

# Convex Optimization: Old Tricks for New Problems

Ryota Tomioka



The University of Tokyo

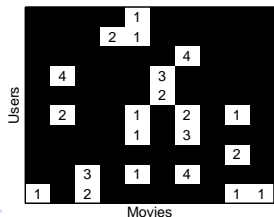
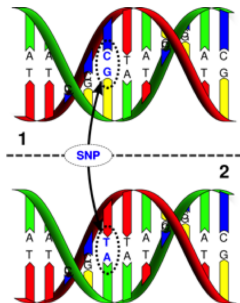
2012-08-15 @ DTU PhD Summer Course

# Introduction

Why care about convex optimization (and sparsity)?

# Why do we care about optimization — sparse estimation

- High dimensional problems (dimension  $\gg$  # samples)
  - ▶ Bioinformatics (microarray , SNP analysis , etc)
  - ▶ Text-mining (POS tagging , )
  - ▶ Magnetic resonance imaging — compressed sensing
- Structure inference
  - ▶ Collaborative filtering — **low-rank structure**
  - ▶ Graphical model inference — **sparse graph structure**



# Ex. 1: SNP (single nucleotide polymorphism) analysis

$\mathbf{x}_i$ : input (SNP) ,  $y_i = 1$ : has the illness ,  $y_i = -1$ : healthy

Goal: Infer the association from genetic variability  $\mathbf{x}_i$  to the illness  $y_i$ .

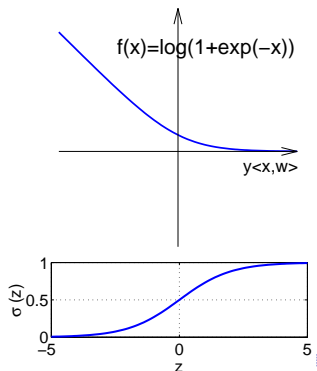
Logistic regression

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \underbrace{\sum_{i=1}^m \log(1 + \exp(-y_i \langle \mathbf{x}_i, \mathbf{w} \rangle))}_{\text{data-fit}} + \underbrace{\lambda \|\mathbf{w}\|_1}_{\text{Regularization}}$$

- E.g., # SNPs  $n = 500,000$ , # subjects  $m = 5,000$
- MAP estimation with the *logistic loss*  $f$ .

$$\log(1 + e^{-y^z}) = -\log P(Y = y|z)$$

$$\text{where } P(Y = +1|z) = \frac{e^z}{1+e^z}.$$

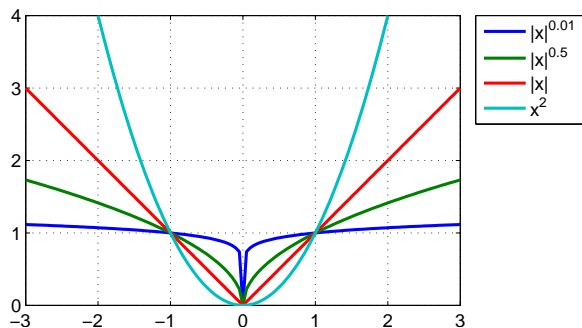


# L1-regularization and sparsity

- Best convex approximation of  $\|\mathbf{w}\|_0$ .

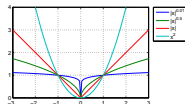
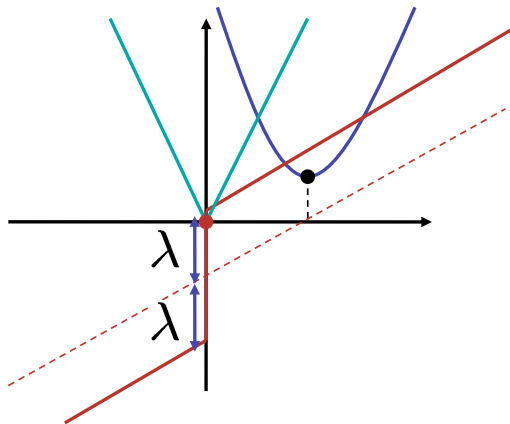
# L1-regularization and sparsity

- Best convex approximation of  $\|\mathbf{w}\|_0$ .



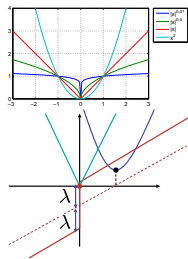
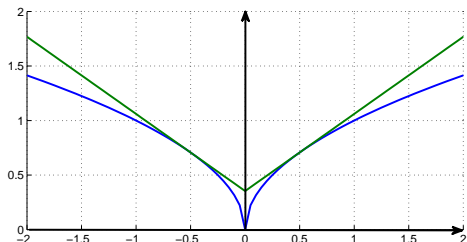
# L1-regularization and sparsity

- Best convex approximation of  $\|\mathbf{w}\|_0$ .
- Threshold occurs for finite  $\lambda$ .



# L1-regularization and sparsity

- Best convex approximation of  $\|\mathbf{w}\|_0$ .
- Threshold occurs for finite  $\lambda$ .
- Non-convex cases ( $p < 1$ ) can be solved by re-weighted L1 minimization





## Ex. 2: Compressed sensing [Candes, Romberg, & Tao 06]

Signal (MRI image) recovery from (noisy) low-dimensional measurements.

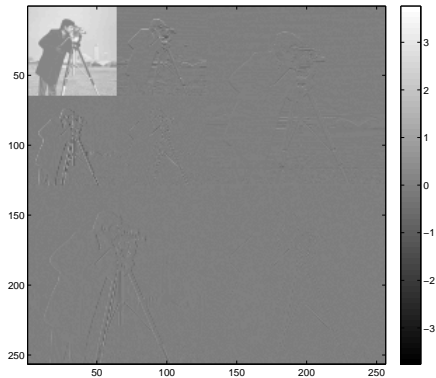
$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{\Omega} \mathbf{w}\|_2^2 + \lambda \|\mathbf{\Phi} \mathbf{w}\|_1$$

- $\mathbf{y}$ : Noisy signal
- $\mathbf{w}$ : Original signal
- $\mathbf{\Omega}: \mathbb{R}^n \rightarrow \mathbb{R}^m$ : Observation matrix (random, fourier transform)
- $\mathbf{\Phi}$ : Transformation s.t. the original signal is sparse

NB: If  $\mathbf{\Phi}^{-1}$  exists, we can solve instead

$$\underset{\tilde{\mathbf{w}} \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{A} \tilde{\mathbf{w}}\|_2^2 + \lambda \|\tilde{\mathbf{w}}\|_1,$$

where  $\mathbf{A} = \mathbf{\Omega} \mathbf{\Phi}^{-1}$ .



## Ex. 3: Estimation of a low-rank matrix [Fazel+ 01; Srebro+ 05]

Goal: Recover a low-rank matrix  $\mathbf{X}$  from partial (noisy) measurement  $\mathbf{Y}$

$$\underset{\mathbf{X}}{\text{minimize}} \quad \frac{1}{2} \|\Omega(\mathbf{X} - \mathbf{Y})\|^2 + \lambda \|\mathbf{X}\|_{S_1}$$

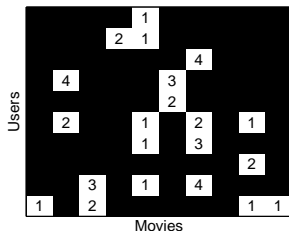
where  $\|\mathbf{X}\|_{S_1} := \sum_{j=1}^r \sigma_j(\mathbf{X})$  (Schatten 1-norm)

Aka trace norm, nuclear norm

⇒ Linear sum of singular values

⇒ Sparsity in the SV spectrum

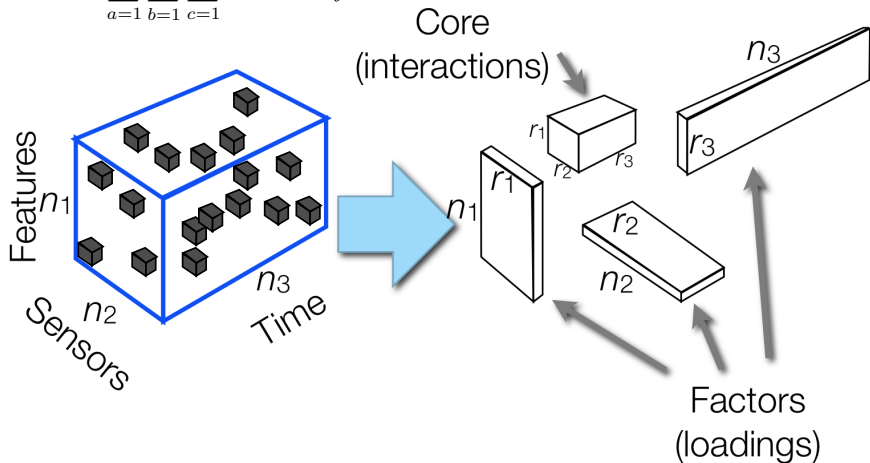
⇒ Low-rank



## Ex. 4: Low-rank tensor completion [Tomioka+11]

$$X_{ijk} = \sum_{a=1}^{r_1} \sum_{b=1}^{r_2} \sum_{c=1}^{r_3} C_{abc} U_{ia}^{(1)} U_{jb}^{(2)} U_{kc}^{(c)}$$

Tucker decomposition



# Simple vs. structured sparse estimation problems

- Simple sparse estimation problem

$$\underset{\mathbf{w}}{\text{minimize}} \quad L(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

- ▶ SNP analysis
- ▶ Compressed sensing with  $\Phi^{-1}$  (e.g., wavelet)
- ▶ Collaborative filtering (matrix completion)

- Structured sparse estimation problem

$$\underset{\mathbf{w}}{\text{minimize}} \quad L(\mathbf{w}) + \lambda \|\Phi \mathbf{w}\|_1$$

- ▶ Compressed sensing without  $\Phi^{-1}$  (e.g., total variation)
- ▶ Low-rank tensor completion

# Common criticisms

- Convex optimization is another developed field (and it is boring). We can just use it as a black box.
  - ▶ Yes, but we can do much better by *knowing the structure of our problems*.

# Common criticisms

- Convex optimization is another developed field (and it is boring). We can just use it as a black box.
  - ▶ Yes, but we can do much better by *knowing the structure of our problems*.
- Convexity is too restrictive.
  - ▶ Convexity depends on parametrization. A seemingly non-convex problem could be *reformulated into a convex problem*.

# Common criticisms

- Convex optimization is another developed field (and it is boring). We can just use it as a black box.
  - ▶ Yes, but we can do much better by *knowing the structure of our problems*.
- Convexity is too restrictive.
  - ▶ Convexity depends on parametrization. A seemingly non-convex problem could be *reformulated into a convex problem*.
- I am only interested in making things work.
  - ▶ Yes, convex optimization works. But it can also be used for *analyzing how algorithms perform at the end*.



# Bayesian inference as a convex optimization

$$\begin{aligned} \underset{q}{\text{minimize}} \quad & \underbrace{\mathbb{E}_q[f(w)]}_{\text{average energy}} + \underbrace{\mathbb{E}_q[\log q(w)]}_{\text{entropy}} \\ \text{s.t.} \quad & q(w) \geq 0, \quad \int q(w)dw = 1 \end{aligned}$$

where

$$f(w) = \underbrace{-\log P(D|w)}_{\text{neg. log likelihood}} - \underbrace{\log P(w)}_{\text{neg. log prior}}$$

# Bayesian inference as a convex optimization

$$\begin{aligned} \underset{q}{\text{minimize}} \quad & \underbrace{\mathbb{E}_q[f(w)]}_{\text{average energy}} + \underbrace{\mathbb{E}_q[\log q(w)]}_{\text{entropy}} \\ \text{s.t.} \quad & q(w) \geq 0, \quad \int q(w)dw = 1 \end{aligned}$$

where

$$\begin{aligned} f(w) &= \underbrace{-\log P(D|w)}_{\text{neg. log likelihood}} - \underbrace{\log P(w)}_{\text{neg. log prior}} \\ \Rightarrow \quad q(w) &= \frac{1}{Z} e^{-f(w)} \quad (\text{Bayesian posterior}) \end{aligned}$$

# Bayesian inference as a convex optimization

$$\begin{aligned} \underset{q}{\text{minimize}} \quad & \underbrace{\mathbb{E}_q[f(w)]}_{\text{average energy}} + \underbrace{\mathbb{E}_q[\log q(w)]}_{\text{entropy}} \\ \text{s.t.} \quad & q(w) \geq 0, \quad \int q(w)dw = 1 \end{aligned}$$

where

$$\begin{aligned} f(w) &= \underbrace{-\log P(D|w)}_{\text{neg. log likelihood}} - \underbrace{\log P(w)}_{\text{neg. log prior}} \\ \Rightarrow \quad q(w) &= \frac{1}{Z} e^{-f(w)} \quad (\text{Bayesian posterior}) \end{aligned}$$

Inner approximations



- Variational Bayes
- Empirical Bayes

# Bayesian inference as a convex optimization

$$\underset{q}{\text{minimize}} \quad \underbrace{\mathbb{E}_q[f(w)]}_{\text{average energy}} + \underbrace{\mathbb{E}_q[\log q(w)]}_{\text{entropy}}$$

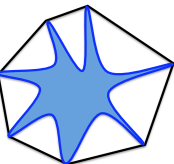
$$\text{s.t.} \quad q(w) \geq 0, \quad \int q(w)dw = 1$$

where

$$f(w) = \underbrace{-\log P(D|w)}_{\text{neg. log likelihood}} - \underbrace{\log P(w)}_{\text{neg. log prior}}$$

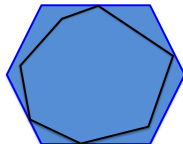
$$\Rightarrow q(w) = \frac{1}{Z} e^{-f(w)} \quad (\text{Bayesian posterior})$$

Inner approximations



- Variational Bayes
- Empirical Bayes

Outer approximations



- Belief propagation

See Wainwright & Jordan 08.

# Overview

- Convex optimization basics
  - ▶ Convex sets
  - ▶ Convex function
  - ▶ Conditions that guarantee convexity
  - ▶ Convex optimization problem
- Looking into more structures
  - ▶ Proximity operators
  - ▶ Conjugate duality and dual ascent
  - ▶ Augmented Lagrangian and ADMM

## References:



Boyd & Vandenberghe. (2004) Convex optimization.



Bertsekas (1999) Nonlinear Programming.



Rockafellar (1970) Convex Analysis.

Moreau (1965) Proximité et dualité dans un espace Hilbertien.

# Convexity

## Learning objectives

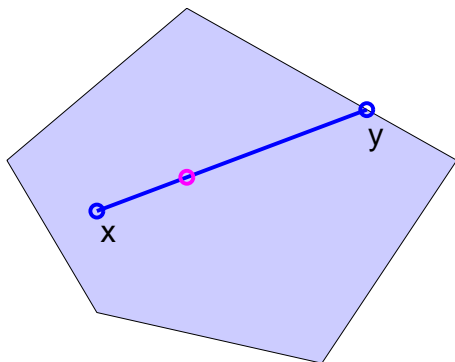
- Convex sets
- Convex function
- Conditions that guarantee convexity
- Convex optimization problem

# Convex set

A subset  $V \subseteq \mathbb{R}^n$  is a **convex set**

$\Leftrightarrow$  line segment between two arbitrary points  $\mathbf{x}, \mathbf{y} \in V$  is included in  $V$ ;  
that is,

$$\forall \mathbf{x}, \mathbf{y} \in V, \forall \lambda \in [0, 1], \quad \lambda \mathbf{x} + (1 - \lambda) \mathbf{y} \in V.$$



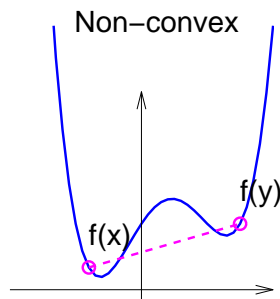
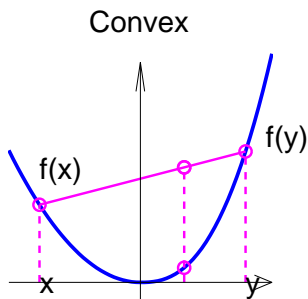
# Convex function

A function  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  is a **convex function**

$\Leftrightarrow$  the function  $f$  is below any line segment between two points on  $f$ ; that is,

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \forall \lambda \in [0, 1], \quad f((1 - \lambda)\mathbf{x} + \lambda\mathbf{y}) \leq (1 - \lambda)f(\mathbf{x}) + \lambda f(\mathbf{y})$$

(Jensen's inequality)



Johan Jensen  
1859 – 1925

NB: when the strict inequality  $<$  holds,  $f$  is called **strictly convex**.

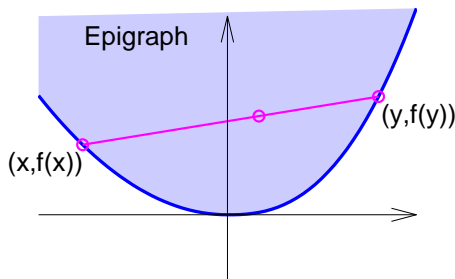


# Convex function

A function  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  is a **convex function**

$\Leftrightarrow$  the **epigraph of  $f$**  is a **convex set**; that is

$V_f := \{(t, \mathbf{x}) : (t, \mathbf{x}) \in \mathbb{R}^{n+1}, t \geq f(\mathbf{x})\}$  is convex.

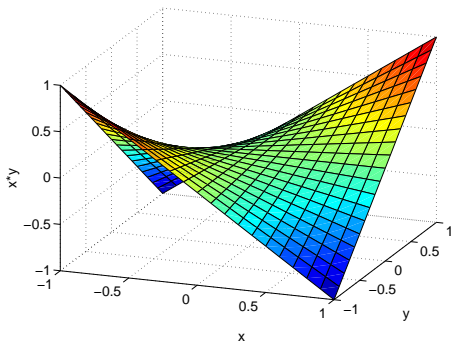


## Jointly convex

- A function  $f(x, y)$  can be convex wrt  $x$  ( $y$ ) for any fixed  $y$  ( $x$ ), respectively, but can fail to be convex for  $x$  and  $y$  simultaneously.

## Jointly convex

- A function  $f(x, y)$  can be convex wrt  $x$  ( $y$ ) for any fixed  $y$  ( $x$ ), respectively, but can fail to be convex for  $x$  and  $y$  simultaneously.



$f(x, y)$  is convex  $\Rightarrow (\not\Leftarrow)$   $f(x, y)$  is convex for  $x$  and  $y$  individually

- To be more explicit, we sometimes say **jointly** convex.

# Why do we allow infinity?

- $f(x) = 1/x$  is convex for  $x > 0$ .

$$f(x) = \begin{cases} 1/x & \text{if } x > 0, \\ +\infty & \text{otherwise.} \end{cases}$$

and we can forget about the domain.

# Why do we allow infinity?

- $f(x) = 1/x$  is convex for  $x > 0$ .

$$f(x) = \begin{cases} 1/x & \text{if } x > 0, \\ +\infty & \text{otherwise.} \end{cases}$$

and we can forget about the domain.

- The **indicator function**  $\delta_C(\mathbf{x})$  of a set  $C$ :

$$\delta_C(\mathbf{x}) = \begin{cases} 0 & \text{if } \mathbf{x} \in C, \\ +\infty & \text{otherwise.} \end{cases}$$

Is this a convex function? (consider the epigraph)

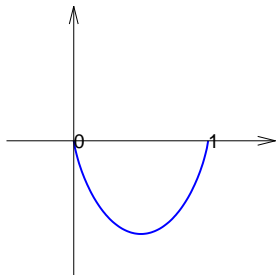
## Condition #1: Hessian

Hessian  $\nabla^2 f(\mathbf{x})$  is positive semidefinite (if  $f$  is differentiable)

### Examples

- (Negative) entropy is a convex function.

$$f(p) = \sum_{i=1}^n p_i \log p_i,$$



## Condition #1: Hessian

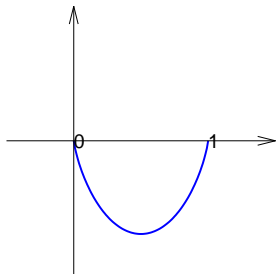
Hessian  $\nabla^2 f(\mathbf{x})$  is positive semidefinite (if  $f$  is differentiable)

### Examples

- (Negative) entropy is a convex function.

$$f(p) = \sum_{i=1}^n p_i \log p_i,$$

$$\nabla^2 f(p) = \text{diag}(1/p_1, \dots, 1/p_n) \succeq 0.$$



## Condition #1: Hessian

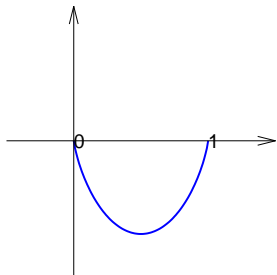
Hessian  $\nabla^2 f(\mathbf{x})$  is positive semidefinite (if  $f$  is differentiable)

### Examples

- (Negative) entropy is a convex function.

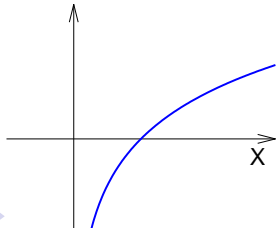
$$f(p) = \sum_{i=1}^n p_i \log p_i,$$

$$\nabla^2 f(p) = \text{diag}(1/p_1, \dots, 1/p_n) \succeq 0.$$



- log determinant is a *concave* ( $-f$  is convex) function

$$f(\mathbf{X}) = \log |\mathbf{X}| \quad (\mathbf{X} \succeq 0),$$
$$\nabla^2 f(\mathbf{X}) = -\mathbf{X}^{-\top} \otimes \mathbf{X}^{-1} \preceq 0$$



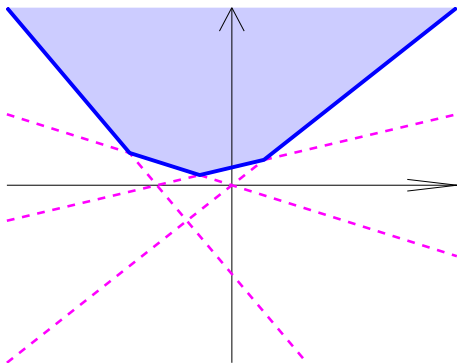


## Condition #2: Maximum over convex functions

Maximum over convex functions  $\{f_j(\mathbf{x})\}_{j=1}^{\infty}$

$$f(\mathbf{x}) := \max_j f_j(\mathbf{x}) \quad (f_j(\mathbf{x}) \text{ is convex for all } j)$$

is convex.



The same as saying “intersection of convex sets is a convex set”

## Condition #2: Maximum over convex functions

Maximum over convex functions  $\{f(\mathbf{x}; \alpha) : \alpha \in \mathbb{R}^m\}$

$$f(\mathbf{x}) := \sup_{\alpha \in \mathbb{R}^m} f(\mathbf{x}; \alpha)$$

is convex.

### Example

- Quadratic over linear is a convex function

$$f(x, y) = \sup_{\alpha \in \mathbb{R}} \left( -\frac{\alpha^2}{2}x + \alpha y \right) \quad (x > 0)$$

## Condition #2: Maximum over convex functions

Maximum over convex functions  $\{f(x; \alpha) : \alpha \in \mathbb{R}^m\}$

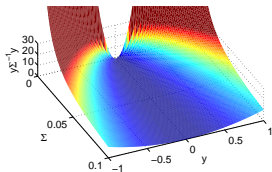
$$f(\mathbf{x}) := \sup_{\alpha \in \mathbb{R}^m} f(\mathbf{x}; \alpha)$$

is convex.

### Example

- Quadratic over linear is a convex function

$$\begin{aligned} f(x, y) &= \sup_{\alpha \in \mathbb{R}} \left( -\frac{\alpha^2}{2}x + \alpha y \right) \quad (x > 0) \\ &= \frac{y^2}{2x} \end{aligned}$$



## Condition #2: Maximum over convex functions

Maximum over convex functions  $\{f(x; \alpha) : \alpha \in \mathbb{R}^m\}$

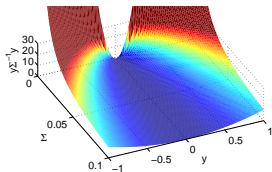
$$f(\mathbf{x}) := \sup_{\alpha \in \mathbb{R}^m} f(\mathbf{x}; \alpha)$$

is convex.

### Example

- Quadratic over linear is a convex function

$$\begin{aligned} f(x, y) &= \sup_{\alpha \in \mathbb{R}} \left( -\frac{\alpha^2}{2}x + \alpha y \right) \quad (x > 0) \\ &= \frac{y^2}{2x} \end{aligned}$$



- Similarly

$$f(\Sigma, \mathbf{y}) = \frac{1}{2} \mathbf{y}^\top \Sigma^{-1} \mathbf{y} \quad (\Sigma \succ 0) \quad \text{is a convex function (show it!)}$$

## Condition #3: Partial minimum

Partial minimum of a convex function  $f(x, y)$

$$f(x) := \min_{y \in \mathbb{R}^n} f(x, y) \quad \text{is convex.}$$

### Examples

- Hierarchical prior minimization

$$f(\mathbf{x}) = \min_{d_1, \dots, d_n \geq 0} \frac{1}{2} \sum_{j=1}^n \left( \frac{x_j^2}{d_j} + \frac{d_j^p}{p} \right) \quad (p \geq 1)$$

### Condition #3: Partial minimum

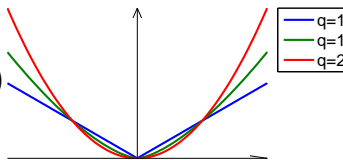
Partial minimum of a convex function  $f(x, y)$

$$f(x) := \min_{y \in \mathbb{R}^n} f(x, y) \quad \text{is convex.}$$

### Examples

- Hierarchical prior minimization

$$\begin{aligned} f(\mathbf{x}) &= \min_{d_1, \dots, d_n \geq 0} \frac{1}{2} \sum_{j=1}^n \left( \frac{x_j^2}{d_j} + \frac{d_j^p}{p} \right) \quad (p \geq 1) \\ &= \frac{1}{q} \sum_{j=1}^n |x_j|^q \quad \left( q = \frac{2p}{1+p} \right) \end{aligned}$$



### Condition #3: Partial minimum

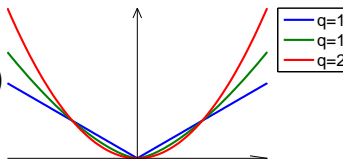
Partial minimum of a convex function  $f(x, y)$

$$f(x) := \min_{y \in \mathbb{R}^n} f(x, y) \quad \text{is convex.}$$

### Examples

- Hierarchical prior minimization

$$\begin{aligned} f(\mathbf{x}) &= \min_{d_1, \dots, d_n \geq 0} \frac{1}{2} \sum_{j=1}^n \left( \frac{x_j^2}{d_j} + \frac{d_j^p}{p} \right) \quad (p \geq 1) \\ &= \frac{1}{q} \sum_{j=1}^n |x_j|^q \quad \left( q = \frac{2p}{1+p} \right) \end{aligned}$$



- Schatten 1- norm (sum of singularvalues)

$$f(\mathbf{X}) = \min_{\Sigma \succeq 0} \frac{1}{2} \left( \text{Tr} \left( \mathbf{X} \Sigma^{-1} \mathbf{X}^\top \right) + \text{Tr} (\Sigma) \right)$$

### Condition #3: Partial minimum

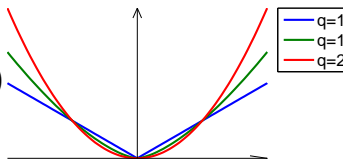
Partial minimum of a convex function  $f(x, y)$

$$f(x) := \min_{y \in \mathbb{R}^n} f(x, y) \quad \text{is convex.}$$

### Examples

- Hierarchical prior minimization

$$\begin{aligned} f(\mathbf{x}) &= \min_{d_1, \dots, d_n \geq 0} \frac{1}{2} \sum_{j=1}^n \left( \frac{x_j^2}{d_j} + \frac{d_j^p}{p} \right) \quad (p \geq 1) \\ &= \frac{1}{q} \sum_{j=1}^n |x_j|^q \quad \left( q = \frac{2p}{1+p} \right) \end{aligned}$$



- Schatten 1- norm (sum of singularvalues)

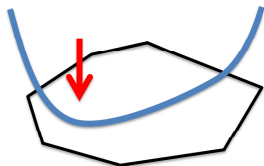
$$f(\mathbf{X}) = \min_{\Sigma \succeq 0} \frac{1}{2} \left( \text{Tr} \left( \mathbf{X} \Sigma^{-1} \mathbf{X}^T \right) + \text{Tr}(\Sigma) \right) = \text{Tr} \left( (\mathbf{X}^T \mathbf{X})^{1/2} \right) = \sum_{j=1}^r \sigma_j(\mathbf{X}).$$



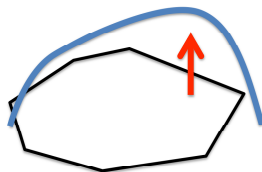
# Convex optimization problem

$f$ : convex function,  $g$ : concave function ( $-g$  is convex),  $C$ : convex set.

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && f(\mathbf{x}), \\ & \text{s.t.} && \mathbf{x} \in C. \end{aligned}$$



$$\begin{aligned} & \underset{\mathbf{y}}{\text{maximize}} && g(\mathbf{y}), \\ & \text{s.t.} && \mathbf{y} \in C. \end{aligned}$$



Why?

- local optimum  $\Rightarrow$  global optimum
- duality (later) can be used to check convergence

$\Rightarrow$  We can be *sure* that we are doing the right thing!

# Coming up next:

- Gradient descent:

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta_t \nabla f(\mathbf{w}^t)$$

- What do we do if we have
    - ▶ Constraints
    - ▶ Non-differentiable terms, like  $\|\mathbf{w}\|_1$
- ⇒ projection/proximity operator

# Proximity operators and iterative shrinkage/thresholding methods

## Learning objectives

- (Projected) gradient method
- Iterative shrinkage/thresholding (IST) method
- Acceleration

# Proximity view on gradient descent

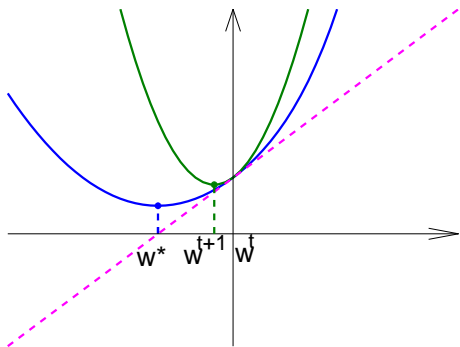
“Linearize and Prox”

$$\begin{aligned}\mathbf{w}^{t+1} &= \operatorname{argmin}_{\mathbf{w}} \left( \nabla f(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right) \\ &= \mathbf{w}^t - \eta_t \nabla f(\mathbf{w}^t)\end{aligned}$$

- Step-size should satisfy  $\eta_t \leq 1/L(f)$ .
- $L(f)$ : the Lipschitz constant

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\| \leq L(f) \|\mathbf{y} - \mathbf{x}\|.$$

- $L(f)$ =upper bound on the maximum eigenvalue of the Hessian



# Constraint minimization problem

- What do we do, if we have a constraint?

$$\begin{aligned} & \underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} && f(\mathbf{w}), \\ & \text{s.t.} && \mathbf{w} \in C. \end{aligned}$$

# Constraint minimization problem

- What do we do, if we have a constraint?

$$\begin{aligned} & \underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} && f(\mathbf{w}), \\ & \text{s.t.} && \mathbf{w} \in C. \end{aligned}$$

- can be equivalently written as

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{w}) + \delta_C(\mathbf{w}),$$

where  $\delta_C(\mathbf{w})$  is the indicator function of the set  $C$ .

# Projected gradient method (Bertsekas 99; Nesterov 03)

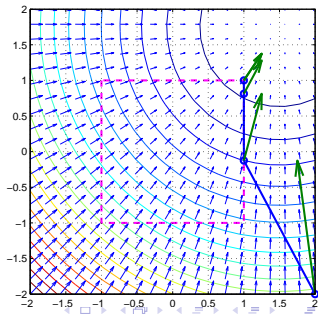
Linearize the objective  $f$ ,  $\delta_C$  is the indicator of the constraint  $C$

$$\begin{aligned}\mathbf{w}^{t+1} &= \operatorname{argmin}_{\mathbf{w}} \left( \nabla f(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) + \delta_C(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 \right) \\ &= \operatorname{argmin}_{\mathbf{w}} \left( \delta_C(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - (\mathbf{w}^t - \eta_t \nabla f(\mathbf{w}^t))\|_2^2 \right) \\ &= \operatorname{proj}_C(\mathbf{w}^t - \eta_t \nabla f(\mathbf{w}^t)).\end{aligned}$$

- Requires  $\eta_t \leq 1/L(f)$ .
- Convergence rate

$$f(\mathbf{w}^k) - f(\mathbf{w}^*) \leq \frac{L(f) \|\mathbf{w}_0 - \mathbf{w}^*\|_2^2}{2k}$$

- Need the projection  $\operatorname{proj}_C$  to be easy to compute



# Ideas for regularized minimization

## Constrained minimization problem

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{w}) + \delta_C(\mathbf{w}).$$

⇒ need to compute the **projection**

$$\mathbf{w}^{t+1} = \underset{\mathbf{w}}{\text{argmin}} \left( \delta_C(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{y}\|_2^2 \right)$$



# Ideas for regularized minimization

## Constrained minimization problem

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{w}) + \delta_C(\mathbf{w}).$$

⇒ need to compute the **projection**

$$\mathbf{w}^{t+1} = \underset{\mathbf{w}}{\text{argmin}} \left( \delta_C(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{y}\|_2^2 \right)$$

## Regularized minimization problem

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

⇒ need to compute the **proximity operator**

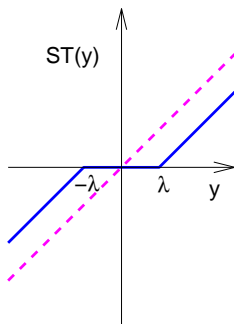
$$\mathbf{w}^{t+1} = \underset{\mathbf{w}}{\text{argmin}} \left( \lambda \|\mathbf{w}\|_1 + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{y}\|_2^2 \right)$$

# Proximal Operator: generalization of projection

$$\text{prox}_g(\mathbf{y}) = \underset{\mathbf{w}}{\text{argmin}} \left( g(\mathbf{w}) + \frac{1}{2} \|\mathbf{w} - \mathbf{y}\|_2^2 \right)$$

- $g = \delta_C$ : Projection onto a convex set  $\text{proj}_C(\mathbf{y})$ .
- $g(\mathbf{w}) = \lambda \|\mathbf{w}\|_1$ : Soft-Threshold

$$\begin{aligned} \text{prox}_\lambda(\mathbf{y}) &= \underset{\mathbf{w}}{\text{argmin}} \left( \lambda \|\mathbf{w}\|_1 + \frac{1}{2} \|\mathbf{w} - \mathbf{y}\|_2^2 \right) \\ &= \begin{cases} y_j + \lambda & (y_j < -\lambda), \\ 0 & (-\lambda \leq y_j \leq \lambda), \\ y_j - \lambda & (y_j > \lambda). \end{cases} \end{aligned}$$



- Prox can be computed easily for a **separable**  $f$ .
- Non-differentiability is OK.

# Exercise

Derive prox operator  $\text{prox}_g$  for

- Ridge regularization

$$g(\mathbf{w}) = \lambda \sum_{j=1}^n w_j^2$$

- Group lasso regularization [Yuan & Lin 2006]

$$g(\mathbf{w}_1, \dots, \mathbf{w}_n) = \lambda \sum_{j=1}^n \|\mathbf{w}_j\|_2$$

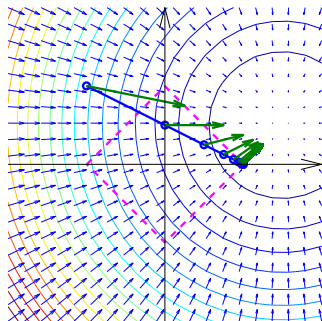
# Iterative Shrinkage Thresholding (IST)

$$\begin{aligned}\mathbf{w}^{t+1} &= \operatorname{argmin}_{\mathbf{w}} \left( \nabla f(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) + \lambda \|\mathbf{w}\|_1 + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 \right) \\ &= \operatorname{argmin}_{\mathbf{w}} \left( \lambda \|\mathbf{w}\|_1 + \frac{1}{2\eta_t} \|\mathbf{w} - (\mathbf{w}^t - \eta_t \nabla f(\mathbf{w}^t))\|_2^2 \right) \\ &= \operatorname{prox}_{\lambda\eta_t}(\mathbf{w}^t - \eta_t \nabla f(\mathbf{w}^t)).\end{aligned}$$

- The same condition for  $\eta_t$ , the same  $O(1/k)$  convergence (Beck & Teboulle 09)

$$f(\mathbf{w}^k) - f(\mathbf{w}^*) \leq \frac{L(f) \|\mathbf{w}_0 - \mathbf{w}^*\|^2}{2k}$$

- If the **Prox operator**  $\operatorname{prox}_\lambda$  is easy, it is simple to implement.
- AKA Forward-Backward Splitting (Lions & Mercier 76)



# IST summary

Solve minimization problem

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

by iteratively computing

$$\mathbf{w}^{t+1} = \text{prox}_{\lambda\eta_t}(\mathbf{w}^t - \eta_t \nabla f(\mathbf{w}^t)),$$

where

$$\text{prox}_{\lambda}(\mathbf{y}) = \underset{\mathbf{w}}{\text{argmin}} \left( \lambda \|\mathbf{w}\|_1 + \frac{1}{2} \|\mathbf{w} - \mathbf{y}\|_2^2 \right).$$

# FISTA: accelerated version of IST (Beck & Teboulle 09; Nesterov 07)

- 1 Initialize  $\mathbf{w}^0$  appropriately,  $\mathbf{z}^1 = \mathbf{w}^0$ ,  $s_1 = 1$ .
- 2 Update  $\mathbf{w}^t$ :

$$\mathbf{w}^t = \text{prox}_{\lambda\eta_t}(\mathbf{z}^t - \eta_t \nabla f(\mathbf{z}^t)).$$

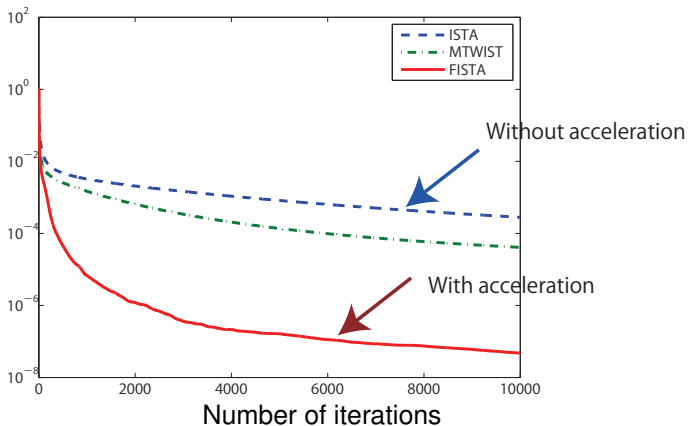
- 3 Update  $\mathbf{z}^t$ :

$$\mathbf{z}^{t+1} = \mathbf{w}^t + \left( \frac{s_t - 1}{s_{t+1}} \right) (\mathbf{w}^t - \mathbf{w}^{t-1}),$$

where  $s_{t+1} = (1 + \sqrt{1 + 4s_t^2})/2$ .

- The same per iteration complexity. Converges as  $O(1/k^2)$ .
- Roughly speaking,  $\mathbf{z}^t$  predicts where the IST step should be computed.

# Effect of acceleration



From Beck & Teboulle 2009 SIAM J. IMAGING SCIENCES

Vol. 2, No. 1, pp. 183-202

# MATLAB Exercise 1: implement an L1 regularized logistic regression via IST

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \underbrace{\sum_{i=1}^m \log(1 + \exp(-y_i \langle \mathbf{x}_i, \mathbf{w} \rangle))}_{\text{data-fit}} + \underbrace{\lambda \sum_{j=1}^n |w_j|}_{\text{Regularization}}$$

Hint: define

$$f_{\ell}(\mathbf{z}) = \sum_{i=1}^m \log(1 + \exp(-z_i)).$$

Then the problem is

$$\underset{\mathbf{w}}{\text{minimize}} \quad f_{\ell}(\mathbf{A}\mathbf{w}) + \lambda \sum_{j=1}^n |w_j| \quad \text{where} \quad \mathbf{A} = \begin{pmatrix} y_1 \mathbf{x}_1^{\top} \\ y_2 \mathbf{x}_2^{\top} \\ \vdots \\ y_m \mathbf{x}_m^{\top} \end{pmatrix}$$



## Some more hints

- 1 Compute the gradient of the loss term

$$\nabla_{\mathbf{w}} f_{\ell}(\mathbf{A}\mathbf{w}) = -\mathbf{A}^{\top} \left( \frac{\exp(-z_i)}{1 + \exp(-z_i)} \right)_{i=1}^m \quad (\mathbf{z} = \mathbf{A}\mathbf{w})$$

- 2 The gradient step becomes

$$\mathbf{w}^{t+\frac{1}{2}} = \mathbf{w}^t + \eta_t \mathbf{A}^{\top} \left( \frac{\exp(-z_i)}{1 + \exp(-z_i)} \right)_{i=1}^m$$

- 3 Then compute the proximity operator

$$\begin{aligned} \mathbf{w}^{t+1} &= \text{prox}_{\lambda\eta_t}(\mathbf{w}^{t+\frac{1}{2}}) \\ &= \begin{cases} \mathbf{w}_j^{t+\frac{1}{2}} + \lambda\eta_t & (\mathbf{w}_j^{t+\frac{1}{2}} < -\lambda\eta_t), \\ 0 & (-\lambda\eta_t \leq \mathbf{w}_j^{t+\frac{1}{2}} \leq \lambda\eta_t), \\ \mathbf{w}_j^{t+\frac{1}{2}} - \lambda\eta_t & (\mathbf{w}_j^{t+\frac{1}{2}} > \lambda\eta_t). \end{cases} \end{aligned}$$

# Matrix completion via IST (Mazumder et al. 10)

Loss function:

$$f(\mathbf{X}) = \frac{1}{2} \|\Omega(\mathbf{X} - \mathbf{Y})\|^2.$$

Regularization:

$$g(\mathbf{X}) = \lambda \sum_{j=1}^r \sigma_j(\mathbf{X}) \quad (\text{S}_1\text{-norm}).$$

# Matrix completion via IST (Mazumder et al. 10)

Loss function:

$$f(\mathbf{X}) = \frac{1}{2} \|\Omega(\mathbf{X} - \mathbf{Y})\|^2.$$

gradient:

$$\nabla f(\mathbf{X}) = \Omega^\top(\Omega(\mathbf{X} - \mathbf{Y}))$$

Regularization:

$$g(\mathbf{X}) = \lambda \sum_{j=1}^r \sigma_j(\mathbf{X}) \quad (\text{S}_1\text{-norm}).$$

# Matrix completion via IST (Mazumder et al. 10)

Loss function:

$$f(\mathbf{X}) = \frac{1}{2} \|\Omega(\mathbf{X} - \mathbf{Y})\|^2.$$

gradient:

$$\nabla f(\mathbf{X}) = \Omega^\top(\Omega(\mathbf{X} - \mathbf{Y}))$$

Regularization:

$$g(\mathbf{X}) = \lambda \sum_{j=1}^r \sigma_j(\mathbf{X}) \quad (\mathcal{S}_1\text{-norm}).$$

Prox operator (Singular Value Thresholding):

$$\text{prox}_\lambda(\mathbf{Z}) = \mathbf{U} \max(\mathbf{S} - \lambda \mathbf{I}, 0) \mathbf{V}^\top.$$

# Matrix completion via IST (Mazumder et al. 10)

Loss function:

$$f(\mathbf{X}) = \frac{1}{2} \|\Omega(\mathbf{X} - \mathbf{Y})\|^2.$$

gradient:

$$\nabla f(\mathbf{X}) = \Omega^\top(\Omega(\mathbf{X} - \mathbf{Y}))$$

Regularization:

$$g(\mathbf{X}) = \lambda \sum_{j=1}^r \sigma_j(\mathbf{X}) \quad (\text{S}_1\text{-norm}).$$

Prox operator (Singular Value Thresholding):

$$\text{prox}_\lambda(\mathbf{Z}) = \mathbf{U} \max(\mathbf{S} - \lambda \mathbf{I}, 0) \mathbf{V}^\top.$$

Iteration:

$$\mathbf{X}^{t+1} = \text{prox}_{\lambda\eta_t} \left( \underbrace{(\mathbf{I} - \eta_t \Omega^\top \Omega)(\mathbf{X}^t)}_{\text{fill in missing}} + \underbrace{\eta_t \Omega^\top \Omega(\mathbf{Y}^t)}_{\text{observed}} \right)$$

- When  $\eta_t = 1$ , fill missings with predicted values  $\mathbf{X}^t$ , overwrite the observed with observed values, then soft-threshold.

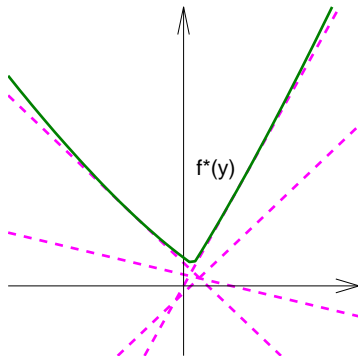
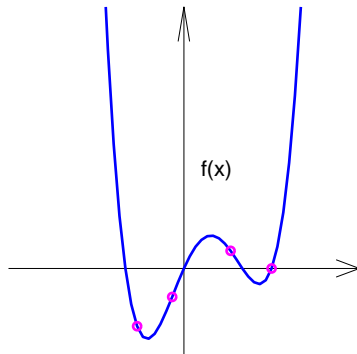
# Conjugate duality and dual ascent

- Convex conjugate function
- Lagrangian relaxation and dual problem

# Conjugate duality

The convex conjugate  $f^*$  of a function  $f$ :

$$f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \mathbb{R}^n} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}))$$



Since the maximum over linear functions is always convex,  $f$  need not be convex.

# Demo

## Try

- `demo_conjugate (@ (x) x.^2/2, -5:0.1:5);`
- `demo_conjugate (@ (x) abs (x), -5:0.1:5);`
- `demo_conjugate (@ (x) x.*log (x) + (1-x) .*log (1-x), ...  
0.001:0.001:0.999);`



# Conjugate duality (dual view)

## Convex conjugate function

Every pair  $(\mathbf{y}, f^*(\mathbf{y}))$  corresponds to a tangent line  $\langle \mathbf{x}, \mathbf{y} \rangle - f^*(\mathbf{y})$  of the original function  $f(\mathbf{x})$ .

Because

$f^*(\mathbf{y}) = \sup_{\mathbf{x}} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}))$   
implies

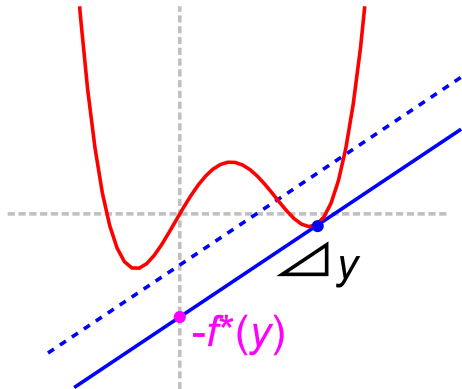
- If  $t < f^*(\mathbf{y})$ , there is a  $\mathbf{x}$  s.t.

$$f(\mathbf{x}) < \langle \mathbf{x}, \mathbf{y} \rangle - t.$$

- If  $t \geq f^*(\mathbf{y})$ ,

$$f(\mathbf{x}) \geq \langle \mathbf{x}, \mathbf{y} \rangle - t$$

for every  $\mathbf{x}$ .



# Conjugate duality (dual view)

## Convex conjugate function

Every pair  $(\mathbf{y}, f^*(\mathbf{y}))$  corresponds to a tangent line  $\langle \mathbf{x}, \mathbf{y} \rangle - f^*(\mathbf{y})$  of the original function  $f(\mathbf{x})$ .

Because

$f^*(\mathbf{y}) = \sup_{\mathbf{x}} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}))$   
implies

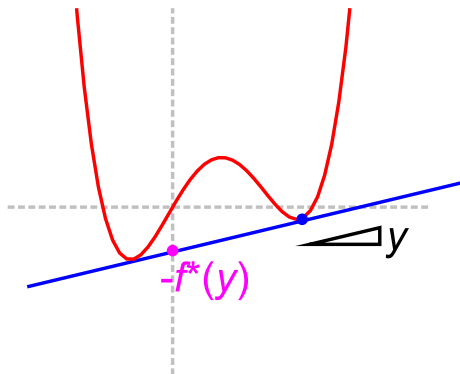
- If  $t < f^*(\mathbf{y})$ , there is a  $\mathbf{x}$  s.t.

$$f(\mathbf{x}) < \langle \mathbf{x}, \mathbf{y} \rangle - t.$$

- If  $t \geq f^*(\mathbf{y})$ ,

$$f(\mathbf{x}) \geq \langle \mathbf{x}, \mathbf{y} \rangle - t$$

for every  $\mathbf{x}$ .



# Conjugate duality (dual view)

## Convex conjugate function

Every pair  $(\mathbf{y}, f^*(\mathbf{y}))$  corresponds to a tangent line  $\langle \mathbf{x}, \mathbf{y} \rangle - f^*(\mathbf{y})$  of the original function  $f(\mathbf{x})$ .

Because

$f^*(\mathbf{y}) = \sup_{\mathbf{x}} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}))$   
implies

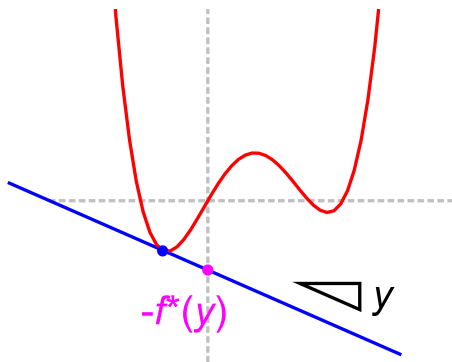
- If  $t < f^*(\mathbf{y})$ , there is a  $\mathbf{x}$  s.t.

$$f(\mathbf{x}) < \langle \mathbf{x}, \mathbf{y} \rangle - t.$$

- If  $t \geq f^*(\mathbf{y})$ ,

$$f(\mathbf{x}) \geq \langle \mathbf{x}, \mathbf{y} \rangle - t$$

for every  $\mathbf{x}$ .



# Conjugate duality (dual view)

## Convex conjugate function

Every pair  $(\mathbf{y}, f^*(\mathbf{y}))$  corresponds to a tangent line  $\langle \mathbf{x}, \mathbf{y} \rangle - f^*(\mathbf{y})$  of the original function  $f(\mathbf{x})$ .

Because

$f^*(\mathbf{y}) = \sup_{\mathbf{x}} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}))$   
implies

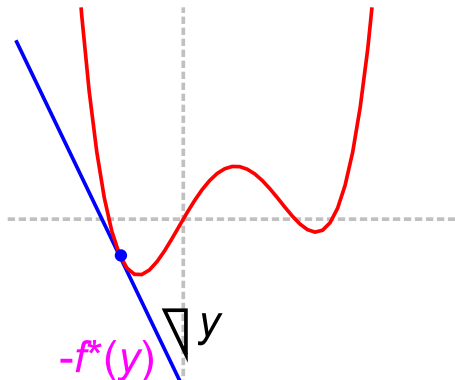
- If  $t < f^*(\mathbf{y})$ , there is a  $\mathbf{x}$  s.t.

$$f(\mathbf{x}) < \langle \mathbf{x}, \mathbf{y} \rangle - t.$$

- If  $t \geq f^*(\mathbf{y})$ ,

$$f(\mathbf{x}) \geq \langle \mathbf{x}, \mathbf{y} \rangle - t$$

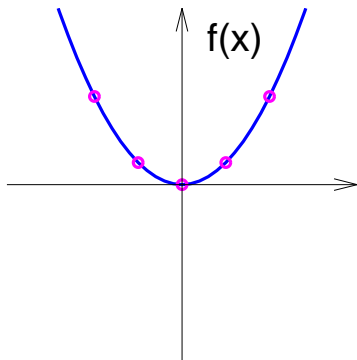
for every  $\mathbf{x}$ .



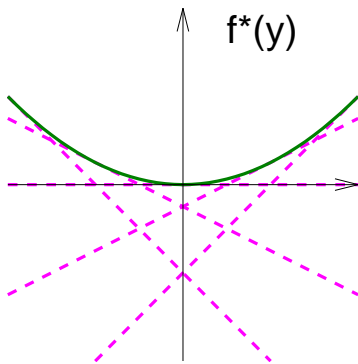
# Example of conjugate duality $f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \mathbb{R}^n} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}))$

- Quadratic function

$$f(x) = \frac{x^2}{2\sigma^2}$$



$$f^*(y) = \frac{\sigma^2 y^2}{2}$$



# Example of conjugate duality $f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \mathbb{R}^n} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}))$

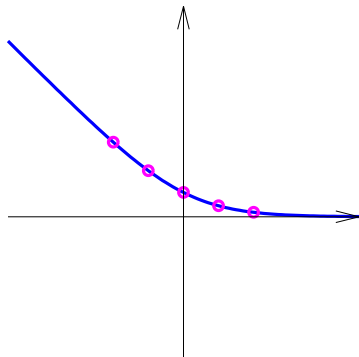
- Logistic loss function

$$f(x) = \log(1 + \exp(-x))$$

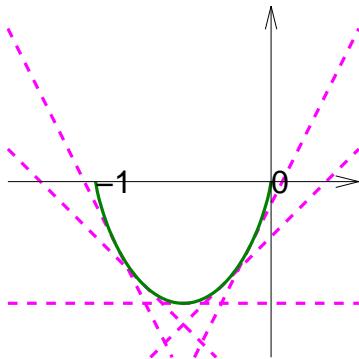
# Example of conjugate duality $f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \mathbb{R}^n} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}))$

- Logistic loss function

$$f(x) = \log(1 + \exp(-x))$$



$$f^*(-y) = y \log(y) + (1 - y) \log(1 - y)$$



# Example of conjugate duality $f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \mathbb{R}^n} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}))$

- L1 regularizer

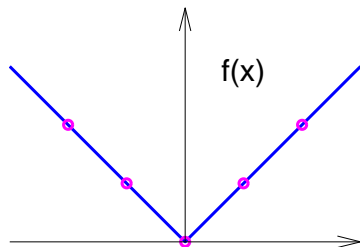
$$f(x) = |x|$$



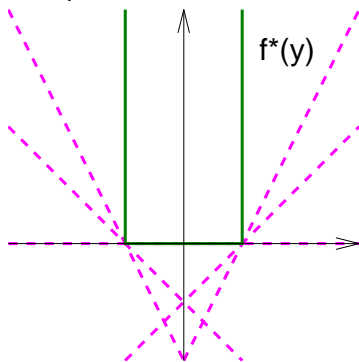
# Example of conjugate duality $f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \mathbb{R}^n} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}))$

- L1 regularizer

$$f(x) = |x|$$

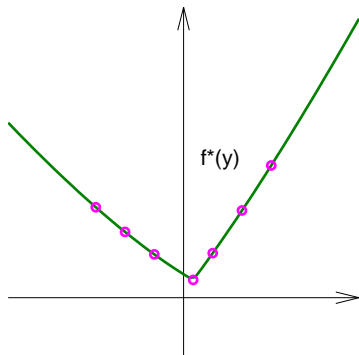
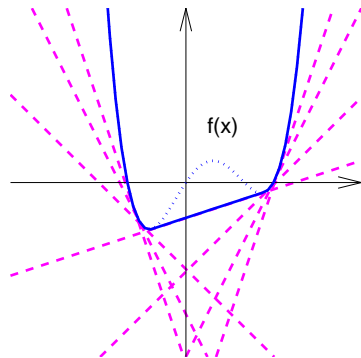


$$f^*(y) = \begin{cases} 0 & (-1 \leq y \leq 1) \\ +\infty & (\text{otherwise}) \end{cases}$$



# Bi-conjugate $f^{**}$ may be different from $f$

For nonconvex  $f$ ,



# Lagrangian relaxation

Our optimization problem:

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{A}\mathbf{w}) + g(\mathbf{w})$$

$$\left( \begin{array}{l} \text{For example} \\ f(\mathbf{z}) = \frac{1}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 \\ \text{(squared loss)} \end{array} \right)$$

# Lagrangian relaxation

Our optimization problem:

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{A}\mathbf{w}) + g(\mathbf{w})$$

$$\left( \begin{array}{l} \text{For example} \\ f(\mathbf{z}) = \frac{1}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 \\ \text{(squared loss)} \end{array} \right)$$

Equivalently written as

$$\underset{\mathbf{z} \in \mathbb{R}^m, \mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{z}) + g(\mathbf{w}),$$

$$\text{s.t.} \quad \mathbf{z} = \mathbf{A}\mathbf{w} \quad (\text{equality constraint})$$

# Lagrangian relaxation

Our optimization problem:

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{A}\mathbf{w}) + g(\mathbf{w})$$

$$\left( \begin{array}{l} \text{For example} \\ f(\mathbf{z}) = \frac{1}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 \\ \text{(squared loss)} \end{array} \right)$$

Equivalently written as

$$\underset{\mathbf{z} \in \mathbb{R}^m, \mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{z}) + g(\mathbf{w}),$$

$$\text{s.t.} \quad \mathbf{z} = \mathbf{A}\mathbf{w} \quad (\text{equality constraint})$$

Lagrangian relaxation

$$\underset{\mathbf{z}, \mathbf{w}}{\text{minimize}} \quad \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha) = f(\mathbf{z}) + g(\mathbf{w}) + \alpha^\top (\mathbf{z} - \mathbf{A}\mathbf{w})$$

# Lagrangian relaxation

Our optimization problem:

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{A}\mathbf{w}) + g(\mathbf{w})$$

$$\left( \begin{array}{l} \text{For example} \\ f(\mathbf{z}) = \frac{1}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 \\ \text{(squared loss)} \end{array} \right)$$

Equivalently written as

$$\underset{\mathbf{z} \in \mathbb{R}^m, \mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{z}) + g(\mathbf{w}),$$

$$\text{s.t.} \quad \mathbf{z} = \mathbf{A}\mathbf{w} \quad (\text{equality constraint})$$

Lagrangian relaxation

$$\underset{\mathbf{z}, \mathbf{w}}{\text{minimize}} \quad \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha) = f(\mathbf{z}) + g(\mathbf{w}) + \alpha^\top (\mathbf{z} - \mathbf{A}\mathbf{w})$$

- As long as  $\mathbf{z} = \mathbf{A}\mathbf{w}$ , the relaxation is exact.
- $\sup_{\alpha} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha)$  recovers the original problem.
- Minimum of  $\mathcal{L}$  is no greater than the minimum of the original.

# Weak duality

$$\inf_{\mathbf{z}, \mathbf{w}} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha) \leq \inf_{\mathbf{w}} (f(\mathbf{A}\mathbf{w}) + g(\mathbf{w})) =: p^*$$

proof

$$\begin{aligned} \inf_{\mathbf{z}, \mathbf{w}} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha) &= \min \left( \inf_{\mathbf{z}=\mathbf{A}\mathbf{w}} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha), \inf_{\mathbf{z} \neq \mathbf{A}\mathbf{w}} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha) \right) \\ &= \min \left( p^*, \inf_{\mathbf{z} \neq \mathbf{A}\mathbf{w}} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha) \right) \\ &\leq p^* \end{aligned}$$

# Dual problem

From the above argument

$$d(\alpha) := \inf_{\mathbf{z}, \mathbf{w}} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha)$$

is a lower bound for  $p^*$  for any  $\alpha$ . Why don't we maximize over  $\alpha$ ?



# Dual problem

From the above argument

$$d(\alpha) := \inf_{\mathbf{z}, \mathbf{w}} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha)$$

is a lower bound for  $p^*$  for any  $\alpha$ . Why don't we maximize over  $\alpha$ ?

## Dual problem

$$\underset{\alpha \in \mathbb{R}^m}{\text{maximize}} \quad d(\alpha)$$

Note

$$\sup_{\alpha} \inf_{\mathbf{z}, \mathbf{w}} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha) = d^* \leq p^* = \inf_{\mathbf{z}, \mathbf{w}} \sup_{\alpha} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha)$$

If  $d^* = p^*$ , **strong duality** holds. This is the case if  $f$  and  $g$  both closed and convex.

# Fenchel's duality

For convex<sup>1</sup> functions  $f$  and  $g$ , and a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$

$$\sup_{\alpha \in \mathbb{R}^m} \left( -f^*(-\alpha) - g^*(\mathbf{A}^\top \alpha) \right) = \inf_{\mathbf{w} \in \mathbb{R}^n} \left( f(\mathbf{A}\mathbf{w}) + g(\mathbf{w}) \right)$$



Werner Fenchel

1905 – 1988

- Only need **conjugate functions  $f^*$  and  $g^*$**  to compute the dual.
- We can make a list of them (like Laplace transform)

MATLAB Exercise 1.5:

- Compute the Fenchel dual of L1-logistic regression problem in Ex.1 and implement the stopping criterion: stop optimization if

$$(\text{obj}_{\text{prim}} - \text{obj}_{\text{dual}}) / \text{obj}_{\text{prim}} < \epsilon \quad (\text{relative duality gap}).$$

---

<sup>1</sup>More precisely, proper, closed, and convex.

# Derivation of Fenchel's duality theorem

$$\begin{aligned}d(\alpha) &= \inf_{\mathbf{z}, \mathbf{w}} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha) \\&= \inf_{\mathbf{z}, \mathbf{w}} \left( f(\mathbf{z}) + g(\mathbf{w}) + \alpha^\top (\mathbf{z} - \mathbf{A}\mathbf{w}) \right) \\&= \inf_{\mathbf{z}} (f(\mathbf{z}) + \langle \alpha, \mathbf{z} \rangle) + \inf_{\mathbf{w}} \left( g(\mathbf{w}) - \langle \mathbf{A}^\top \alpha, \mathbf{w} \rangle \right) \\&= -\sup_{\mathbf{z}} (\langle -\alpha, \mathbf{z} \rangle - f(\mathbf{z})) - \sup_{\mathbf{w}} \left( \langle \mathbf{A}^\top \alpha, \mathbf{w} \rangle - g(\mathbf{w}) \right) \\&= -f^*(-\alpha) - g^*(\mathbf{A}^\top \alpha)\end{aligned}$$

# Augmented Lagrangian and ADMM

## Learning objectives

- Structured sparse estimation
- Augmented Lagrangian
- Alternating direction method of multipliers (ADMM)

# Recap: Simple vs. structured sparse estimation problems

- Simple sparse estimation problem

$$\underset{\mathbf{w}}{\text{minimize}} \quad f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

- ▶ SNP analysis
- ▶ Compressed sensing with  $\Phi^{-1}$  (e.g., wavelet)
- ▶ Collaborative filtering (matrix completion)

- Structured sparse estimation problem

$$\underset{\mathbf{w}}{\text{minimize}} \quad f(\mathbf{w}) + \lambda \|\Phi \mathbf{w}\|_1$$

- ▶ Compressed sensing without  $\Phi^{-1}$  (e.g., total variation)
- ▶ Low-rank tensor completion

# Total Variation based image denoising [Rudin, Osher, Fatemi 92]

$$\underset{W}{\text{minimize}} \quad \frac{1}{2} \|W - M\|_F^2 + \lambda \sum_{i,j} \left\| \begin{pmatrix} \partial_x W_{ij} \\ \partial_y W_{ij} \end{pmatrix} \right\|_2$$

Original  $W_0$



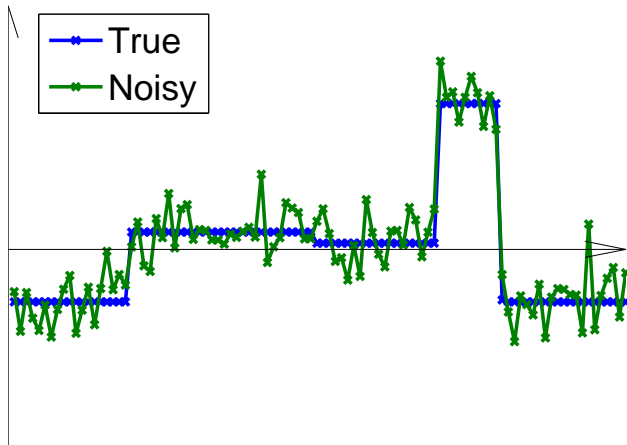
Observed  $M$



# In one dimension

- Fused lasso [Tibshirani et al. 05]

$$\underset{\mathbf{w}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \sum_{j=1}^{n-1} |w_{j+1} - w_j|$$



# Structured sparse estimation

- TV denoising

$$\underset{W}{\text{minimize}} \quad \frac{1}{2} \|W - M\|_F^2 + \lambda \sum_{i,j} \left\| \begin{pmatrix} \partial_x W_{ij} \\ \partial_y W_{ij} \end{pmatrix} \right\|_2$$

- Fused lasso

$$\underset{\mathbf{w}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \sum_{j=1}^{n-1} |w_{j+1} - w_j|$$



# Structured sparse estimation

- TV denoising

$$\underset{W}{\text{minimize}} \quad \frac{1}{2} \|W - M\|_F^2 + \lambda \sum_{i,j} \left\| \begin{pmatrix} \partial_x W_{ij} \\ \partial_y W_{ij} \end{pmatrix} \right\|_2$$

- Fused lasso

$$\underset{\mathbf{w}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \sum_{j=1}^{n-1} |w_{j+1} - w_j|$$

## Structured sparse estimation problem

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \underbrace{f(\mathbf{w})}_{\text{data-fit}} + \underbrace{\lambda \|\mathbf{A}\mathbf{w}\|_1}_{\text{regularization}}$$

# Structured sparse estimation problem

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \underbrace{f(\mathbf{w})}_{\text{data-fit}} + \underbrace{\lambda \|\mathbf{A}\mathbf{w}\|_1}_{\text{regularization}}$$

- Not easy to compute prox operator (because it is **non-separable**)  
⇒ difficult to apply **IST-type methods**.

# Structured sparse estimation problem

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \underbrace{f(\mathbf{w})}_{\text{data-fit}} + \underbrace{\lambda \|\mathbf{A}\mathbf{w}\|_1}_{\text{regularization}}$$

- Not easy to compute prox operator (because it is **non-separable**)  
⇒ difficult to apply **IST-type methods**.

Can we use the Lagrangian relaxation trick?

# Forming the Lagrangian

## Structured sparsity problem

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \underbrace{f(\mathbf{w})}_{\text{data-fit}} + \underbrace{\lambda \|\mathbf{A}\mathbf{w}\|_1}_{\text{regularization}}$$

Equivalently written as

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{w}) + \underbrace{\lambda \|\mathbf{z}\|_1}_{\text{separable!}},$$

s.t.  $\mathbf{z} = \mathbf{A}\mathbf{w}$  (equality constraint)

# Forming the Lagrangian

## Structured sparsity problem

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \underbrace{f(\mathbf{w})}_{\text{data-fit}} + \underbrace{\lambda \|\mathbf{A}\mathbf{w}\|_1}_{\text{regularization}}$$

Equivalently written as

$$\begin{aligned} \underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad & f(\mathbf{w}) + \underbrace{\lambda \|\mathbf{z}\|_1}_{\text{separable!}}, \\ \text{s.t.} \quad & \mathbf{z} = \mathbf{A}\mathbf{w} \quad (\text{equality constraint}) \end{aligned}$$

## Lagrangian function

$$\mathcal{L}(\mathbf{w}, \mathbf{z}, \boldsymbol{\alpha}) = f(\mathbf{w}) + \lambda \|\mathbf{z}\|_1 + \boldsymbol{\alpha}^\top (\mathbf{z} - \mathbf{A}\mathbf{w}).$$

$\boldsymbol{\alpha}$ : Lagrangian multiplier vector.

# Dual ascent

## Dual problem

$$\max_{\alpha} \inf_{\mathbf{z}, \mathbf{w}} \left( f(\mathbf{w}) + \lambda \|\mathbf{z}\|_1 + \alpha^\top (\mathbf{z} - \mathbf{A}\mathbf{w}) \right)$$

We can compute the dual objective  $d(\alpha)$  by separately minimizing

$$(1) \quad \min_{\mathbf{w}} \left( f(\mathbf{w}) - \alpha^\top \mathbf{A}\mathbf{w} \right)$$

$$(2) \quad \min_{\mathbf{z}} \left( \lambda \|\mathbf{z}\|_1 + \alpha^\top \mathbf{z} \right)$$

# Dual ascent

## Dual problem

$$\max_{\alpha} \inf_{\mathbf{z}, \mathbf{w}} \left( f(\mathbf{w}) + \lambda \|\mathbf{z}\|_1 + \alpha^\top (\mathbf{z} - \mathbf{A}\mathbf{w}) \right)$$

We can compute the dual objective  $d(\alpha)$  by separately minimizing

$$(1) \quad \min_{\mathbf{w}} \left( f(\mathbf{w}) - \alpha^\top \mathbf{A}\mathbf{w} \right) = -f^*(\mathbf{A}^\top \alpha),$$

$$(2) \quad \min_{\mathbf{z}} \left( \lambda \|\mathbf{z}\|_1 + \alpha^\top \mathbf{z} \right) = -(\lambda \|\cdot\|_1)^*(-\alpha).$$

# Dual ascent

## Dual problem

$$\max_{\alpha} \inf_{\mathbf{z}, \mathbf{w}} \left( f(\mathbf{w}) + \lambda \|\mathbf{z}\|_1 + \alpha^\top (\mathbf{z} - \mathbf{A}\mathbf{w}) \right)$$

We can compute the dual objective  $d(\alpha)$  by separately minimizing

$$(1) \quad \min_{\mathbf{w}} \left( f(\mathbf{w}) - \alpha^\top \mathbf{A}\mathbf{w} \right) = -f^*(\mathbf{A}^\top \alpha),$$

$$(2) \quad \min_{\mathbf{z}} \left( \lambda \|\mathbf{z}\|_1 + \alpha^\top \mathbf{z} \right) = -(\lambda \|\cdot\|_1)^*(-\alpha).$$

But also we get the gradient of  $d(\alpha)$  (for free) as follows:

$$\nabla_{\alpha} d(\alpha) = \mathbf{z}^* - \mathbf{A}\mathbf{w}^*,$$

where  $\mathbf{w}^*$ : argmin of (1),  $\mathbf{z}^*$ : argmin of (2). See Chapter 6, Bertsekas 1999.

Gradient ascent (in the dual)!



# Dual ascent (Arrow, Hurwicz, & Uzawa 1958)

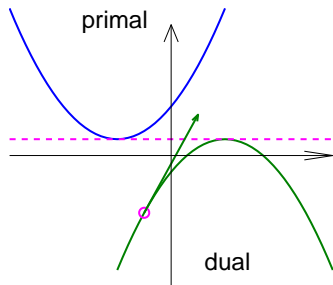
$$\left\{ \begin{array}{l} \text{Minimize the Lagrangian wrt } \mathbf{x} \text{ and } \mathbf{z}: \\ \mathbf{w}^{t+1} = \operatorname{argmin}_{\mathbf{w}} (f(\mathbf{w}) - \alpha^\top \mathbf{A}\mathbf{w}). \\ \\ \mathbf{z}^{t+1} = \operatorname{argmin}_{\mathbf{z}} (\lambda \|\mathbf{z}\|_1 + \alpha^\top \mathbf{z}), \\ \text{Update the Lagrangian multiplier } \alpha^t: \\ \alpha^{t+1} = \alpha^t + \eta_t (\mathbf{z}^{t+1} - \mathbf{A}\mathbf{w}^{t+1}). \end{array} \right.$$

- **Pro:** Very simple.
- **Con:** When  $f^*$  or  $g^*$  is non-differentiable, it is a dual subgradient method (convergence more tricky)

NB:  $f^*$  is differentiable  $\Leftrightarrow f$  is strictly convex.



H. Uzawa



# Forming the *augmented* Lagrangian

Structured sparsity problem

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \underbrace{f(\mathbf{w})}_{\text{data-fit}} + \underbrace{\lambda \|\mathbf{A}\mathbf{w}\|_1}_{\text{regularization}}$$

Equivalently written as (for any  $\eta > 0$ )

$$\begin{aligned} \underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad & f(\mathbf{w}) + \underbrace{\lambda \|\mathbf{z}\|_1}_{\text{separable!}} + \underbrace{\frac{\eta}{2} \|\mathbf{z} - \mathbf{A}\mathbf{w}\|_2^2}_{\text{penalty term}}, \\ \text{s.t.} \quad & \mathbf{z} = \mathbf{A}\mathbf{w} \quad (\text{equality constraint}) \end{aligned}$$

# Forming the *augmented* Lagrangian

Structured sparsity problem

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \underbrace{f(\mathbf{w})}_{\text{data-fit}} + \underbrace{\lambda \|\mathbf{A}\mathbf{w}\|_1}_{\text{regularization}}$$

Equivalently written as (for any  $\eta > 0$ )

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{w}) + \underbrace{\lambda \|\mathbf{z}\|_1}_{\text{separable!}} + \underbrace{\frac{\eta}{2} \|\mathbf{z} - \mathbf{A}\mathbf{w}\|_2^2}_{\text{penalty term}},$$

s.t.  $\mathbf{z} = \mathbf{A}\mathbf{w}$  (equality constraint)

## Augmented Lagrangian function

$$\mathcal{L}_\eta(\mathbf{w}, \mathbf{z}, \boldsymbol{\alpha}) = f(\mathbf{w}) + \lambda \|\mathbf{z}\|_1 + \boldsymbol{\alpha}^\top (\mathbf{z} - \mathbf{A}\mathbf{w}) + \frac{\eta}{2} \|\mathbf{z} - \mathbf{A}\mathbf{w}\|_2^2$$

$\boldsymbol{\alpha}$ : Lagrangian multiplier,  $\eta$ : penalty parameter

# Augmented Lagrangian Method

## Augmented Lagrangian function

$$\mathcal{L}_\eta(\mathbf{w}, \mathbf{z}, \alpha) = f(\mathbf{w}) + \lambda \|\mathbf{z}\|_1 + \alpha^\top (\mathbf{z} - \mathbf{A}\mathbf{w}) + \frac{\eta}{2} \|\mathbf{z} - \mathbf{A}\mathbf{w}\|^2.$$

## Augmented Lagrangian method (Hestenes 69, Powell 69)

$$\left\{ \begin{array}{l} \text{Minimize the AL function wrt } \mathbf{w} \text{ and } \mathbf{z}: \\ (\mathbf{w}^{t+1}, \mathbf{z}^{t+1}) = \underset{\mathbf{w} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{R}^m}{\operatorname{argmin}} \mathcal{L}_\eta(\mathbf{w}, \mathbf{z}, \alpha^t). \\ \\ \text{Update the Lagrangian multiplier:} \\ \alpha^{t+1} = \alpha^t + \eta(\mathbf{z}^{t+1} - \mathbf{A}\mathbf{w}^{t+1}). \end{array} \right.$$

- **Pro**: The dual is **always** differentiable due to the penalty term.
- **Con**: Cannot minimize over  $\mathbf{w}$  and  $\mathbf{z}$  independently

# Alternating Direction Method of Multipliers (ADMM; Gabay & Mercier 76)

$$\left\{ \begin{array}{l} \text{Minimize the AL function } \mathcal{L}_\eta(\mathbf{w}, \mathbf{z}^t, \alpha^t) \text{ wrt } \mathbf{w}: \\ \mathbf{w}^{t+1} = \underset{\mathbf{w} \in \mathbb{R}^n}{\operatorname{argmin}} \left( f(\mathbf{w}) + \frac{\eta}{2} \|\mathbf{z}^t - \mathbf{A}\mathbf{w} + \alpha^t/\eta\|_2^2 \right). \\ \\ \text{Minimize the AL function } \mathcal{L}_\eta(\mathbf{w}^{t+1}, \mathbf{z}, \alpha^t) \text{ wrt } \mathbf{z}: \\ \mathbf{z}^{t+1} = \underset{\mathbf{z} \in \mathbb{R}^m}{\operatorname{argmin}} \left( \lambda \|\mathbf{z}\|_1 + \frac{\eta}{2} \|\mathbf{z} - \mathbf{A}\mathbf{w}^{t+1} + \alpha^t/\eta\|_2^2 \right). \\ \\ \text{Update the Lagrangian multiplier:} \\ \alpha^{t+1} = \alpha^t + \eta(\mathbf{z}^{t+1} - \mathbf{A}\mathbf{w}^{t+1}). \end{array} \right.$$

- Looks ad-hoc but convergence can be shown rigorously.
- Stability does not rely on the choice of step-size  $\eta$ .
- The newly updated  $\mathbf{w}^{t+1}$  enters the computation of  $\mathbf{z}^{t+1}$ .

# MATLAB Exercise 2: implement an ADMM for fused lasso

## Fused lasso

$$\underset{\mathbf{w}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \sum_{j=1}^{n-1} |w_{j+1} - w_j|$$

- What is the loss function  $f$ ?
- What is the matrix  $\mathbf{A}$  for fused lasso?
- How does the  $\mathbf{w}$ -update step look?
- How does the  $\mathbf{z}$ -update step look?

# Conclusion

- Three approaches for various sparse estimation problems
  - ▶ Iterative shrinkage/thresholding – [proximity operator](#)
  - ▶ Uzawa's method – [convex conjugate function](#)
  - ▶ ADMM – combination of the above two
- Above methods go beyond black-box models (e.g., gradient descent or Newton's method) – takes better care of the problem structures.
- These methods are simple enough to be implemented rapidly, but should not be considered as a *silver bullet*.  
⇒ [Trade-off between:](#)
  - ▶ Quick implementation – test new ideas rapidly
  - ▶ Efficient optimization – more inspection/try-and-error/cross validation

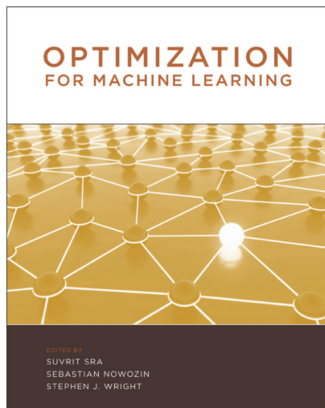
# Topics we did not cover

- Beyond polynomial convergence  $O(1/k^2)$ 
  - ▶ Dual Augmented Lagrangian (DAL) converges super-linearly  $o(\exp(-k))$ . Software  
<http://mloss.org/software/view/183/>  
(This is limited to non-structured sparse estimation.)
- Beyond convexity
  - ▶ Generalized eigenvalue problems.
  - ▶ Difference of convex (DC) programming.
  - ▶ **Dual ascent** (or dual decomposition) for sequence labeling in natural language processing; see [Wainwright, Jaakkola, Willsky 05; Koo et al. 10]
- Stochastic optimization
  - ▶ Good tutorial by Nathan Srebro (ICML2010)



# Optimization for Machine Learning

A new book “Optimization for Machine Learning” (2011)



# Possible projects

- 1 Compare the three approaches, namely IST, dual ascent, and ADMM, and discuss empirically (and theoretically) their pros and cons.
- 2 Apply one of the methods discussed in the lecture to model some real problem with (structured) sparsity or low-rank matrix.

# References

## Recent surveys

- Tomioka, Suzuki, & Sugiyama (2011) Augmented Lagrangian Methods for Learning, Selecting, and Combining Features. In Sra, Nowozin, Wright., editors, *Optimization for Machine Learning*, MIT Press.
- Combettes & Pesquet (2010) Proximal splitting methods in signal processing. In *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*. Springer-Verlag.
- Boyd, Parikh, Peleato, & Eckstein (2010) Distributed optimization and statistical learning via the alternating direction method of multipliers.

## Textbooks

- Rockafellar (1970) *Convex Analysis*. Princeton University Press.
- Bertsekas (1999) *Nonlinear Programming*. Athena Scientific.
- Nesterov (2003) *Introductory Lectures on Convex Optimization: A Basic Course*. Springer.
- Boyd & Vandenberghe. (2004) *Convex optimization*, Cambridge University Press.

# References

## IST/FISTA

- Moreau (1965) Proximité et dualité dans un espace Hilbertien. Bul letin de la S. M. F.
- Nesterov (2007) Gradient Methods for Minimizing Composite Objective Function.
- Beck & Teboulle (2009) A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems. SIAM J Imag Sci 2, 183–202.

## Dual ascent

- Arrow, Hurwicz, & Uzawa (1958) Studies in Linear and Non-Linear Programming. Stanford University Press.
- Chapter 6 in Bertsekas (1999).
- Wainwright, Jaakkola, & Willsky (2005) Map estimation via agreement on trees: message-passing and linear programming. IEEE Trans IT, 51(11).

## Augmented Lagrangian

- Rockafellar (1976) Augmented Lagrangians and applications of the proximal point algorithm in convex programming. Math. of Oper. Res. 1.
- Bertsekas (1982) Constrained Optimization and Lagrange Multiplier Methods. Academic Press.
- Tomioka, Suzuki, & Sugiyama (2011) Super-Linear Convergence of Dual Augmented Lagrangian Algorithm for Sparse Learning. JMLR 12.

# References

## ADMM

- Gabay & Mercier (1976) A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Comput Math Appl* 2, 17–40.
- Lions & Mercier (1979) Splitting Algorithms for the Sum of Two Nonlinear Operators. *SIAM J Numer Anal* 16, 964–979.
- Eckstein & Bertsekas (1992) On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators.

## Matrices

- Srebro, Rennie, & Jaakkola (2005) Maximum-Margin Matrix Factorization. *Advances in NIPS* 17, 1329–1336.
- Cai, Candès, & Shen (2008) A singular value thresholding algorithm for matrix completion.
- Tomioka, Suzuki, Sugiyama, & Kashima (2010) A Fast Augmented Lagrangian Algorithm for Learning Low-Rank Matrices. In *ICML 2010*.
- Mazumder, Hastie, & Tibshirani (2010) Spectral Regularization Algorithms for Learning Large Incomplete Matrices. *JMLR* 11, 2287–2322.

# References

## Multi-task/Multiple kernel learning

- Evgeniou, Micchelli, & Pontil (2005) Learning Multiple Tasks with Kernel Methods. JMLR 6, 615–637.
- Lanckriet, Christiani, Bartlett, Ghaoui, & Jordan (2004) Learning the Kernel Matrix with Semidefinite Programming.
- Bach, Thibaux, & Jordan (2005) Computing regularization paths for learning multiple kernels. Advances in NIPS, 73–80.

## Structured sparsity

- Tibshirani, Saunders, Rosset, Zhu and Knight. (2005) Sparsity and smoothness via the fused lasso. J. Roy. Stat. Soc. B, 67.
- Rudin, Osher, Fetemi. (1992) Nonlinear total variation based noise removal algorithms. Physica D: Nonlinear Phenomena, 60.
- Goldstein & Osher (2009) Split Bregman method for L1 regularization problems. SIAM J. Imag. Sci. 2.
- Mairal, Jenatton, Obozinski, & Bach. (2011) Convex and network flow optimization for structured sparsity.

## Bayes & Probabilistic Inference

- Wainwright & Jordan (2008) Graphical Models, Exponential Families, and Variational Inference.