

Convex Optimization: Old Tricks for New Problems

Ryota Tomioka¹

¹The University of Tokyo

2011-08-26 @ DTU PhD Summer Course

Why care about convex optimization (and sparsity)?

A typical machine learning problem (1/2)

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \underbrace{\frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2}_{\text{data-fit}} + \underbrace{\phi_\lambda(\mathbf{w})}_{\text{Regularization}}$$

Design
 $m \times n$

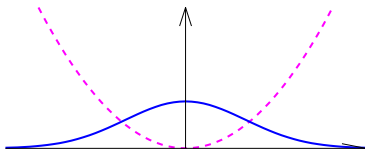
Variables
 $n \times 1$

Targets
 $m \times 1$



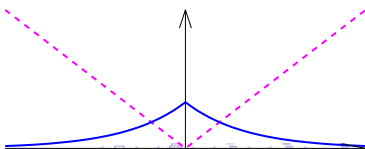
Ridge penalty

$$\phi_\lambda = \frac{\lambda}{2} \sum_{j=1}^n w_j^2.$$



L1 penalty

$$\phi_\lambda = \lambda \sum_{j=1}^n |w_j|.$$



A typical machine learning problem (2/2)

Logistic regression for binary ($y_i \in \{-1, +1\}$) classification:

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

The logistic loss function

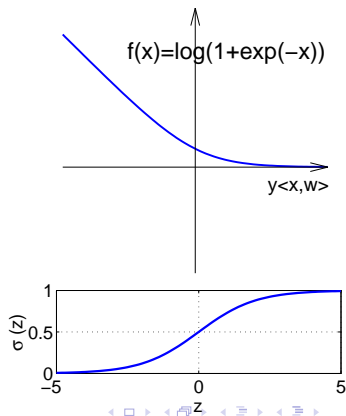
$$\log(1 + e^{-yz}) = -\log P(Y = y|z)$$

negative log-likelihood

where

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

logistic function



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

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

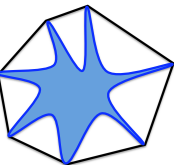
Bayesian inference as a convex optimization

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

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

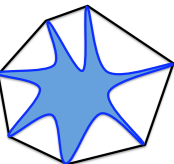
Bayesian inference as a convex optimization

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

where

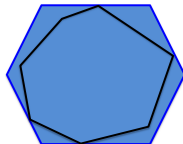
$$\begin{aligned} f(w) &= \underbrace{-\log P(D|w)}_{\text{neg. log likelihood}} - \underbrace{\log P(w)}_{\text{neg. log prior}} \\ \Rightarrow \quad q(w) &= \frac{1}{Z} e^{-f(w)} \quad (\text{Bayesian posterior}) \end{aligned}$$

Inner approximations



- Variational Bayes
- Empirical Bayes

Outer approximations



- Belief propagation

See Wainwright & Jordan 08.

Convex optimization = standard forms (boring?)

Example: Linear Programming (LP)

Primal problem

$$\begin{aligned} \text{(P)} \quad & \min \quad \mathbf{c}^\top \mathbf{x}, \\ & \text{s.t.} \quad \mathbf{Ax} = \mathbf{b}, \mathbf{x} \geq 0. \end{aligned}$$

Dual problem

$$\begin{aligned} \text{(D)} \quad & \max \quad \mathbf{b}^\top \mathbf{y}, \\ & \text{s.t.} \quad \mathbf{A}^\top \mathbf{y} \leq \mathbf{c}. \end{aligned}$$

Quadratic Programming (QP), Second Order Cone Programming (SOCP), Semidefinite Programming (SDP), etc...

Convex optimization = standard forms (boring?)

Example: Linear Programming (LP)

Primal problem

$$\begin{aligned} \text{(P)} \quad & \min \quad \mathbf{c}^\top \mathbf{x}, \\ & \text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{x} \geq 0. \end{aligned}$$

Dual problem

$$\begin{aligned} \text{(D)} \quad & \max \quad \mathbf{b}^\top \mathbf{y}, \\ & \text{s.t.} \quad \mathbf{A}^\top \mathbf{y} \leq \mathbf{c}. \end{aligned}$$

Quadratic Programming (QP), Second Order Cone Programming (SOCP), Semidefinite Programming (SDP), etc...

- **Pro:** “Efficient” (but complicated) solvers are already available.
- **Con:** Have to *rewrite* your problem into one of them.

Easy problems (that we don't discuss)

- Objective f is differentiable & no constraint
 - ▶ L-BFGS quasi-Newton method
 - ★ requires only gradient.
 - ★ scales well.
 - ▶ Newton's method
 - ★ requires also Hessian.
 - ★ very accurate.
 - ★ for medium sized problems.
- Differentiable f & simple box constraint
 - ▶ L-BFGS-B quasi-Newton method

Non-differentiability is everywhere

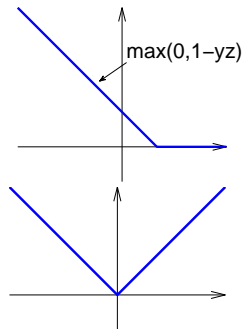
- Support Vector Machine

$$\underset{\mathbf{w}}{\text{minimize}} \quad C \sum_{i=1}^m \ell_H(y_i \langle \mathbf{x}_i, \mathbf{w} \rangle) + \frac{1}{2} \|\mathbf{w}\|^2$$

- Lasso (least absolute **shrinkage** and **selection** operator)

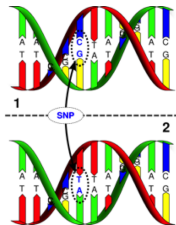
$$\underset{\mathbf{w}}{\text{minimize}} \quad L(\mathbf{w}) + \lambda \sum_{j=1}^n |w_j|$$

⇒ Leads to sparse (most of w_j will be zero) solutions

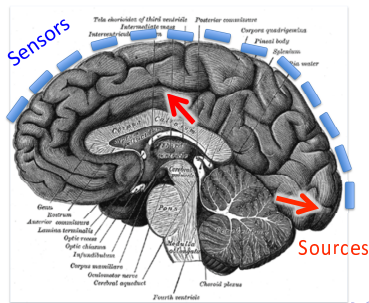


Why we need sparsity

- Genome-wide association studies
 - ▶ Hundreds of thousands of genetic variations (SNPs), small number of participants (samples).
 - ▶ Number of genes responsible for the disease is small.
 - ▶ Solve classification problem (disease/healthy) with sparsity constraint.
- EEG/MEG source localization
 - ▶ Number of possible sources \gg number of sensors
 - ▶ Needs sparsity at a group level



$$\phi_\lambda(\mathbf{w}) = \lambda \sum_{g \in \mathcal{G}} \|\mathbf{w}_g\|_2$$
$$(\mathbf{w}_g \in \mathbb{R}^3)$$

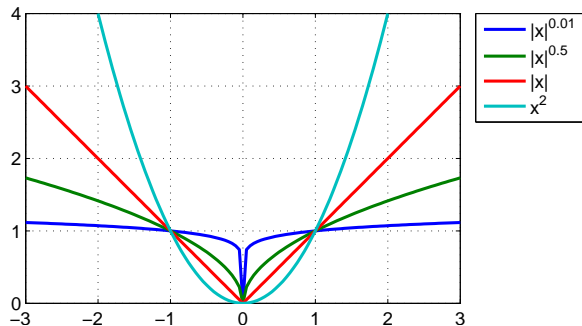


L1-regularization and sparsity

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

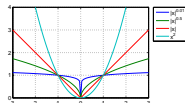
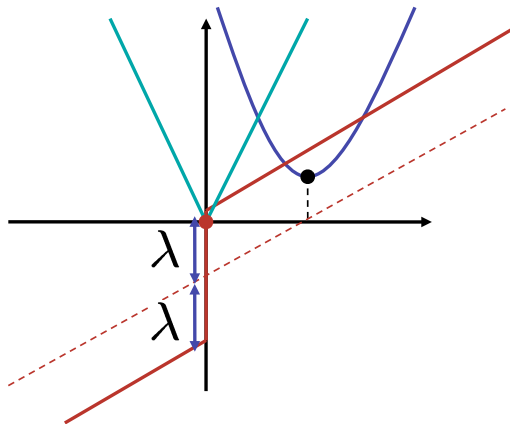
L1-regularization and sparsity

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



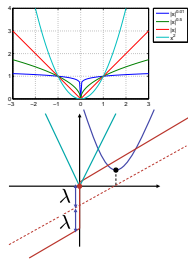
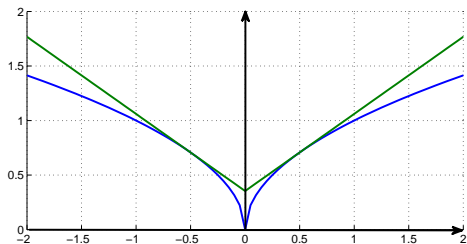
L1-regularization and sparsity

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



L1-regularization and sparsity

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



Multiple kernels & multiple tasks

- Multiple kernel learning [Lanckriet et al., 04; Bach et al., 04;...]
 - ▶ Given: kernel functions $k_1(x, x'), \dots, K_M(x, x')$
 - ▶ How do we optimally select and combine “good” kernels?

$$\begin{array}{l} \text{minimize} \\ f_1 \in \mathcal{H}_1, \\ f_2 \in \mathcal{H}_2, \\ \dots, f_M \in \mathcal{H}_M \end{array} C \sum_{i=1}^N \ell \left(y_i \sum_{m=1}^M f_m(x_i) \right) + \lambda \sum_{m=1}^M \|f_m\|_{\mathcal{H}_m}$$

Multiple kernels & multiple tasks

- Multiple kernel learning [Lanckriet et al., 04; Bach et al., 04;...]
 - ▶ Given: kernel functions $k_1(x, x'), \dots, K_M(x, x')$
 - ▶ How do we optimally select and combine “good” kernels?

$$\underset{\substack{f_1 \in \mathcal{H}_1, \\ f_2 \in \mathcal{H}_2, \\ \dots, f_M \in \mathcal{H}_M}}{\text{minimize}} \quad C \sum_{i=1}^N \ell \left(y_i \sum_{m=1}^M f_m(x_i) \right) + \lambda \sum_{m=1}^M \|f_m\|_{\mathcal{H}_m}$$

- Multiple task learning [Evgeniou et al 05]
 - ▶ Given: two learning tasks.
 - ▶ Can we do better than solving them individually?

$$\underset{\mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_{12}}{\text{minimize}} \quad \underbrace{L_1(\mathbf{w}_1 + \mathbf{w}_{12})}_{\text{Task 1 loss}} + \underbrace{L_2(\mathbf{w}_2 + \mathbf{w}_{12})}_{\text{Task 2 loss}} + \lambda (\|\mathbf{w}_1\| + \|\mathbf{w}_2\| + \|\mathbf{w}_{12}\|)$$

\mathbf{w}_{12} : shared component, \mathbf{w}_1 : Task 1 only component, \mathbf{w}_2 : Task 2 only component.

Estimation of low-rank matrices (1/2)

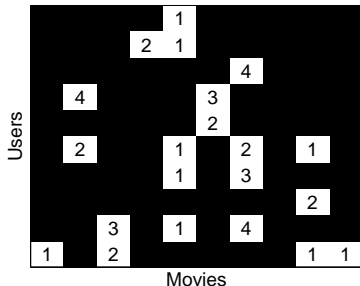
- Completion of partially observed low-rank matrix

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

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

Linear sum of singular-values \Rightarrow sparsity in the singular-values.

- ▶ Collaborative filtering (netflix)
- ▶ Sensor network localization



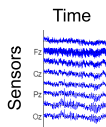
Estimation of low-rank matrices (2/2)

- Classification of matrix shaped data \mathbf{X} .

$$f(\mathbf{X}) = \langle \mathbf{W}, \mathbf{X} \rangle + b$$

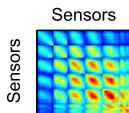
- Multivariate Time Series

$$\mathbf{X} =$$



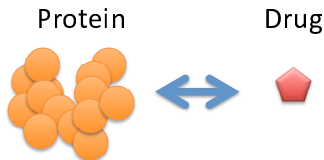
- Second order statistics

$$\mathbf{X} =$$



- Classification of binary relationship between two objects (e.g., protein and drug)

$$f(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{W} \mathbf{y} + b$$



Agenda

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

Convexity

Learning objectives

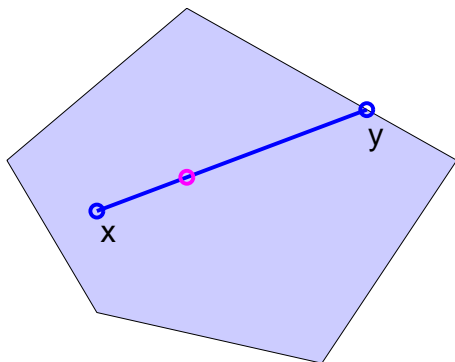
- Convex sets
- Convex function
- Conditions that guarantee convexity
- Convex optimization problem

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,

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



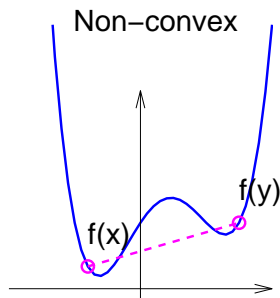
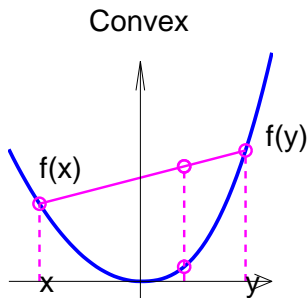
Convex function

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

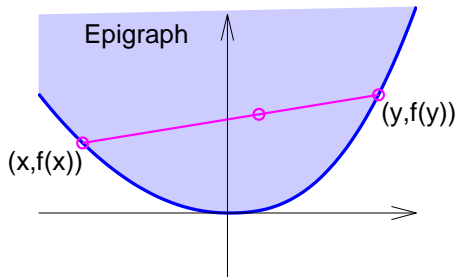
NB: when the strict inequality $<$ holds, f is called **strictly convex**.

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.



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

Exercise

- Show that the indicator function $\delta_C(\mathbf{x})$ of a convex set C is a convex function. Here

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

Conditions that guarantee convexity (1/3)

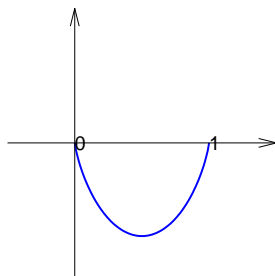
- Hessian $\nabla^2 f(\mathbf{x})$ is positive semidefinite (if f is differentiable)

Examples

- ▶ (Negative) entropy is a convex function.

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

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



Conditions that guarantee convexity (1/3)

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

Examples

- ▶ (Negative) entropy is a convex function.

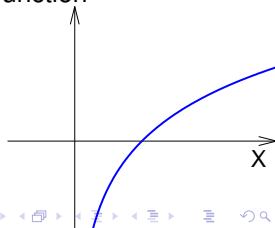
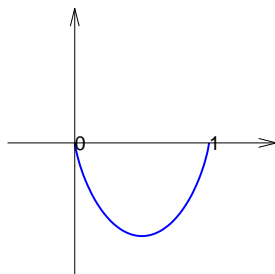
$$f(p) = \sum_{i=1}^n p_i \log p_i,$$

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

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

$$f(\mathbf{X}) = \log |\mathbf{X}| \quad (\mathbf{X} \succeq 0),$$

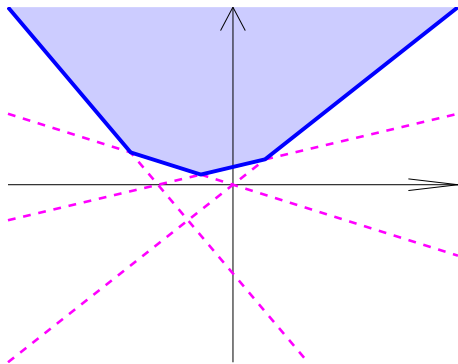
$$\nabla^2 f(\mathbf{X}) = -\mathbf{X}^{-\top} \otimes \mathbf{X}^{-1} \preceq 0$$



Conditions that guarantee convexity (2/3)

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

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



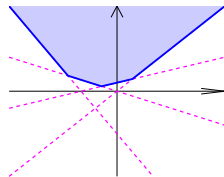
The same as saying “intersection of convex sets is a convex set”

Conditions that guarantee convexity (2/3)

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

$$f(\mathbf{x}) := \max_{\boldsymbol{\alpha} \in \mathbb{R}^n} f(\mathbf{x}; \boldsymbol{\alpha})$$

Example



- Quadratic over linear is a convex function

$$f(\mathbf{y}, \boldsymbol{\Sigma}) = \max_{\boldsymbol{\alpha}} \left[-\frac{1}{2} \boldsymbol{\alpha}^\top \boldsymbol{\Sigma} \boldsymbol{\alpha} + \boldsymbol{\alpha}^\top \mathbf{y} \right] \quad (\boldsymbol{\Sigma} \succ 0)$$

Conditions that guarantee convexity (2/3)

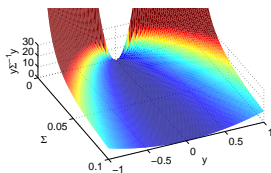
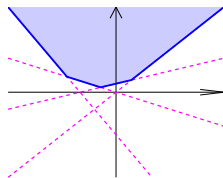
- Maximum over convex functions $\{f(\mathbf{x}; \alpha) : \alpha \in \mathbb{R}^n\}$

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

Example

- Quadratic over linear is a convex function

$$\begin{aligned} f(\mathbf{y}, \Sigma) &= \max_{\alpha} \left[-\frac{1}{2} \alpha^{\top} \Sigma \alpha + \alpha^{\top} \mathbf{y} \right] \quad (\Sigma \succ 0) \\ &= \frac{1}{2} \mathbf{y}^{\top} \Sigma^{-1} \mathbf{y} \end{aligned}$$



Conditions that guarantee convexity (3/3)

- Minimum of **jointly convex** function $f(x, y)$

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

Examples

- ▶ Hierarchical prior minimization

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

Conditions that guarantee convexity (3/3)

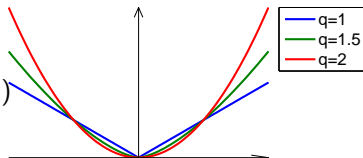
- Minimum of **jointly convex** function $f(x, y)$

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

Examples

- ▶ Hierarchical prior minimization

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



Conditions that guarantee convexity (3/3)

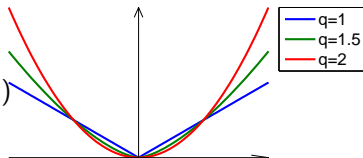
- Minimum of **jointly convex** function $f(x, y)$

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

Examples

- Hierarchical prior minimization

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



- Schatten 1- norm (sum of singularvalues)

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

Conditions that guarantee convexity (3/3)

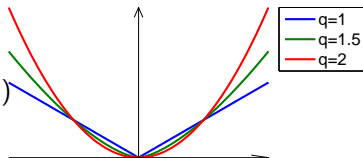
- Minimum of **jointly convex** function $f(x, y)$

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

Examples

- Hierarchical prior minimization

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



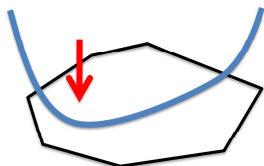
- Schatten 1- norm (sum of singularvalues)

$$\begin{aligned} f(\mathbf{X}) &= \min_{\Sigma \geq 0} \frac{1}{2} \left(\text{Tr}(\mathbf{X}\Sigma^{-1}\mathbf{X}^\top) + \text{Tr}(\Sigma) \right) \\ &= \text{Tr} \left((\mathbf{X}^\top \mathbf{X})^{1/2} \right) = \sum_{j=1}^r \sigma_j(\mathbf{X}). \end{aligned}$$

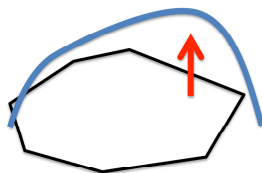
Convex optimization problem

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!

Proximity operators and IST methods

Learning objectives

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

Proximity view on gradient descent

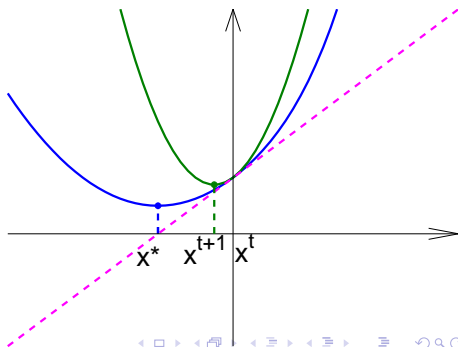
“Linearize and Prox”

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

- 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?

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

Constraint minimization problem

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

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

- can be equivalently written as

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

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

Projected gradient method (Bertsekas 99; Nesterov 03)

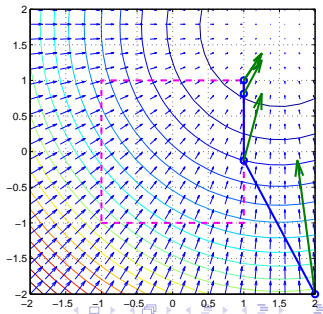
Linearize the objective f , δ_C is the indicator of the constraint C

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

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

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

- Need the projection proj_C to be easy to compute



Ideas for regularized minimization

Constrained minimization problem

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

⇒ need to compute the **projection**

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

Regularized minimization problem

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}) + \phi_\lambda(\mathbf{x})$$

⇒ need to compute the **proximity operator**

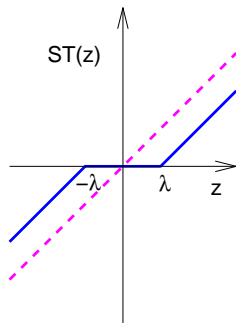
$$\mathbf{x}^{t+1} = \underset{\mathbf{x}}{\text{argmin}} \left(\phi_\lambda(\mathbf{x}) + \frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{y}\|_2^2 \right)$$

Proximal Operator: generalization of projection

$$\text{prox}_{\phi_\lambda}(\mathbf{z}) = \underset{\mathbf{x}}{\text{argmin}} \left(\phi_\lambda(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{z}\|_2^2 \right)$$

- $\phi_\lambda = \delta_C$: Projection onto a convex set
 $\text{prox}_{\delta_C}(\mathbf{z}) = \text{proj}_C(\mathbf{z})$.
- $\phi_\lambda(\mathbf{x}) = \lambda \|\mathbf{x}\|_1$: Soft-Threshold

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



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

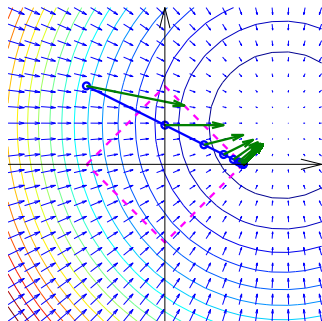
Iterative Shrinkage Thresholding (IST)

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

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

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

- 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

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

by iteratively computing

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

Exercise: Derive prox operator for

- Ridge regularization

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

- Elastic-net regularization

$$\phi_\lambda(\mathbf{w}) = \lambda \sum_{j=1}^n \left((1 - \theta) |w_j| + \theta w_j^2 \right).$$

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 hints

- 1 Compute the gradient of the loss term

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

- 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:

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

Regularization:

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

Matrix completion via IST (Mazumder et al. 10)

Loss function:

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

Regularization:

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

gradient:

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

Matrix completion via IST (Mazumder et al. 10)

Loss function:

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

gradient:

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

Regularization:

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

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:

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

gradient:

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

Regularization:

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

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.

FISTA: accelerated version of IST (Beck & Teboulle 09;

Nesterov 07)

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

$$\mathbf{x}^t = \text{prox}_{\lambda\eta_t}(\mathbf{y}^t - \eta_t \nabla L(\mathbf{y}^t)).$$

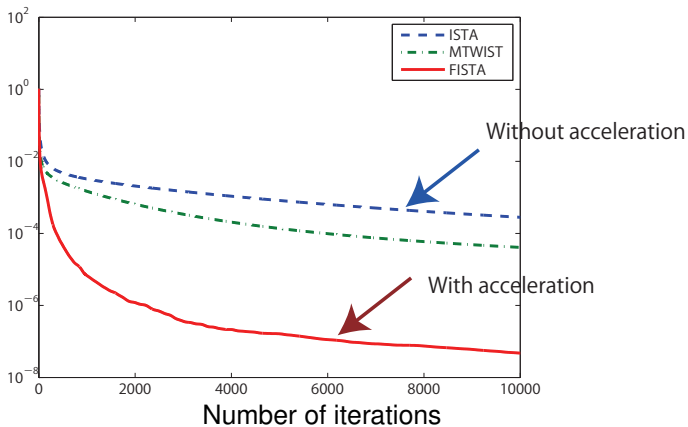
- 3 Update \mathbf{y}^t :

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

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{y}^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

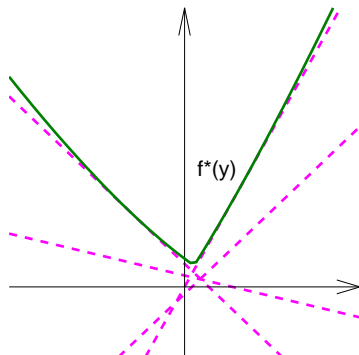
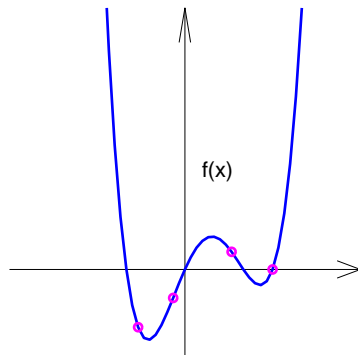
Conjugate duality and dual ascent

- Convex conjugate function
- Lagrangian relaxation and dual problem
- Dual ascent

Conjugate duality

The convex conjugate f^* of a function f :

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



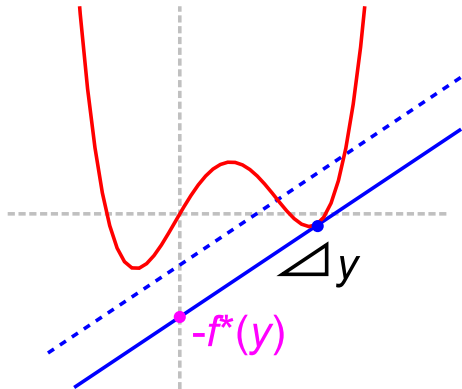
Since the maximum over linear functions is always convex, f need not be convex.

Conjugate duality (dual view)

Convex conjugate function

$-f^*(\mathbf{y})$ is the minimum y -intercept of the hyperplanes that has slope \mathbf{y} and have intersection with the graph of $f(\mathbf{x})$.

$$\begin{aligned} f^*(\mathbf{y}) &= \sup_{\mathbf{x}} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x})) \\ \Leftrightarrow -f^*(\mathbf{y}) &= \inf_{\mathbf{x}} (f(\mathbf{x}) - \langle \mathbf{x}, \mathbf{y} \rangle) \\ &= \inf_{\mathbf{x}, b} b, \\ \text{s.t. } f(\mathbf{x}) &= \langle \mathbf{x}, \mathbf{y} \rangle + b. \end{aligned}$$

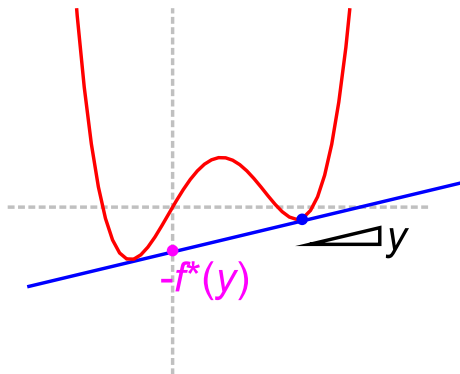


Conjugate duality (dual view)

Convex conjugate function

$-f^*(\mathbf{y})$ is the minimum y -intercept of the hyperplanes that has slope \mathbf{y} and have intersection with the graph of $f(\mathbf{x})$.

$$\begin{aligned} f^*(\mathbf{y}) &= \sup_{\mathbf{x}} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x})) \\ \Leftrightarrow -f^*(\mathbf{y}) &= \inf_{\mathbf{x}} (f(\mathbf{x}) - \langle \mathbf{x}, \mathbf{y} \rangle) \\ &= \inf_{\mathbf{x}, b} b, \\ \text{s.t. } f(\mathbf{x}) &= \langle \mathbf{x}, \mathbf{y} \rangle + b. \end{aligned}$$

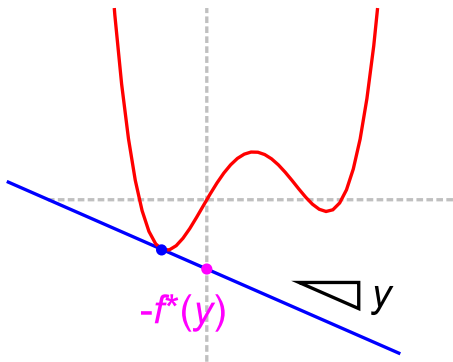


Conjugate duality (dual view)

Convex conjugate function

$-f^*(\mathbf{y})$ is the minimum y -intercept of the hyperplanes that has slope \mathbf{y} and have intersection with the graph of $f(\mathbf{x})$.

$$\begin{aligned} f^*(\mathbf{y}) &= \sup_{\mathbf{x}} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x})) \\ \Leftrightarrow -f^*(\mathbf{y}) &= \inf_{\mathbf{x}} (f(\mathbf{x}) - \langle \mathbf{x}, \mathbf{y} \rangle) \\ &= \inf_{\mathbf{x}, b} b, \\ \text{s.t. } f(\mathbf{x}) &= \langle \mathbf{x}, \mathbf{y} \rangle + b. \end{aligned}$$

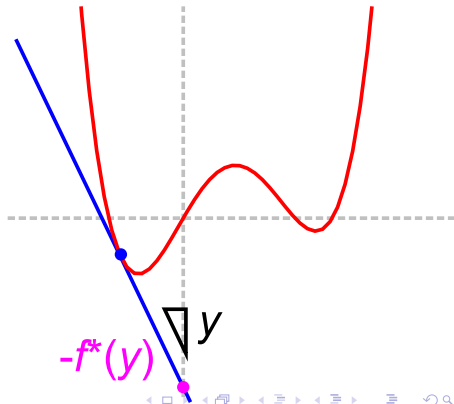


Conjugate duality (dual view)

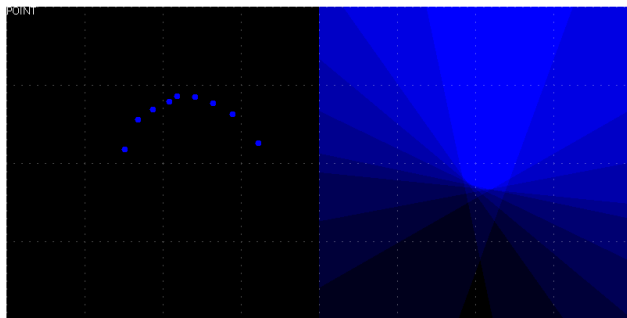
Convex conjugate function

$-f^*(\mathbf{y})$ is the minimum y -intercept of the hyperplanes that has slope \mathbf{y} and have intersection with the graph of $f(\mathbf{x})$.

$$\begin{aligned} f^*(\mathbf{y}) &= \sup_{\mathbf{x}} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x})) \\ \Leftrightarrow -f^*(\mathbf{y}) &= \inf_{\mathbf{x}, b} (f(\mathbf{x}) - \langle \mathbf{x}, \mathbf{y} \rangle) \\ &= \inf_{\mathbf{x}, b} b, \\ \text{s.t. } f(\mathbf{x}) &= \langle \mathbf{x}, \mathbf{y} \rangle + b. \end{aligned}$$



Demo

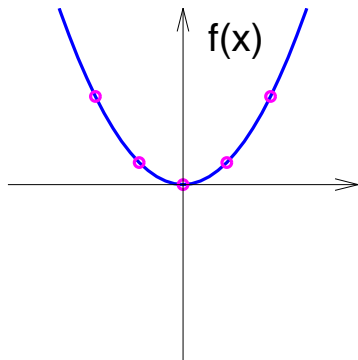


<http://www.ibis.t.u-tokyo.ac.jp/ryotat/applets/pld/>

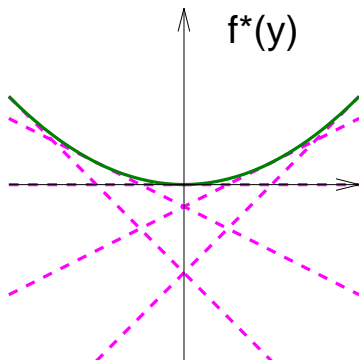
Example of conjugate duality $f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \mathbb{R}^n} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}))$

- Quadratic function

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



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



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

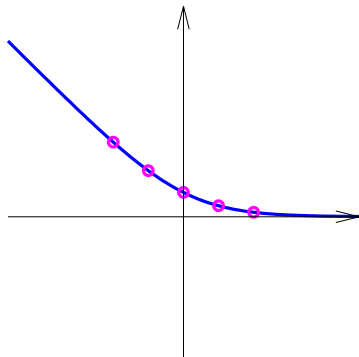
- Logistic loss function

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

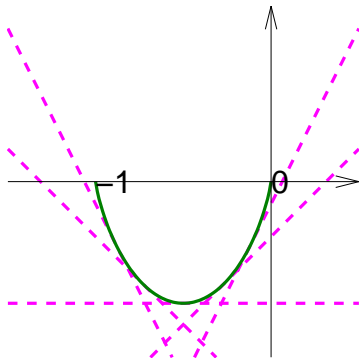
Example of conjugate duality $f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \mathbb{R}^n} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}))$

- Logistic loss function

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



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



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

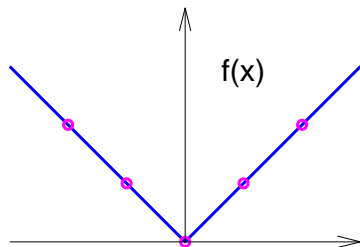
- L1 regularizer

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

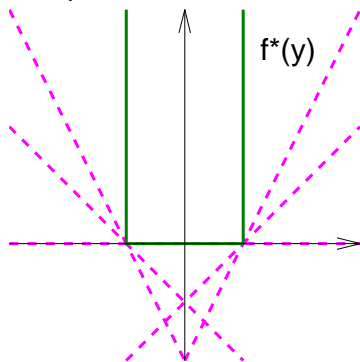
Example of conjugate duality $f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \mathbb{R}^n} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}))$

- L1 regularizer

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

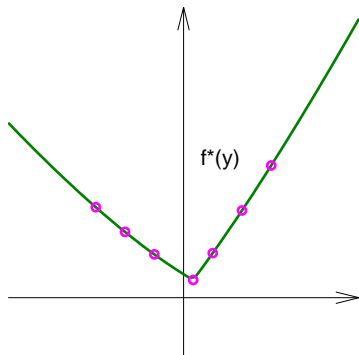
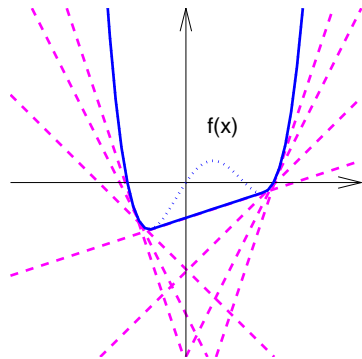


$$f^*(y) = \begin{cases} 0 & (-1 \leq y \leq 1) \\ +\infty & (\text{otherwise}) \end{cases}$$



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.
- Minimum of \mathcal{L} is no greater than the minimum of the original.

Weak duality

$$\inf_{\mathbf{z}, \mathbf{w}} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha) \leq p^* \quad (\text{primal optimal})$$

proof

$$\inf_{\mathbf{z}, \mathbf{w}} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha) = \inf \left(\inf_{\mathbf{z}=\mathbf{Aw}} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha), \inf_{\mathbf{z} \neq \mathbf{Aw}} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha) \right)$$

Weak duality

$$\inf_{\mathbf{z}, \mathbf{w}} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha) \leq p^* \quad (\text{primal optimal})$$

proof

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

Weak duality

$$\inf_{\mathbf{z}, \mathbf{w}} \mathcal{L}(\mathbf{z}, \mathbf{w}, \alpha) \leq p^* \quad (\text{primal optimal})$$

proof

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

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 α . Why don't we maximize over \mathbf{w} ?

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 α . Why don't we maximize over \mathbf{w} ?

Dual problem

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

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.

Dual problem

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

Dual problem

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

Dual problem

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

Dual problem

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

Dual problem

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

Fenchel's duality



M. W. Fenchel

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

Examples

- Logistic regression with L1 regularization

$$f(\mathbf{z}) = \sum_{i=1}^m \log(1 + \exp(-z_i)), \quad g(\mathbf{w}) = \lambda \|\mathbf{w}\|_1.$$

- Support vector machine (SVM)

$$f(\mathbf{z}) = C \sum_{i=1}^m \max(0, 1 - z_i), \quad g(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_2^2.$$

Example 1: Logistic regression with L1 regularization

Primal

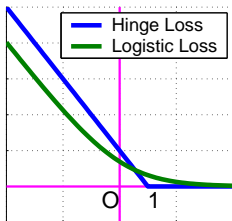
$$\min_{\mathbf{w}} f(\mathbf{y} \circ \mathbf{X}\mathbf{w}) + \phi_{\lambda}(\mathbf{w})$$

$$\begin{cases} f(\mathbf{z}) = \sum_{i=1}^m \log(1 + \exp(-z_i)), \\ \phi_{\lambda}(\mathbf{w}) = \lambda \|\mathbf{w}\|_1. \end{cases}$$

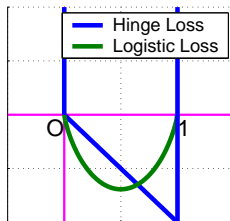
Dual

$$\max_{\alpha} -f^*(-\alpha) - \phi_{\lambda}^*(\mathbf{X}^T(\alpha \circ \mathbf{y}))$$

$$\begin{cases} f^*(-\alpha) = \sum_{i=1}^m \alpha_i \log(\alpha_i) \\ \quad + (1 - \alpha_i) \log(1 - \alpha_i), \\ \phi_{\lambda}^*(\mathbf{v}) = \begin{cases} 0 & (\|\mathbf{w}\|_{\infty} \leq \lambda), \\ +\infty & (\text{otherwise}). \end{cases} \end{cases}$$



(a) primal losses



(b) dual losses

Example 2: Support vector machine

Primal

$$\min_{\mathbf{w}} f(\mathbf{y} \circ \mathbf{X}\mathbf{w}) + \phi_{\lambda}(\mathbf{w})$$

$$\begin{cases} f(\mathbf{z}) = C \sum_{i=1}^m \max(0, 1 - z_i), \\ \phi_{\lambda}(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2. \end{cases}$$

Dual

$$\max_{\alpha} -f^*(-\alpha) - \phi_{\lambda}^*(\mathbf{X}^{\top}(\alpha \circ \mathbf{y}))$$

$$\begin{cases} f^*(-\alpha) = \begin{cases} \sum_{i=1}^m -\alpha_i & (0 \leq \alpha \leq C), \\ +\infty & (\text{otherwise}), \end{cases} \\ \phi_{\lambda}^*(\mathbf{v}) = \frac{1}{2} \|\mathbf{v}\|^2. \end{cases}$$

Dual ascent

Assume for a moment that the dual $d(\alpha)$ is differentiable.

For a given α^t

$$d(\alpha^t) = \inf_{\mathbf{z}, \mathbf{w}} (f(\mathbf{z}) + g(\mathbf{w}) + \langle \alpha^t, \mathbf{z} - \mathbf{A}\mathbf{w} \rangle)$$

and one can show that (Chapter 6, Bertsekas 99)

$$\nabla_{\alpha} d(\alpha^t) = \mathbf{z}^{t+1} - \mathbf{A}\mathbf{w}^{t+1}$$

where

$$\mathbf{z}^{t+1} = \underset{\mathbf{z}}{\operatorname{argmin}} (f(\mathbf{z}) + \langle \alpha^t, \mathbf{z} \rangle)$$

$$\mathbf{w}^{t+1} = \underset{\mathbf{w}}{\operatorname{argmin}} (g(\mathbf{w}) - \langle \mathbf{A}^{\top} \alpha^t, \mathbf{w} \rangle)$$

Dual ascent (Uzawa's method)

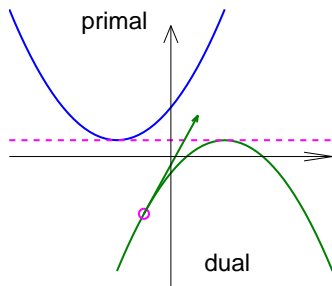
$$\left\{ \begin{array}{l} \text{Minimize the Lagrangian wrt } \mathbf{x} \text{ and } \mathbf{z}: \\ \mathbf{z}^{t+1} = \operatorname{argmin}_{\mathbf{z}} (f(\mathbf{z}) + \langle \alpha^t, \mathbf{z} \rangle), \\ \mathbf{w}^{t+1} = \operatorname{argmin}_{\mathbf{w}} (g(\mathbf{w}) - \langle \mathbf{A}^T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

- **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



Exercise 2: Matrix completion via dual ascent (Cai et al. 08)

$$\begin{aligned} \underset{\mathbf{X}}{\text{minimize}} \quad & \underbrace{\frac{1}{2\lambda} \|\mathbf{z} - \mathbf{y}\|^2}_{\text{Strictly convex}} + \underbrace{\left(\tau \|\mathbf{X}\|_{\text{tr}} + \frac{1}{2} \|\mathbf{X}\|^2 \right)}_{\text{Strictly convex}}, \\ \text{s.t.} \quad & \Omega(\mathbf{X}) = \mathbf{z}. \end{aligned}$$

Exercise 2: Matrix completion via dual ascent (Cai et al. 08)

$$\begin{aligned} \underset{\mathbf{X}}{\text{minimize}} \quad & \underbrace{\frac{1}{2\lambda} \|\mathbf{z} - \mathbf{y}\|^2}_{\text{Strictly convex}} + \underbrace{\left(\tau \|\mathbf{X}\|_{\text{tr}} + \frac{1}{2} \|\mathbf{X}\|^2 \right)}_{\text{Strictly convex}}, \\ \text{s.t.} \quad & \Omega(\mathbf{X}) = \mathbf{z}. \end{aligned}$$

⇓

Lagrangian:

$$\mathcal{L}(\mathbf{X}, \mathbf{z}, \alpha) = \underbrace{\frac{1}{2\lambda} \|\mathbf{z} - \mathbf{y}\|^2}_{=f(\mathbf{z})} + \underbrace{\left(\tau \|\mathbf{X}\|_{S_1} + \frac{1}{2} \|\mathbf{X}\|^2 \right)}_{=g(\mathbf{X})} + \alpha^\top (\mathbf{z} - \Omega(\mathbf{X})).$$

Dual ascent

$$\begin{cases} \mathbf{X}^{t+1} = \text{prox}_\tau (\Omega^\top(\alpha^t)) & (\text{Singular-Value Thresholding}) \\ \mathbf{z}^{t+1} = \mathbf{y} - \lambda \alpha^t \\ \alpha^{t+1} = \alpha^t + \eta_t (\mathbf{z}^{t+1} - \Omega(\mathbf{X}^{t+1})) \end{cases}$$

Augmented Lagrangian and ADMM

Learning objectives

- Structured sparse estimation
- Augmented Lagrangian
- Alternating direction method of multipliers

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

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

Original X_0



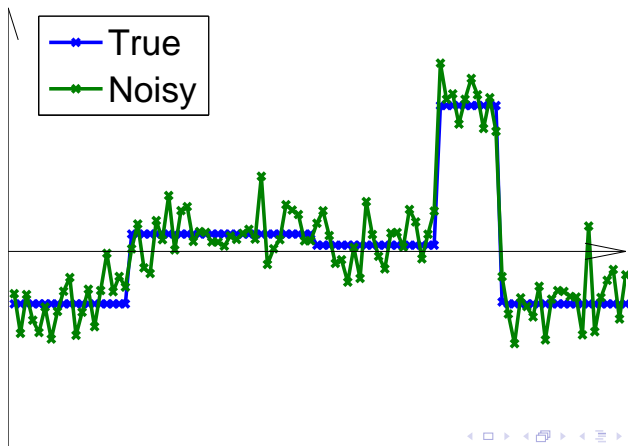
Observed Y



In one dimension

- Fused lasso [Tibshirani et al. 05]

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \sum_{j=1}^{n-1} |x_{j+1} - x_j|$$



Structured sparsity estimation

- TV denoising

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

- Fused lasso

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \sum_{j=1}^{n-1} |x_{j+1} - x_j|$$

Structured sparsity estimation

- TV denoising

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

- Fused lasso

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \sum_{j=1}^{n-1} |x_{j+1} - x_j|$$

Structured sparse estimation problem

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

Structured sparse estimation problem

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad \underbrace{f(\mathbf{x})}_{\text{data-fit}} + \underbrace{\phi_\lambda(\mathbf{Ax})}_{\text{regularization}}$$

- Not easy to compute prox operator (because it is **non-separable**)
⇒ difficult to apply **IST-type methods**.
- Dual is not necessarily differentiable
⇒ difficult to apply **dual ascent**.

Forming the *augmented* Lagrangian

Structured sparsity problem

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad \underbrace{f(\mathbf{x})}_{\text{data-fit}} + \underbrace{\phi_\lambda(\mathbf{Ax})}_{\text{regularization}}$$

Equivalently written as

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

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

Forming the *augmented* Lagrangian

Structured sparsity problem

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad \underbrace{f(\mathbf{x})}_{\text{data-fit}} + \underbrace{\phi_\lambda(\mathbf{Ax})}_{\text{regularization}}$$

Equivalently written as

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

Augmented Lagrangian function

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

Augmented Lagrangian Method

Augmented Lagrangian function

$$\mathcal{L}_\eta(\mathbf{x}, \mathbf{z}, \boldsymbol{\alpha}) = f(\mathbf{x}) + \phi_\lambda(\mathbf{z}) + \boldsymbol{\alpha}^\top (\mathbf{z} - \mathbf{Ax}) + \frac{\eta}{2} \|\mathbf{z} - \mathbf{Ax}\|^2.$$

Augmented Lagrangian method (Hestenes 69, Powell 69)

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

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

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

- Minimize the AL function $\mathcal{L}_\eta(\mathbf{x}, \mathbf{z}^t, \boldsymbol{\alpha}^t)$ wrt \mathbf{x} :
- Minimize the AL function $\mathcal{L}_\eta(\mathbf{x}^{t+1}, \mathbf{z}, \boldsymbol{\alpha}^t)$ wrt \mathbf{z} :
- Update the Lagrangian multiplier:
 $\boldsymbol{\alpha}^{t+1} = \boldsymbol{\alpha}^t + \eta(\mathbf{z}^{t+1} - \mathbf{A}\mathbf{x}^{t+1})$.

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

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

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

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

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

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

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

Exercise: implement an ADMM for fused lasso

Fused lasso

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{A}\mathbf{x}\|_1$$

- What is the loss function f ?
- What is the regularizer g ?
- What is the matrix \mathbf{A} for fused lasso?
- What is the prox operator for the regularizer g ?

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

- Stopping criterion
 - ▶ Care must be taken when making a comparison.
- 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
 - ▶ Dual problem is always convex. It provides a lower-bound of the original problem. If $p^* = d^*$, you are done!
 - ▶ **Dual ascent** (or dual decomposition) for sequence labeling in natural language processing; see [Wainwright, Jaakkola, Willsky 05; Koo et al. 10]
 - ▶ Difference of convex (DC) programming.
 - ▶ Eigenvalue problem.
- Stochastic optimization
 - ▶ Good tutorial by Nathan Srebro (ICML2010)

A new book “Optimization for Machine Learning” is coming out from the MIT press.

amazon.co.uk Hello. Sign in to get personalised recommendations. New Customer? [Start here](#)

Your Amazon.co.uk Today's Deals Gift Cards Gifts & Wish Lists

Shop All Departments Search books

Books Advanced Search Browse Genres Bestsellers New & Future Releases Paperbacks Seasonal Offers

Optimization for Machine Learning (Neural Information Processing Series) [Hardcover]
Paul H. Geor (Author)

No image available

Price: **£34.95**
Price: **£33.20** & this item **Delivered FREE in the UK** with Super Saver Delivery. [See details and conditions](#)

You Save: **£1.75 (5%)**

Pre-order Price Guarantee. [Learn more.](#)

This title has not yet been released.
You may pre-order it now and we will deliver it to you when it arrives.
Dispatched from and sold by **Amazon.co.uk**. Gift-wrap available.

Seasonal Offers See great savings on 1000s of books in our [Seasonal Offers](#)

[Publisher: learn how customers can search inside this book.](#)

Contributed authors including: A. Nemirovksi, D. Bertsekas, L. Vandenberghe, and more.

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.