)})$$

$$\begin{aligned} \frac{\partial L}{\partial z^{(i)}} : \quad \lambda^{(i)} &= \frac{\exp(-z^{(i)})}{1 + \exp(-z^{(i)})} \quad \implies \lambda^{(i)} \geq 0 \\ &\implies z^{(i)} = \log \frac{1 - \lambda^{(i)}}{\lambda^{(i)}} \quad \implies \lambda^{(i)} \leq 1 \end{aligned}$$

Dual function:

$$g(\lambda) = -\frac{1}{2C} \left\| \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} \lambda^{(i)} y^{(i)} \phi(\mathbf{x}^{(i)}) \right\|_2^2 + \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} H(\lambda^{(i)})$$

with binary entropy  $H(\lambda^{(i)})$

Minimize Lagrangian w.r.t. primal variables ( $\min_{\mathbf{w}, z} L()$ ):

$$\frac{\partial L}{\partial \mathbf{w}} : \quad C\mathbf{w} = \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \lambda^{(i)} y^{(i)} \phi(x^{(i)})$$

$$\begin{aligned} \frac{\partial L}{\partial z^{(i)}} : \quad \lambda^{(i)} &= \frac{\exp(-z^{(i)})}{1 + \exp(-z^{(i)})} \quad \implies \lambda^{(i)} \geq 0 \\ &\implies z^{(i)} = \log \frac{1 - \lambda^{(i)}}{\lambda^{(i)}} \quad \implies \lambda^{(i)} \leq 1 \end{aligned}$$

Dual function:

$$g(\lambda) = -\frac{1}{2C} \left\| \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \lambda^{(i)} y^{(i)} \phi(x^{(i)}) \right\|_2^2 + \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} H(\lambda^{(i)})$$

with binary entropy  $H(\lambda^{(i)})$

Dual program:

Minimize Lagrangian w.r.t. primal variables ( $\min_{\mathbf{w}, z} L()$ ):

$$\frac{\partial L}{\partial \mathbf{w}} : \quad C\mathbf{w} = \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} \lambda^{(i)} y^{(i)} \phi(\mathbf{x}^{(i)})$$

$$\begin{aligned} \frac{\partial L}{\partial z^{(i)}} : \quad \lambda^{(i)} &= \frac{\exp(-z^{(i)})}{1 + \exp(-z^{(i)})} \implies \lambda^{(i)} \geq 0 \\ &\implies z^{(i)} = \log \frac{1 - \lambda^{(i)}}{\lambda^{(i)}} \implies \lambda^{(i)} \leq 1 \end{aligned}$$

Dual function:

$$g(\lambda) = -\frac{1}{2C} \left\| \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} \lambda^{(i)} y^{(i)} \phi(\mathbf{x}^{(i)}) \right\|_2^2 + \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} H(\lambda^{(i)})$$

with binary entropy  $H(\lambda^{(i)})$

Dual program:

$$\max_{\lambda} g(\lambda) \quad \text{s.t.} \quad 0 \leq \lambda^{(i)} \leq 1 \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}$$

Instead of optimizing the primal we can optimize the dual and convert the result

Instead of optimizing the primal we can optimize the dual and convert the result

Why is this useful?

Instead of optimizing the primal we can optimize the dual and convert the result

Why is this useful?

- Sometimes less constraints

Instead of optimizing the primal we can optimize the dual and convert the result

Why is this useful?

- Sometimes less constraints
- Sometimes easier to optimize



Instead of optimizing the primal we can optimize the dual and convert the result

Why is this useful?

- Sometimes less constraints
- Sometimes easier to optimize
- Sometimes interesting insights

Instead of optimizing the primal we can optimize the dual and convert the result

Why is this useful?

- Sometimes less constraints
- Sometimes easier to optimize
- Sometimes interesting insights
- Sometimes lower bounds

## Properties of Dual Program

$$\max_{\lambda, \nu} g(\lambda, \nu) \quad \text{s.t.} \quad \lambda_i \geq 0 \quad \forall i$$

## Properties of Dual Program

$$\max_{\lambda, \nu} g(\lambda, \nu) \quad \text{s.t.} \quad \lambda_i \geq 0 \quad \forall i$$

- May only have simple constraints if at all

## Properties of Dual Program

$$\max_{\lambda, \nu} g(\lambda, \nu) \quad \text{s.t.} \quad \lambda_i \geq 0 \quad \forall i$$

- May only have simple constraints if at all
- Can be used for sensitivity analysis

## Properties of Dual Program

$$\max_{\lambda, \nu} g(\lambda, \nu) \quad \text{s.t.} \quad \lambda_i \geq 0 \quad \forall i$$

- May only have simple constraints if at all
- Can be used for sensitivity analysis
- Lower-bounds the optimal primal value

## Properties of Dual Program

$$\max_{\lambda, \nu} g(\lambda, \nu) \quad \text{s.t.} \quad \lambda_i \geq 0 \quad \forall i$$

- May only have simple constraints if at all
- Can be used for sensitivity analysis
- Lower-bounds the optimal primal value
- Dual Program is always concave:

## Properties of Dual Program

$$\max_{\lambda, \nu} g(\lambda, \nu) \quad \text{s.t.} \quad \lambda_i \geq 0 \quad \forall i$$

- May only have simple constraints if at all
- Can be used for sensitivity analysis
- Lower-bounds the optimal primal value
- Dual Program is always concave:

$$g(\lambda, \nu) = \min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}, \lambda, \nu) := f_0(\mathbf{w}) + \sum_{i=1}^{C_1} \lambda_i f_i(\mathbf{w}) + \sum_{i=1}^{C_2} \nu_i h_i(\mathbf{w})$$



## Properties of Dual Program

$$\max_{\lambda, \nu} g(\lambda, \nu) \quad \text{s.t.} \quad \lambda_i \geq 0 \quad \forall i$$

- May only have simple constraints if at all
- Can be used for sensitivity analysis
- Lower-bounds the optimal primal value
- Dual Program is always concave:

$$g(\lambda, \nu) = \min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}, \lambda, \nu) := f_0(\mathbf{w}) + \sum_{i=1}^{C_1} \lambda_i f_i(\mathbf{w}) + \sum_{i=1}^{C_2} \nu_i h_i(\mathbf{w})$$

- ▶ Pointwise minimum

## Properties of Dual Program

$$\max_{\lambda, \nu} g(\lambda, \nu) \quad \text{s.t.} \quad \lambda_i \geq 0 \quad \forall i$$

- May only have simple constraints if at all
- Can be used for sensitivity analysis
- Lower-bounds the optimal primal value
- Dual Program is always concave:

$$g(\lambda, \nu) = \min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w}, \lambda, \nu) := f_0(\mathbf{w}) + \sum_{i=1}^{C_1} \lambda_i f_i(\mathbf{w}) + \sum_{i=1}^{C_2} \nu_i h_i(\mathbf{w})$$

- ▶ Pointwise minimum
- ▶ Affine functions in  $\lambda, \nu$

Weak duality:

$$f(\mathbf{w}^*) \geq g(\lambda^*, \nu^*)$$

Weak duality:

$$f(\mathbf{w}^*) \geq g(\lambda^*, \nu^*)$$

- Always holds (for convex and non-convex problems)

Weak duality:

$$f(\mathbf{w}^*) \geq g(\lambda^*, \nu^*)$$

- Always holds (for convex and non-convex problems)
- Can be used to find nontrivial lower bounds

Weak duality:

$$f(\mathbf{w}^*) \geq g(\lambda^*, \nu^*)$$

- Always holds (for convex and non-convex problems)
- Can be used to find nontrivial lower bounds

Strong duality:

$$f(\mathbf{w}^*) = g(\lambda^*, \nu^*)$$

Weak duality:

$$f(\mathbf{w}^*) \geq g(\lambda^*, \nu^*)$$

- Always holds (for convex and non-convex problems)
- Can be used to find nontrivial lower bounds

Strong duality:

$$f(\mathbf{w}^*) = g(\lambda^*, \nu^*)$$

- Does not hold in general

Weak duality:

$$f(\mathbf{w}^*) \geq g(\lambda^*, \nu^*)$$

- Always holds (for convex and non-convex problems)
- Can be used to find nontrivial lower bounds

Strong duality:

$$f(\mathbf{w}^*) = g(\lambda^*, \nu^*)$$

- Does not hold in general
- (Usually) holds for convex problems



## Quiz:

## Quiz:

- What to do before computing the Lagrangian?

## Quiz:

- What to do before computing the Lagrangian?
- How to obtain the dual program?

## Quiz:

- What to do before computing the Lagrangian?
- How to obtain the dual program?
- Why duality?

## Important topics of this lecture

## Important topics of this lecture

- Lagrangian

## Important topics of this lecture

- Lagrangian
- Dual program

## Important topics of this lecture

- Lagrangian
- Dual program

## Up next:

- Support vector machines