



1.5 : Algorithms Primer

① Unconstrained Optimization

For convexity: $\nabla f(x)^T \Delta x < 0$

step size
search direction

$$A: \text{Descent Method: } x^{k+1} = x^k + t^k \Delta x^k, \quad f(x^{k+1}) < f(x^k)$$



line search

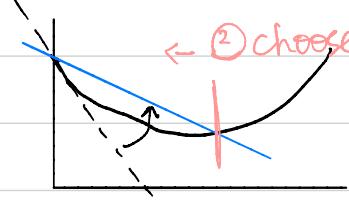
1) Exact line search: $t = \underset{t > 0}{\operatorname{argmin}} f(x + t \Delta x)$

2) Backtracking: 1. start at $t = 1$ (full step)

2. until $f(x + t \Delta x) < f(x) + \alpha t \nabla f(x)^T \Delta x$

3. $t := \beta t, \beta \in (0, 1)$

① lower slope (cut through func): smaller α



Gradient Descent Method: $\Delta x = -\nabla f(x)$

→ needs twice-differentiable func.

B: Newton's Method: Minimizing the second order approximation

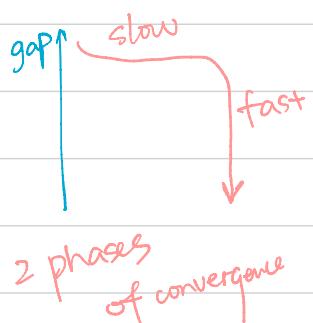
Derivation: $\hat{f}(x+v) = f(x) + \nabla f(x)^T v + \frac{1}{2} v^T \nabla^2 f(x) v$

choose v to minimize \hat{f} : $D\hat{f} = Df + D^2 v = 0 \Rightarrow v = -D^{-2} Df$

Steps: 1) Compute Newton Step and decrement

$$\Delta x_{nt}^k = -\nabla^2 f(x^k)^{-1} \nabla f(x^k), \quad \lambda(x^k)^2 = -\nabla f(x^k)^T \Delta x_{nt}^k$$

tells how far from optimum



2) stop criterion: if $\lambda(x^k)^2 / 2 \leq \epsilon$

3) line search

4) update: $x^{k+1} = x^k + t^k \Delta x_{nt}^k$

5) $k = k + 1$

② Constrained Optimization

→ or KKT

A: Equality Constrained Optimization (Eliminating equality constraints)

$$Ax = b \Rightarrow Fz + x_0, F \text{ is nullspace of } A \text{ i.e. } AF = 0$$

$$\begin{cases} \min f(x) \\ \text{s.t. } Ax = b \end{cases} \rightarrow \min \tilde{f}(z) = f(Fz + x_0) \text{ is convex} \because f \text{ is assumed convex}$$

↳ unconstrained: Newton

$$\text{chain rule: } \nabla \tilde{f}(z) = F^T \nabla f(x)$$

$$\text{"also gradient on } z \text{ is OK; in } B \text{ is gradient on } x" \quad \nabla^2 \tilde{f}(z) = F^T \nabla^2 f(x) F$$

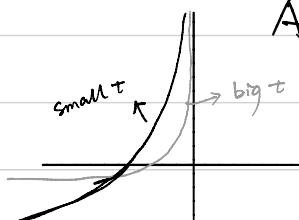
L5 : cont'd.

B. Gradient Projection Method (both equality or inequality constraints)

$$\min f(x) \Rightarrow x^{k+1} = [x^k - \alpha^k \nabla f(x^k)]_X^{\text{projection (closest feasible point)}} \\ \text{s.t. } x \in X \quad \downarrow \\ \text{another optimization prob.}$$

C. Interior - Point Methods (IPM)

$$\begin{aligned} \min_x f_0(x) & \quad \min_x f_0(x) + \sum_{i=1}^m I_-(f_i(x)) \\ \text{s.t. } f_i(x) \leq 0 & \quad \text{s.t. } Ax = b \\ Ax = b & \quad \text{indicator func: } \begin{cases} 0, \text{ if } u \leq 0 \\ \infty, \text{ otherwise.} \end{cases} \\ -\left(\frac{1}{t}\right) \log(-x) & \quad \text{approximate with a smooth function} \end{aligned}$$



Barrier method:

= called Interior - Point

- Step:
- ① Very initial point \rightarrow each iteration the point is feasible
 - ② Start with small t (at the begining, large $t \rightarrow$ Newton slow converge)
 - ③ use Newton (fast phase better), use last x as initial point
 - ④ stop at m/t (duality gap) $< \epsilon$
 - ⑤ increase t

Feasibility problem: find x such that $f_i(x) \leq 0, i = 1, \dots, m; Ax = b$

Phase 1 Method: minimize s

x, s

s.t. $f_i(x) \leq s; Ax = b$

choose a $s > \max f_i(x)$ and decrease. If feasible $s^* < 0$, then x^* feasible

initial point for this optimization problem (IPM to solve)

L5 . cont'd

③ Block Coordinate Algorithms

A: Block Coordinate Descent

update parameters block sequentially

Feature: 1) non-increasing objective value

2) each subproblem may be much easier to solve (even closed-form)

Application: Lasso ($L_2 - L_1$ optimization) via BCD

$$\text{prob: } \underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) \triangleq \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|,$$

closed form BCD: $\underset{x_i}{\text{minimize}} \quad f_i(x_i) \triangleq \frac{1}{2} \|\tilde{\mathbf{y}}_i^k - a_i x_i\|_2^2 + \lambda |x_i|$

vector \mathbf{x}

scalar λ

$$\text{where } \tilde{\mathbf{y}}_i^k \triangleq \mathbf{y} - \sum_{j < i} a_j x_j^{k+1} - \sum_{j > i} a_j x_j^k$$

(将相乘运算拆成 scalar 形式，求偏导)

closed form solution: $x_i^{k+1} = \text{soft}_\lambda(a_i^\top \tilde{\mathbf{y}}_i^k) / \|a_i\|^2$

this can also be solved by MM

B. Jacobi Algorithm:

update parameters block in parallel

Feature: 1) no non-increasing objective value

L9 Portfolio Optimization

Premier on Financial Data

log-prices $y_t \triangleq \log P_t$. follows random walk $y_t = \mu + y_{t-1} + \varepsilon_t$

simple return $R_t \triangleq \frac{P_t - P_{t-1}}{P_{t-1}} = \frac{P_t}{P_{t-1}} - 1$

log-return $r_t \triangleq y_t - y_{t-1} = \log\left(\frac{P_t}{P_{t-1}}\right) = \log(1 + R_t)$ 插绘 price 差分

why log-return?

- ① prices are assumed to be distributed log normally $\Rightarrow r_t$ normally distributed
- ② $r_t \approx R_t$ when returns are small. good approximation
- ③ Compound return: $\prod (1+R_i) = (1+R_1)(1+R_2)\dots(1+R_n)$
 $\log \prod (1+R_i) = r_1 + r_2 + \dots + r_n \rightarrow$ time-additivity, normally distributed

各种 model.

i.i.d 最简单 $\begin{cases} \mu_t \triangleq E[r_t | F_{t-1}] = \mu \\ \Sigma_t \triangleq Cov[r_t | F_{t-1}] = E[(r_t - \mu_t)(r_t - \mu_t)^T | F_{t-1}] \end{cases}$ 未来 mean, var \rightarrow iid 以为 const \rightarrow sample value 来代替

Portfolio Basics

$$\begin{cases} \text{portfolio } w: \quad \bar{w}^T w = 1 \\ \text{portfolio return: } R_t^P \approx w^T r_t \end{cases}$$

performance measures expected return: $w^T \mu$

volatility: $\sqrt{w^T \Sigma w}$

Sharpe Ratio (SR): $SR = \frac{w^T \mu - r_f}{\sqrt{w^T \Sigma w}}$ \leftarrow risk-free rate

L9 cont'd.

Markowitz's Mean - Variance Portfolio (MVP)

$$\begin{bmatrix} \underset{\omega}{\text{maximize}} & \omega^T \mu - \lambda \omega^T \Sigma \omega \\ \text{s.t.} & \mathbf{1}^T \omega = 1 \end{bmatrix} \quad \text{convex QP}$$

Global minimum variance portfolio (GMVP)

$$\begin{bmatrix} \underset{\omega}{\text{minimize}} & \omega^T \Sigma \omega \\ \text{s.t.} & \mathbf{1}^T \omega = 1 \end{bmatrix} \quad \begin{array}{l} \text{ignores the expected return :} \\ \text{convex QP} \end{array}$$

Maximum Sharpe ratio portfolio (MSRP)

$$\begin{bmatrix} \underset{\omega}{\text{maximize}} & \frac{\omega^T \mu - r_f}{\sqrt{\omega^T \Sigma \omega}} \\ \text{s.t.} & \mathbf{1}^T \omega = 1 \end{bmatrix} \quad \begin{array}{l} \text{nonconvex} \\ \text{fractional programming (FP)} \end{array}$$

$$\begin{bmatrix} \underset{\omega, t}{\text{maximize}} & t \\ \text{s.t.} & t \leq \frac{\omega^T \mu - r_f}{\sqrt{\omega^T \Sigma \omega}} \\ & \mathbf{1}^T \omega = 1 \end{bmatrix} \stackrel{\downarrow}{\Rightarrow} \begin{bmatrix} \text{find } \omega \\ \text{s.t.} & t \|\Sigma^{1/2} \omega\|_2 \leq \omega^T \mu - r_f \\ & \mathbf{1}^T \omega = 1 \end{bmatrix}$$

Bisection + *SOCPr*

Fractional Programming

$$\underset{x}{\text{maximize}} \quad \frac{f(x)}{g(x)} \quad \begin{array}{l} \rightarrow \text{concave} \\ \rightarrow \text{convex} \end{array}$$

solution

$$\text{s.t.} \quad x \in \mathcal{X}$$

① bisection: the problem is quasi-convex \rightsquigarrow epigraph form

$$\begin{array}{ll} \underset{x, t}{\text{maximize}} & t \\ \text{s.t.} & t \leq \frac{f(x)}{g(x)}, x \in \mathcal{X} \end{array} \quad \begin{array}{l} \xrightarrow{\text{fix } t} \\ \xrightarrow{\text{i) feasible: } t \uparrow} \\ \xrightarrow{\text{ii) otherwise: } t \downarrow} \end{array} \quad \begin{array}{ll} \underset{x}{\text{maximize}} & 0 \\ \text{s.t.} & \frac{t g(x)}{f(x)} \leq 1 \end{array} \quad \begin{array}{l} \text{positive} \\ \text{convex problem} \end{array}$$

② Dinkelbach transform

$$\underset{x}{\text{maximize}} \quad f(x) - y g(x) \quad \xrightarrow{\text{iteratively update}} \quad y^{(k)} = \frac{f(x^{(k)})}{g(x^{(k)})}$$

2.11. MM.

① MM Algorithm

MM in a Nutshell:

difficult optimization problem successively minimize a surrogate func.

$$\begin{aligned} \text{minimize}_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{x} \in \mathcal{X} \end{aligned} \Rightarrow \quad \mathbf{x}^{k+1} = \underset{\mathbf{x} \in \mathcal{X}}{\operatorname{argmin}} u(\mathbf{x}, \mathbf{x}^k)$$

rule of surrogate func. $\left\{ \begin{array}{l} u(y, y) = f(y), \forall y \in \mathcal{X} \quad \text{touch at } \mathbf{x}^k \\ u(x, y) \geq f(x), \forall x, y \in \mathcal{X} \quad \text{upper bound} \\ u'(x, y; d)|_{x=y} = f'(y; d), \forall d \text{ with } y+d \in \mathcal{X} \quad \text{same slope at } \mathbf{x}^k \\ u(x, y) \text{ is continuous in } x \text{ and } y \end{array} \right.$

Applications:

(1) nonnegative least squares: $\underset{\mathbf{x} \geq 0}{\text{minimize}} \|A\mathbf{x} - \mathbf{b}\|_2^2 ; A \in \mathbb{R}_{++}^{m \times n}, \mathbf{b} \in \mathbb{R}_+^m$

$$\begin{aligned} f(\mathbf{x}) &= \frac{1}{2} \|A\mathbf{x} - \mathbf{b}\|_2^2 \\ u(\mathbf{x}, \mathbf{x}^k) &= f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^T (\mathbf{x} - \mathbf{x}^k) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^k)^T \phi(\mathbf{x}^k) (\mathbf{x} - \mathbf{x}^k) \end{aligned}$$

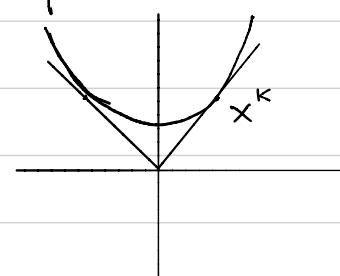
(Hessian of u) $\phi(\mathbf{x}^k) = \underbrace{\text{diag}(\dots)}_{\mathbf{x}^k} \simeq \mathbf{A}^T \mathbf{A}$ (Hessian of f)

$\underset{\mathbf{x} \geq 0}{\text{minimize}} u(\mathbf{x}, \mathbf{x}^k)$ is easier because each x_i is decoupled.

(2) reweighted LS for L_1 -norm minimization. \rightarrow LP solver

$$\underset{\mathbf{x}}{\text{minimize}} \|A\mathbf{x} - \mathbf{b}\|_1 \text{ is LP with no closed form solution}$$

$L_1 \rightarrow L_2$ if L_2 -norm \rightarrow LS (easy to solve)



(3) Sparse generalized eigenvalue problem:

$$\begin{aligned} \text{e.g. } \max_{\mathbf{x}} \quad & \mathbf{x}^T \mathbf{A} \mathbf{x} - \rho \|\mathbf{x}\|_0 \\ \text{s.t.} \quad & \mathbf{x}^T \mathbf{B} \mathbf{x} = 1 \end{aligned}$$

surrogate with quadratic not linear

Generalized eigenvalue problem

$$\begin{bmatrix} \max_{\mathbf{x}} & \mathbf{x}^T (\mathbf{A} - \rho \mathbf{B}) \mathbf{x} \\ \text{s.t.} & \mathbf{x}^T \mathbf{B} \mathbf{x} = 1 \end{bmatrix}$$

LII cont'd

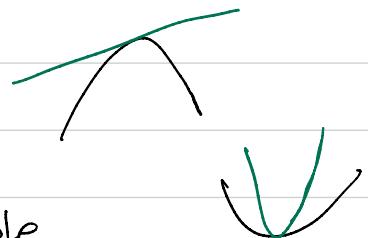
Construction of surrogate functions :

(1) Construction by convexity : Jensen inequality

(2) Construction by Taylor expansion:

① $k(x)$ is concave and differentiable:

$$k(x) \leq k(x^k) + \nabla k(x^k)^T(x - x^k)$$



② $k(x)$ is convex and twice differentiable.

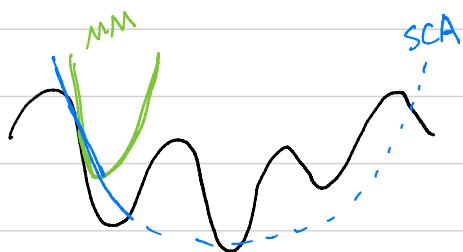
$$k(x) \leq k(x^k) + \nabla k(x^k)^T(x - x^k) + \frac{1}{2}(x - x^k)^T M(x - x^k)$$

$$\text{if } M - \nabla^2 k(x) \succeq 0 \quad \forall x$$

EM is a kind of MM

Connection to Successive Convex Approximation (SCA)

- △ surrogate function
 - { may not be convex in MM
 - { should be convex in SCA
 - { should be upper bound in MM
 - { doesn't need to be upper bound in SCA



△ MM can be easily extended to nonconvex X

△ SCA naturally has a parallel update

② Block MM : BCD + MM (with respect to that block)



when finding surrogate jointly is difficult

L12 Risk Parity Portfolio

只考虑 risk 分布, 不考虑收益

Problem formulation:

one number for total portfolio volatility: $\sigma(w) = \sqrt{w^T \Sigma w}$ → total
 volatility → decompose into a ↓
 sum of risk corresponding decompose: $\sigma(w) = \sum_{i=1}^N w_i \frac{\partial \sigma}{\partial w_i} = \sum_{i=1}^N w_i \cdot \frac{(\sum w)_i}{\sqrt{w^T \Sigma w}}$ ↑ each asset
 to each asset (Euler's theorem)
 marginal risk contribution (MRC) $MRC_i = \frac{\partial \sigma}{\partial w_i} = \frac{(\sum w)_i}{\sqrt{w^T \Sigma w}}$

risk contribution (RC) $RC_i = w_i \frac{\partial \sigma}{\partial w_i} = \frac{w_i (\sum w)_i}{\sqrt{w^T \Sigma w}}$

relative risk contribution (RRC) $RRC_i = \frac{RC_i}{\sigma(w)} = \frac{w_i (\sum w)_i}{w^T \Sigma w}$
 "normalize RC" → $\sum_{i=1}^N RRC_i = 1$

Risk parity portfolio: $RC_i = \sigma(w)/N$ or $RRC_i = 1/N$

general ↴
 Risk budgeting portfolio: $RRC_i = b_i \rightarrow w_i (\sum w)_i = b_i w^T \Sigma w$
 (enable different weights on risk, RPP follows $b_i = \frac{1}{N}$) ↓
 feasibility problem (nonlinear constraint)

Solutions:

① assume Σ is diagonal: $w_i^2 \sigma_i^2 = b_i \sum_{j=1}^N w_j^2 \sigma_j^2 \propto b_i$
 $\rightarrow w_i = \frac{\sqrt{b_i}}{\sum_{j=1}^N \sqrt{b_j}}$

② vanilla convex formulation:

let $x = w / \sqrt{w^T \Sigma w}$, then root follows: $x_i (\sum x)_i = b_i$ ↗
 Vector Form: $\sum x = b/x$

considering convex function: $f(x) = \frac{1}{2} x^T \Sigma x - b^T \log(x)$ ↑ same!!

minimization problem: $\min_{x \geq 0} f(x)$ by Newton's method / BCD

Find the root ⇒ Solve a minimization problem

L₁₂. cont'd.

③ General nonconvex formulation : because ② only when $\vec{w} = 1$
 $w \geq 0$

measure the differences \leftarrow add more constraints
between the terms $w_i(\Sigma w)_i$ (cannot achieve $w_i(\Sigma w)_i = \vec{w}^T \Sigma_i w$)

$$\begin{array}{ll}\text{minimize}_{w, \theta} & \sum_{i=1}^N (w_i(\Sigma w)_i - \theta)^2 \quad \dots \text{many other versions} \\ \text{s.t.} & \vec{w} = 1 \quad \text{nonconvex}\end{array}$$

L13 Sparsity

Optimization with Sparsity

$$\min_x f(x)$$

s.t. $x \in \mathcal{X}, \text{card}(x) \leq K$

NP hard. solve with convex relaxation

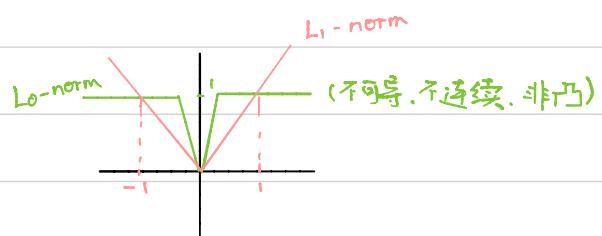
$\rightarrow \ell_0\text{-norm: } \|x\|_0$ (虽然不是 norm)

$\left\{ \begin{array}{l} \text{card}(x): \text{number of nonzero elements} \\ \text{supp}(x): \text{position with nonzero values} \end{array} \right.$



Algorithm for Sparsity Problems

- ① $\ell_1\text{-norm}$:
 1. apply $\ell_1\text{-norm}$
 2. set the very small elements to zero
 3. re-solve the new problem with fixed zero pattern



interpretation: ℓ_1 -norm is convex envelope of card on $[-1, 1]$
(largest convex function as an underestimator)

② Iterative Reweighted ℓ_1 -norm Heuristic

Suppose problem: $\min \|x\|_0$
s.t. $x \in \mathcal{X}$

Algorithm: set $w = 1$ repeat

$$\text{minimize}_x \|\text{Diag}(w)x\|_1, \text{ s.t. } x \in \mathcal{X}$$

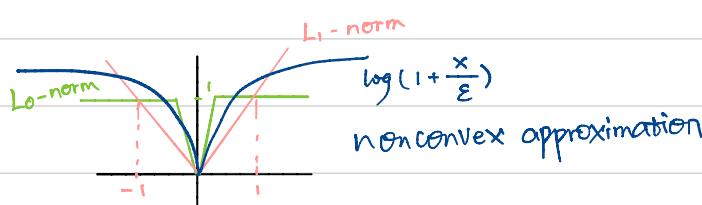
$$w_i = 1 / (\varepsilon + |x_i|)$$

until convergence to local point

for small $|x_i|$, w_i increases. $|x_i|$ even smaller

for larger $|x_i|$, w_i decreases. allow $|x_i|$ to be larger if necessary

Derivation: ① better approximation \rightarrow ② MM (Taylor. linearize)



assume $x \geq 0$

$$\sum_{i=1}^n \log(1 + x_i/\varepsilon) \approx \sum_{i=1}^n \log(1 + x_i^{(k)}/\varepsilon) + \sum_{i=1}^n \frac{x_i - x_i^{(k)}}{\varepsilon + x_i^{(k)}}$$

$$\begin{aligned} & \text{minimize}_x \sum_{i=1}^n w_i x_i \\ & \text{s.t. } x \in \mathcal{X} \end{aligned}$$

$w_i = \frac{1}{\varepsilon + x_i^{(k)}}$

L13. cont'd

Application of Sparsity:

① Compressed Sensing: $y = Ax$ with $A \in \mathbb{R}^{m \times n}$

- if $m \geq n$, and A is fullrank \Rightarrow unique or no solution
- if $m < n$, infinite solutions

$\hat{x} = A^+y$ is the solution of $\hat{x} = \operatorname{argmin}_{y=Ax} \|x\|_2$

Sparsity problem: $x^* = \operatorname{argmin}_{y=Ax} \|x\|_0$

② Estimation with Outliers

measurements: $y_i = a_i^T x + v_i$, v_i is Gaussian noise

add sparsity: $y_i = a_i^T x + v_i + w_i$

$$\underset{x, w}{\text{minimize}} \quad \|y - Ax - w\|_2$$

$$\text{s.t.} \quad \text{card}(w) \leq k$$

③ Feature Selection: choose subset of k regressors that best fit y

$$\underset{\beta}{\text{minimize}} \quad \|y - X^T \beta\|_2^2$$

$$\text{s.t.} \quad \text{card}(\beta) \leq k$$

④ LASSO: $\hat{\beta}_{\text{LASSO}} = \operatorname{argmin}_{\beta} \|y - X^T \beta\|_2^2 + \gamma \|\beta\|_1$

\rightarrow QP. when N is extremely large requires faster solution

L14 index tracking

Definition: Returns of an index in T days: $r^b = [r_1^b, \dots, r_T^b]^T$

Returns of N assets in T days: $X = [r_1, \dots, r_T]^T$

Company index \rightarrow Design a sparse portfolio w to track

$$\begin{cases} b > 0 \\ b^T 1 = 1 \\ Xb = r^b \end{cases}$$

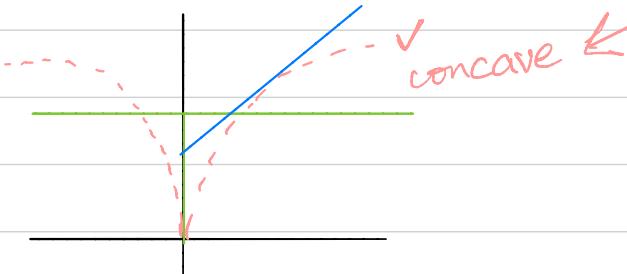
$$\begin{cases} w \geq 0 \\ w^T 1 = 1 \\ Xw \approx r^b \end{cases}$$

Sparse index tracking: via MM

$$\begin{aligned} & \underset{w}{\text{minimize}} \quad \frac{1}{T} \|Xw - r^b\|_2^2 + \lambda \|w\|_0 \\ & \text{s.t.} \quad w \in W \end{aligned}$$

l_0 -norm approximation

Step 1: Approximation of the l_0 -norm: $P_{p,r}(w) = \frac{\log(1 + |w|/p)}{\log(1 + \gamma/p)}$



$$\begin{aligned} & \underset{w}{\text{minimize}} \quad \frac{1}{T} \|Xw - r^b\|_2^2 + \lambda^T P_{p,u}(w) \\ & \text{s.t.} \quad w \in W \end{aligned}$$

Step 2: Linear approximation of $P_{p,r}(w)$

$$\begin{aligned} & \underset{w}{\text{minimize}} \quad \frac{1}{T} \|Xw - r^b\|_2^2 + \lambda d_{p,u}^{(k)^T} w \\ & \text{s.t.} \quad w \in W \end{aligned}$$

Step 3: Upper-bound the norm-square

$$\begin{aligned} & \underset{w}{\text{minimize}} \quad \frac{1}{T} w^T X^T X w + (\lambda d_{p,u}^{(k)} - \frac{2}{T} X^T r^b)^T w + \text{const.} \\ & \text{s.t.} \quad w \in W \end{aligned}$$

there is a closed-form update algorithm

Extension: if error term is not l_2 -norm, we can use MM to find a quadratic majorizer (L_2) first.

L15 - SVM

Classification:

prob: a set of input points $x_i \in \mathbb{R}^n \rightarrow$ binary labels $y_i \in \{-1, 1\}$

↓
separate with a linear model: $\hat{y} = \beta^\top x + \beta_0$

{ predict "Pos Class", if $\text{sign}(\hat{y}) = +1$
 { predict "Neg Class", if $\text{sign}(\hat{y}) = -1$

↓
optimization: $\begin{array}{ll} \underset{\beta_0, \beta, \{\hat{y}_i\}}{\text{minimize}} & \sum_{i=1}^N \mathbf{1}\{\text{sign}(\hat{y}_i) \neq y_i\} \rightarrow \text{nonconvex, nondifferentiable} \\ \text{s.t.} & \hat{y}_i = \beta^\top x_i + \beta_0, \forall i \end{array}$

Alternatives:

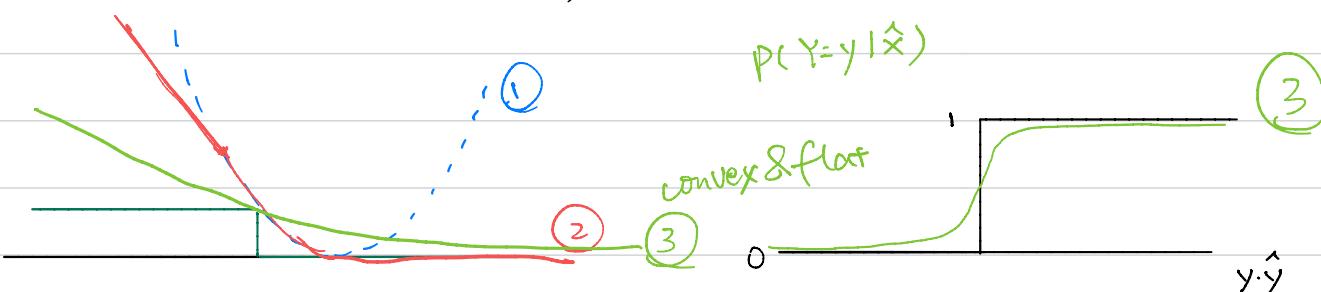
① Linear Regression

$$\left[\begin{array}{ll} \underset{\beta_0, \beta, \{\hat{y}_i\}}{\text{minimize}} & \sum_{i=1}^N (\hat{y}_i - y_i)^2 \Rightarrow \sum_{i=1}^N (1 - y_i \cdot \hat{y}_i)^2 \\ \text{s.t.} & \hat{y}_i = \beta^\top x_i + \beta_0, \forall i \end{array} \right]$$

② Regression with Huberized Loss

$$\phi_{\text{hub}}(x) = \begin{cases} |x|^2 & . |x| < M \text{ linear} \\ M(2|x|-M) & . |x| \geq M \text{ quadratic} \end{cases}$$

$$\left[\begin{array}{ll} \underset{\beta_0, \beta, \{\hat{y}_i\}}{\text{minimize}} & \sum_{i=1}^N \phi_{\text{hub-pos}}(1 - y_i \cdot \hat{y}_i)^2 \\ \text{s.t.} & \hat{y}_i = \beta^\top x_i + \beta_0, \forall i \end{array} \right]$$



③ Logistic Model model the probability of $y \in \{-1, 1\}$ as $P = \frac{1}{1+e^{-y \cdot \hat{y}}}$

$$\left[\begin{array}{ll} \underset{\beta_0, \beta, \{\hat{y}_i\}}{\text{minimize}} & \sum_{i=1}^N \log(1 + e^{-y_i \cdot \hat{y}_i}) \leftarrow \text{negative log-likelihood} \\ \text{s.t.} & \hat{y}_i = \beta^\top x_i + \beta_0, \forall i \end{array} \right]$$

L15 cont'd.

SVM:

① Linear Separable SVM:

Optimal Separating Hyperplane: ① Separate two classes

② Maximize the distance to closest point

$$\begin{aligned} \text{signed distance: } & \frac{1}{\|\beta\|} \beta^T (x - x_0) = \frac{1}{\|\beta\|} (\beta^T x + \beta_0) \\ & \text{平面法向量} \end{aligned}$$

$$\text{margin: } \frac{1}{\|\beta\|} y_i (\beta^T x_i + \beta_0)$$

$$\left[\begin{array}{l} \text{minimize}_{\beta_0, \beta} M \\ \text{s.t. } \frac{1}{\|\beta\|} y_i (\beta^T x_i + \beta_0) \geq M, \forall i \end{array} \right] \rightarrow \left[\begin{array}{l} \text{minimize}_{\beta_0, \beta} \|\beta\| \\ \text{s.t. } y_i (\beta^T x_i + \beta_0) \geq 1, \forall i \end{array} \right]$$

② Linear Nonseparable SVM

$$\left[\begin{array}{l} \text{minimize}_{\beta_0, \beta, \{\xi_i\}} \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N \xi_i \\ \text{s.t. } y_i (\beta^T x_i + \beta_0) \geq 1 - \xi_i, \xi_i \geq 0, \forall i \end{array} \right]$$

$\|\xi\|$

add sparsity

$$\left[\begin{array}{l} \text{minimize}_{\beta_0, \beta, \{\hat{\gamma}_i\}} \sum_{i=1}^N [1 - y_i \hat{\gamma}_i]^+ + \frac{\lambda}{2} \|\beta\|^2 \\ \text{s.t. } \hat{\gamma}_i = \beta^T x_i + \beta_0, \forall i \end{array} \right] \quad \lambda = \frac{1}{C}$$

③ Nonlinear SVM

Feature Transformation: $\left\{ \begin{array}{l} \mathbb{R}^n: \text{input space} \rightarrow H: \text{feature space} \\ x \rightarrow h(x) \end{array} \right\}$

kernel: $k(x_i, x_j) \triangleq h(x_i)^T h(x_j)$

↑ directly find k instead of h

1.6. Low Rank

Problem Formulation :
$$\begin{array}{l} \underset{x}{\text{minimize}} \quad \text{rank}(x) \\ \text{s.t.} \quad x \in C \leftarrow \text{convex set} \end{array}$$

non-convex

if X is diagonal . $\text{rank}(X) = \|\text{diag}(X)\|_1$.

rank minimization \rightarrow l_1 -norm minimization

Lemma : Let $X \in \mathbb{R}^{m \times n}$ be a given matrix . Then $\text{rank}(X) \leq r$ iff there exist $Y = Y^T \in \mathbb{R}^{m \times m}$ and $Z = Z^T \in \mathbb{R}^{n \times n}$ such that

$$\begin{bmatrix} Y & X \\ X^T & Z \end{bmatrix} \succeq 0, \quad \text{rank}(Y) + \text{rank}(Z) \leq 2r$$

$$\begin{bmatrix} \underset{x, Y, Z}{\text{minimize}} & \frac{1}{2} \text{rank}(\text{blkdiag}(Y, Z)) \\ \text{s.t.} & \begin{bmatrix} Y & X \\ X^T & Z \end{bmatrix} \succeq 0, \quad x \in C \end{bmatrix}$$

rank of rectangular
matrix \rightarrow square

Heuristics for Rank Minimization Problem NP-hard.

Solutions can be categorized into 2 groups :

① Approximate the rank function (non-convex) with surrogate functions

- Nuclear norm
- Log-det

② Solve a sequence of rank-constrained feasibility problems

- Matrix factorization
- Rank constraint via convex iteration

Lib cont'd.

Method 1: Nuclear Norm

$$\begin{bmatrix} \underset{x}{\text{minimize}} & \|X\|_* \\ \text{s.t.} & X \in C \\ (\|X\|_* = \sum_{i=1}^r \sigma_i) \end{bmatrix} \xrightarrow{\text{if } X = X^T \geq 0} \begin{bmatrix} \underset{x}{\text{minimize}} & \text{Tr}(X) \\ \text{s.t.} & X \in C \end{bmatrix}$$

Intuition: Nuclear norm = L_1 -norm of singular value vector
 → sparse singular value vector
 → low rank

$\|X\|_*$ is the convex envelope of $\text{rank}(X)$ on $\{X \mid \|X\|_2 \leq 1\}$

Lemma: For $X \in \mathbb{R}^{m \times n}$ and $t \in \mathbb{R}$, we have $\|X\|_* \leq t$

iff there exist matrices $Y \in \mathbb{R}^{m \times m}$ and $Z \in \mathbb{R}^{n \times n}$ such that

$$\begin{bmatrix} Y & X \\ X^T & Z \end{bmatrix} \geq 0, \quad \text{Tr}(Y) + \text{Tr}(Z) \leq 2t$$

△ $\begin{bmatrix} \underset{x,y,z}{\text{minimize}} & \frac{1}{2} \text{Tr}(Y+Z) \\ \text{s.t.} & \begin{bmatrix} Y & X \\ X^T & Z \end{bmatrix} \geq 0, \quad X \in C \end{bmatrix}$

minimize nuclear norm of rectangular matrix ⇒ write in bigger form (square)

singular values → eigenvalues → trace

Method 2: Log-det

Intuition: determinant = product of singular values

$$\begin{bmatrix} \underset{x}{\text{minimize}} & \log \det(X + \delta I) \\ \text{s.t.} & X \in C \end{bmatrix} \quad \left\{ \begin{array}{l} \log \det(X + \delta I) = \sum_i \log(\sigma_i(X + \delta I)) \\ \text{rank}(X) = \|\sigma(X)\|_1 \end{array} \right.$$

↓
 ⇒ log det can be seen as a surrogate function of $\text{card}(S)$. but it is concave

Iterative linearization & minimization
 (MM)

General for nonsquare X $\begin{bmatrix} \underset{x,y,z}{\text{minimize}} & \log \det(\text{blkdiag}(Y, Z) + \delta I) \\ \text{s.t.} & \begin{bmatrix} Y & X \\ X^T & Z \end{bmatrix} \geq 0, \quad X \in C \end{bmatrix}$

L16 cont'd.

Method 3: Matrix Factorization based Method.

< Previous 2 methods are not for big data (Large matrix) >

Low-rank Factorization:

$$X = F \cdot G$$

$M \times M$ M

$X = F \cdot G$ 本质上还是 rank, 对于 (X, F, G)

是 nonconvex ; 但对于 (X, F) 或 (X, G) 是 convex

当令一个变量 fixed.

[We assume X is convex, F (or G) is affine by X and both convex]

\approx BCD : (F, G)
 ↑ ↑
 optimize alternatively

Method 4: Rank Constraint via Convex Iteration

semidefinite rank-constrained
feasibility problem

find X
st. $X \in C, X \succeq 0$
 $\text{rank}(X) \leq r$ } convex prob.
→ nonconvex

$$= \sum_{i=r+1}^n \lambda_i(X) = \text{sum of smallest } n-r \hat{\lambda}_i(X)$$

$$\begin{bmatrix} \text{minimize}_{\substack{X \\ s.t. \\ X \succeq 0}} \text{Tr}(W^* X) \\ X \in C \end{bmatrix} \Rightarrow$$

$$\begin{bmatrix} \text{minimize}_{\substack{W \\ s.t. \\ ① \\ ② \\ \text{Tr}(W) = n-r}} \text{Tr}(WX^*) \\ ① \\ ② \\ 0 \leq W \leq I \\ ③ \end{bmatrix} \Rightarrow 0 \leq \lambda_i(w) \leq 1 \\ \sum \lambda_i(w) = n-r$$

Intuition:

align X with nullspace of w^*

Intuition: Considering X^* is semidefinite. $\lambda_1(x) \geq \lambda_2(x) \dots$

- (1) if only constraint ①: all $\lambda_i(w) = 0$
- (2) if constraint ①, ③: $\lambda_n(w) = n-r$ $\lambda_{1 \dots n-1}(w) = 0$
- (3) if constraint ①, ②, ③: $\underbrace{\lambda_{r+1}(w) = \lambda_{r+2}(w) = \dots = \lambda_n(w)}_{n-r} = 1, \lambda_{1 \dots r}(w) = 0$

Minimize the smallest $n-r \lambda_i(x) \Leftrightarrow$ Find the smallest $n-r \lambda_i(x)$



To smallest $n-r \lambda_i(x) = 0 \Rightarrow \text{rank}(x) \leq r$

L17 Robust Optimization

Robust Optimization:

Problem formulation: some problem contains parameters θ that are typically estimated in practice ($\hat{\theta}$) $\Rightarrow x^*(\hat{\theta}) \neq x^*(\theta)$

Several ways to make the problem robust to parameters errors:

- manageable {
- ① stochastic robust optimization (involving expectations)
 - ② worst-case robust optimization 太悲观 假设没有概率分布 (普遍情况)
 - ③ chance programming or chance robust optimization best but difficult (nonconvex)

Stochastic robust optimization: Expectations

Instead of using approximated function $f(x; \hat{\theta})$, use $E_\theta[f(x; \theta)]$.

e.g. ① model the estimated value as $\theta = \hat{\theta} + s$, $s \sim N(0, \sigma)$

② consider $f(x; \hat{\theta}) = (\hat{c}^\top x)^2$

$$\begin{aligned} \rightarrow E_\theta[f(x; \theta)] &= E_s[((\hat{c} + s)^\top x)^2] \\ &= E_s[x^\top \hat{c} \hat{c}^\top x + x^\top s s^\top x] \\ &= (\hat{c}^\top x)^2 + \underbrace{x^\top Q x}_{\text{regularizer}} \end{aligned}$$

Convexity: if $f(x; \theta)$ convex, then expectation (nonnegative sum) is convex

Worst-case robust optimization

Assume that true parameter lies in an uncertainty region centered around the estimated value: $\theta \in U$

关键: ① U 好解 ② U 大

sphere region: $U = \{\theta \mid \|\theta - \hat{\theta}\|_2 \leq \delta\}$

box region: $U = \{\theta \mid \|\theta - \hat{\theta}\|_\infty \leq \delta\}$

elliptical region: $U = \{\theta \mid (\theta - \hat{\theta})^\top S^{-1}(\theta - \hat{\theta}) \leq \delta^2\}$

L17 cont'd

- e.g. ① consider a sphere uncertainty region $\mathcal{U} = \{c \mid \|c - \hat{c}\|_2 \leq \delta\}$
 ② consider $f(x; \hat{\theta}) = (\hat{c}^T x)^2$ if the function is the objective
 to be minimized or it is a constraint of the form $f(x; \hat{\theta}) \leq 0$
 (maximized) (\Rightarrow)

$$\rightarrow \text{worst case: } \max_{c \in \mathcal{U}} |c^T x| = \max_{\|e\| \leq \delta} |(\hat{c} + e)^T x|$$

Cauchy-Schwarz inequality

$$|u^T v| \leq \|u\| \|v\|$$

$$\leq \max_{\|e\| \leq \delta} (|\hat{c}^T x| + |e^T x|)$$

$$\leq |\hat{c}^T x| + \delta \|x\| \leftarrow \text{the upper bound can achieve (maximum)}$$

equal iff $u = \lambda v$

要 minimize 的 objective. 考虑它 worst case 则是 $\min \max$

要是 constraint ≤ 0 , 考虑它的 worst case 是 $\max \leq 0$

- 一般要求 max 有闭式解, 否则问题会复杂.

→ - 系列不同 C 的 function

convexity: if C is convex set, f is convex. then pointwise maximum is convex

Application:

1) Robust Beamforming in Wireless Communications

sol 1: (no robust)

$$\max_w \text{SINR} = \frac{|w^H a|^2}{w^H R_{in} w} \quad (\text{signal power})$$

$$w^H R_{in} w \quad (\text{interference \& noise power})$$

a & R_{in} are parameters

required to be estimated

$$\begin{aligned} \text{nonconvex/nonconcave} &= \frac{w^H a a^H w}{w^H R_{in} w} = \frac{\tilde{w}^H R_{in}^{-1/2} a a^H R_{in}^{1/2} \tilde{w}}{\tilde{w}^H \tilde{w}} \quad \text{rank 1 matrix} \\ \text{but has hidden convexity} & \quad \tilde{w} = R_{in}^{1/2} w \quad \text{let it be 1 (normalize)} \end{aligned}$$

$$\Rightarrow \tilde{w} \propto R_{in}^{-1/2} a \leftarrow \text{唯一非零 eigenvector}$$

$$\max_x x^T A x$$

$$\text{s.t. } \|x\| = 1$$

x 找 A 最大 eigenvalue 对应的 eigenvector
 (最小 对应 min 问题)

L1.7 cont'd

sol2: $\min_w w^T R w$ $\Leftrightarrow \min_w w^T R w$ } QP. KKT
 (no robust) s.t. $|w^T a| = 1$ s.t. $w^T a = 1$

$\Delta w^H = w e^{j\phi}$ 多个维度 (对于每个 a 有对应的 ϕ , 不适用于 robust 讨论 - 系列 a)

sol2: add robust

$$A = \{c \mid c = \hat{a} + e, \|e\| \leq \varepsilon\}$$

problem: $\min_w w^T R w$ $\Leftrightarrow \min_w w^T R w$
 s.t. $|w^T a| = 1, \forall a \in A$ s.t. $|w^T a| \geq 1, \forall a \in A$
 - 放松 constraint worst-case $\rightarrow \min_w |w^T a| \geq |w^T \hat{a}| - \varepsilon \|w\|$
 ↓
 只有 \hat{a} , 此时可用 $w^H = w e^{j\phi}$ + trick

2) Robust Portfolio Optimization

(例如 LP/QP \rightarrow SOCP)
 加上 robustness, 形式上会变复杂

① Robust Global Maximum Return Portfolio Optimization

$$\begin{aligned} \max_w w^T \mu &\Rightarrow \max_w \min_{\mu \in U_\mu} w^T \mu && \Rightarrow \max_w w^T \hat{\mu} - \kappa \|S^{1/2} w\|_2 \\ \text{s.t. } 1^T w = 1 &\quad \text{s.t. } 1^T w = 1 && \text{s.t. } 1^T w = 1 \\ \text{LP} &\quad (\text{assume } U_\mu = \{\mu = \hat{\mu} + \kappa S^{1/2} u \mid \|u\|_2 \leq 1\}) && \text{soep(epigraph form)} \end{aligned}$$

② Robust Global Minimum Variance Portfolio Optimization

$$\begin{aligned} \min_w w^T \Sigma w &\Rightarrow \min_w \max_{\Sigma \in U_\Sigma} w^T \Sigma w && \Rightarrow \dots \Rightarrow \min_w \|\hat{x} w\|_2 + \delta_x \|w\|_2 \\ \text{s.t. } 1^T w = 1 &\quad \text{s.t. } 1^T w = 1 && \text{s.t. } 1^T w = 1 \\ \text{QP} &\quad \downarrow && \text{socp(epigraph form)} \end{aligned}$$

Matrix version $\hat{\Sigma} = \frac{1}{T} \mathbf{x}^T \mathbf{x}$, $x = \hat{x} + \Delta$, $U_x = \{x \mid \|x - \hat{x}\|_F \leq \delta_x\}$

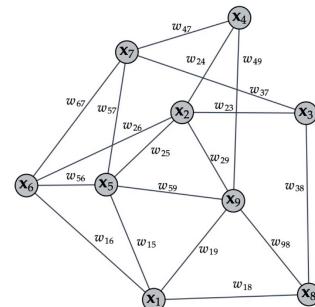
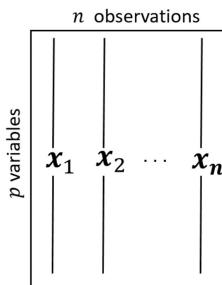
QP 变形: 增加平方 $\min_w \|\hat{x} w\|_2^2 + \delta_x \|w\|_2^2 = w^T (\hat{x}^T \hat{x} + \delta_x I) w$
 (试了才知道) s.t. $1^T w = 1$

L18 Graph Learning

Graph Basics:

data matrix $X =$

$$[x_1, x_2, \dots, x_n] \in \mathbb{R}^{p \times n}$$



(有多种形式，一般对称且
对角为0，非PSD)

weight adjacency
matrix W

Degree matrix D contains the degrees d_1, \dots, d_p along the diagonal

$$D = \text{Diag}(W^T W) \text{ i.e. } d_i = \sum_{j=1}^p W_{ij}$$

Laplacian of a graph: $L = D - W$ {

- ① symmetric & PSD
- ② # zero eigenvalues = # connected components
- ③ measure the smoothness of that vector

$x^T L x = \frac{1}{2} \sum_{ij} W_{ij} (x_i - x_j)^2 \leftarrow$ weighted with the graph weights

Learning Graph from Data

learn W or L both OK.

Method 1: Similarity function based (Learn W)

Intuition: decide the weight between two nodes based on their similarity

Various measures of similarity:

- Thresholded Euclidean distance graph: $\|x_i - x_j\|^2 \leq \gamma$, $W_{ij}=0/1$
- Gaussian graph: $W_{ij} = \exp(-\frac{\|x_i - x_j\|^2}{2\sigma^2})$
- k-NN graph
- Feature correlation graph: $W_{ij} = x_i^T x_j$
- Self-tuned Gaussian graph: $W_{ij} = \exp(-\frac{\|x_i - x_j\|^2}{2\sigma_i \sigma_j})$ distance(x_i , x_j 's kNN)

L18 cont'd

Method 2: Smooth signal based

Intuition: find a L that minimizes $\sum_{i=1}^n x_i^T L x_i = \text{Tr}(X^T L X)$

obj: minimize $\text{Tr}(X^T L X) + \gamma h(L)$ regularization (all convex)
 $\|L\|_1$: sparsity, $\|L\|_F^2$: energy
 $-\log \det(L)$: volume ($\because \det$ is volume)

convex

QP

$$\left[\begin{array}{l} \underset{L}{\text{minimize}} \quad \text{Tr}(X^T L X) + \gamma \|L\|_{F,\text{off}}^2 \quad (\text{off-diagonal}) \\ \text{s.t.} \quad \text{diag}(L) = 1 \quad (\text{可选, 控制 degree}) \\ \quad L \mathbf{1} = \mathbf{0}, \quad L_{ij} = L_{ji} \leq 0 \quad (\text{Laplacian 性质}) \end{array} \right] \xrightarrow{\text{all linear}} \text{写成 } L$$

Linear + QP

$$\left[\begin{array}{l} \underset{L}{\text{minimize}} \quad \text{Tr}(W Z) + \gamma \|W\|_F^2 \\ \text{s.t.} \quad W \mathbf{1} = \mathbf{1} \\ \quad \text{diag}(W) = \mathbf{0}, \quad W = W^T \geq \mathbf{0} \end{array} \right] \xrightarrow{\text{写成 } W} Z_{ij} = \|x_i - x_j\|^2$$

Solution: ignore $W = W^T$, the problem decomposes into each row;

at the end symmetrize by $W \leftarrow \frac{1}{2}(W + W^T)$.

$$\left[\begin{array}{l} \underset{L}{\text{minimize}} \quad z_i^T w_i + \gamma \|w_i\|^2 \\ \text{s.t.} \quad w_i^T \mathbf{1} = 1, \quad w_{ii} = 0, \quad w_i \geq \mathbf{0} \end{array} \right] \xrightarrow{\text{KKT waterfilling}}$$

or

relax $w_i^T \mathbf{1} = 1$

$$\left[\begin{array}{l} \underset{w}{\text{minimize}} \quad \frac{1}{2} z_i^T w_i - \alpha \log(w_i) + \frac{\beta}{2} \|w_i\|^2 \\ \text{s.t.} \quad w_{ii} = 0, \quad w_i \geq \mathbf{0} \end{array} \right]$$

L18 cont'd.

Extension: Learning a k -connected graph

Goal: Given an initial noisy adjacency matrix W_0 , infer a k -connected graph
 = k smallest eigenvalues are zero

Define the Laplacian operator: $L(W) = \text{Diag}(W\mathbf{1}) - W$

problem formulation (nonconvex)

$$\begin{aligned} & \underset{W}{\text{minimize}} \quad \|W - W_0\|_F^2 \\ & \text{s.t.} \quad W\mathbf{1} = \mathbf{1}, \text{diag}(W) = \mathbf{0}, W = W^\top \geq \mathbf{0} \\ & \quad \underbrace{\text{rank}(L(W)) = p-k}_{\downarrow} \end{aligned}$$

relax constraint

$$\lambda_1(L(W)) = \dots = \lambda_k(L(W)) = 0$$

$$\begin{aligned} & \underset{W}{\text{minimize}} \quad \|W - W_0\|_F^2 + \beta \sum_{i=1}^k \lambda_i(L(W)) \\ & \text{s.t.} \quad W\mathbf{1} = \mathbf{1}, \text{diag}(W) = \mathbf{0}, W = W^\top \geq \mathbf{0} \end{aligned}$$

knowing that

$$\sum_{i=1}^k \lambda_i(x) = \begin{cases} \min_{F \in \mathbb{R}^{p \times k}} \text{Tr}(F^\top x F) & (\text{nonconvex for } (F, x)) \\ \text{s.t.} \quad F^\top F = I & (\text{nonconvex}) \end{cases}$$

BCD

- 1) fix F : convex
- 2) fix W : nonconvex

L17 beamforming

F 取 $L(W)$ 最小 k 个 λ_i 对应 eigenvector

Method 3: i.i.d model based (MLE based)

Details can be found in handout.

L19~20 SDP relaxation

BPSK Signal Detection Problem

Linear observation: $y = \begin{matrix} Hs \\ \downarrow \end{matrix} + w \rightarrow$ gaussian noise
 received signal $\xrightarrow{\text{BPSK}}$

problem formulation: $\begin{bmatrix} \underset{s}{\text{minimize}} & \|y - Hs\|_2^2 \\ \text{s.t.} & s \in \{\pm 1\}^n \end{bmatrix} \xrightarrow{\text{convex}} \text{nonconvex}$

Convex Relations (Relax尽可能是 convex)

Type 1: Box relaxation $-1 \leq s_i \leq 1$

Δ Type 2: SDP relaxation

homogeneous form: $\begin{bmatrix} \underset{x}{\text{minimize}} & x^T L x \\ \text{s.t.} & x \in \{\pm 1\}^{n+1} \end{bmatrix}$

equivalent

$$L = \begin{bmatrix} H^T H & -H^T y \\ -y^T H & y^T y \end{bmatrix}$$

$$x = \begin{bmatrix} s \\ 1 \end{bmatrix}$$

$x_{n+1} = 1$ can be removed

equivalent

define $X = xx^T$.

$$X_{ii} = x_i^2 = 1$$

$$\underset{x.x}{\text{minimize}} \quad \text{Tr}(LX)$$

$$\text{s.t.} \quad \text{diag}(X) = \mathbf{1}_{n+1}$$

$$X = xx^T$$

nonconvex for (X, x)

relaxation

$X = xx^T$ equals to $X \succeq 0$

and $\text{rank}(X) = 1 \leftarrow$ nonconvex;

remove rank constraint

$$\underset{x.x}{\text{minimize}} \quad \text{Tr}(LX)$$

$$\text{s.t.} \quad \text{diag}(X) = \mathbf{1}_{n+1}$$

$$X \succeq 0$$

SDP

Reconstruct Binary Solutions

Goal: Get X , reconstruct $x(s)$. If $\text{rank}(X)=1$ then it means there is no relaxation (solved), which is not the common case.

Method 1: Simple Quantization $X \rightarrow s \rightarrow \text{sign}(s)$

Method 2: Eigenvalue Decomposition take the eigenvector with the largest λ as s

Method 3: Randomization $X = xx^T$; use X as the covariance matrix to generate random points, and keep the \hat{x} with minimum objective value

L 19 ~ 20 cont'd

Application: Multiuser Downlink Beamforming

problem formulation

$$\begin{aligned} & \underset{\{w_l\}}{\text{minimize}} \quad \sum_{l=1}^L \|w_l\|^2 \quad \rightarrow \text{overall transmitted power} \\ & \text{s.t.} \quad \frac{w_m^H R_{mm} w_m}{\sum_{l \neq m} w_l^H R_{ml} w_l + \sigma_m^2} \geq p_m \quad \rightarrow \text{SINR constraint} \end{aligned}$$

Method 1: SOCP formulation

trick: $w_m^H = w_m \cdot e^{j\phi}$ to make w_m^H real

SOCP

$$\begin{aligned} & \underset{\{w_l\}}{\text{minimize}} \quad \sum_{l=1}^L \|w_l\|^2 \\ & \text{s.t.} \quad w_m^H h_m \geq p_m \|\tilde{R}_m^{1/2} \tilde{w}_m\| \\ & \quad w_m^H h_m \geq 0 \end{aligned}$$

assume $R_{mm} = h_m h_m^H$
add $w_m^H h_m \geq 0$

Method 2: SDP relaxation

define the rank-1 matrix: $X_L = w_L w_L^H$

notation: $w_L^H A w_L = \text{Tr}(w_L^H A w_L) = \text{Tr}(A w_L w_L^H) = \text{Tr}(A X_L) \triangleq A \cdot X_L$

SDP

$$\begin{aligned} & \underset{X_1, \dots, X_L}{\text{minimize}} \quad \sum_{l=1}^L I \cdot X_l \\ & \text{s.t.} \quad R_{mm} \cdot X_m - p_m \sum_{l \neq m} R_{ml} \cdot X_l \geq p_m \sigma_m^2 \quad \forall m \\ & \quad X_l \succeq 0 \end{aligned}$$

$\text{rank}(X_L) = 1$ ignore it \rightarrow convex problem

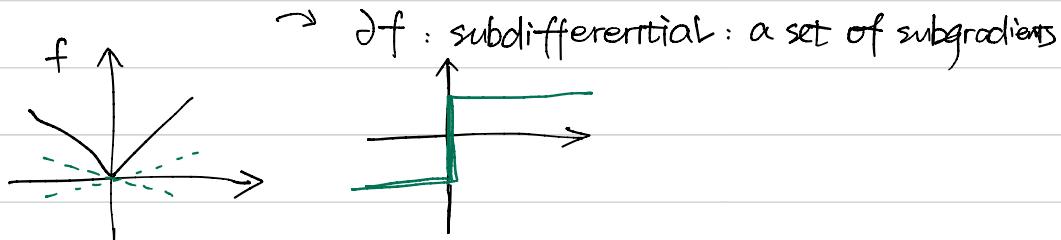
transfer nonconvexity in constraint into rank

L21 Primal Dual Decomposition

Subgradient definition: g is a subgradient of f (not necessarily convex) at x if

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

support hyperplane.



optimal conditions Unconstrained Case: $\nabla f(x^*) = 0 \rightarrow 0 \in \partial f(x^*)$

Constrained Case: $\min f_0(x)$ (KKT cond. 3)
 s.t. $f_i(x) \leq 0 \rightarrow 0 \in \partial f_0(x^*) + \sum_{i=1}^m \lambda_i^* \partial f_i(x^*)$

Subgradient Method minimize a nondifferentiable convex function f

similar to gradient descent Step 1: $x^{(k+1)} = x^{(k)} - \alpha_k g^{(k)}$ subgradient
 Step 2: $f_{\text{best}} = \min_{i=1 \dots k} f(x^{(i)})$

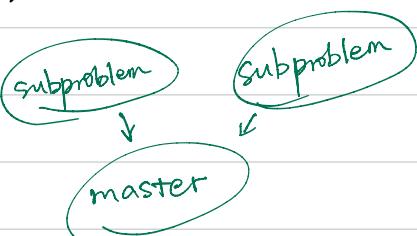
not guarantee to decrease each step (keep the best)
 guarantee to converge to a neighbourhood.

Projected Subgradient Method for constrained optimization $\min f(x) \Rightarrow x^{(k+1)} = [x^{(k)} - \alpha_k g^{(k)}]_X$
 s.t. $x \in X$ step + project

Primal Decomposition (for coupling variable) problem: $\begin{bmatrix} \min_{x,y} f_1(x_1, y) + f_2(x_2, y) \\ \text{s.t. } x_1 \in \mathcal{X}_1, x_2 \in \mathcal{X}_2, y \in \mathcal{Y} \end{bmatrix}$

Step 1: subproblem: fixed y $\begin{cases} f_1^*(y) = \min_{x_1 \in \mathcal{X}_1} f_1(x_1, y) \\ f_2^*(y) = \min_{x_2 \in \mathcal{X}_2} f_2(x_2, y) \end{cases}$

Step 2: master (primal) problem: $\min_{y \in \mathcal{Y}} f_1^*(y) + f_2^*(y)$



L21 cont'd.

To solve the master problem:

1. bisection (if y is scalar)
2. gradient or Newton method (if f_i^* differentiable)
3. subgradient. cutting-plane. or ellipsoid method "key"
 - ↳ projected subgradient method: $y^{(k+1)} = [y^{(k)} - c_k(s_1(k) + s_2(k))]$

Calculate subgradients (free):

Lemma 1) Let $f^*(y)$ be the optimal value of the convex problem

$$\underset{x}{\text{minimize}} \quad f_0(x)$$

$$\text{s.t.} \quad h_i(x) \leq y_i \quad \text{A subgradient of } f^*(y) \text{ is } -\nabla^*(y)$$

Dual Decomposition

(coupling constraints)

problem:
$$\begin{bmatrix} \min_x \quad f_1(x_1) + f_2(x_2) \\ \text{s.t.} \quad x_1 \in X_1, x_2 \in X_2 \\ h_1(x_1) + h_2(x_2) \leq h_0 \end{bmatrix}$$

Partial Lagrangian $L(x, \lambda) = f_1(x_1) + f_2(x_2) + \lambda^\top (h_1(x_1) + h_2(x_2) - h_0)$



Max the dual function $g(\lambda) = \inf_{x \in X} L(x, \lambda)$

$$= \inf_{x_1 \in X_1} \{f_1(x_1) + \lambda^\top h_1(x_1)\}$$

$$+ \inf_{x_2 \in X_2} \{f_2(x_2) + \lambda^\top h_2(x_2)\} - \lambda^\top h_0$$

Step 1: subproblem: fixed y

$$\begin{cases} g_1(\lambda) = \underset{x_1 \in X_1}{\text{minimize}} \quad f_1(x_1) + \lambda^\top h_1(x_1) \\ g_2(\lambda) = \underset{x_2 \in X_2}{\text{minimize}} \quad f_2(x_2) + \lambda^\top h_2(x_2) \end{cases}$$

Step 2: master dual problem: $\underset{\lambda \geq 0}{\text{maximize}} \quad g_1(\lambda) + g_2(\lambda) - \lambda^\top h_0$

→ projected subgradient method. $\lambda^{(k+1)} = [\lambda^{(k)} + c_k(s_1(k) + s_2(k) - h_0)]^+$

Calculate subgradients (free):

Lemma: Let $g(\lambda)$ be the dual function corresponding to the problem

$$\underset{x}{\text{minimize}} \quad f_0(x)$$

$$\text{s.t.} \quad h_i(x) \leq 0, i = 1, \dots, m$$

A subgradient of $g(\lambda)$ is $h(x^*(\lambda))$

L21 cont'd.

benefit: no need for projection

(ADMM → dual decomposition → subgradient → non-differentiable) solve

ADMM (Alternating Direction Method of Multipliers)

Background:

problem: $\begin{bmatrix} \underset{x}{\text{minimize}} & f(x) \\ \text{s.t.} & Ax = b \end{bmatrix}$

Lagrangian: $L(x, y) = f(x) + y^T(Ax - b)$
 dual function: $g(y) = \inf_x L(x, y)$
 dual problem: $\underset{y}{\text{maximize}} g(y)$
 recover: $x^* = \underset{x}{\text{argmin}} L(x, y^*)$

Dual Ascent Idea: $x^{k+1} := \underset{x}{\text{argmin}} L(x, y^k)$

$$y^{k+1} := y^k + \rho^k \nabla g(y^k) = y^k + \rho^k (Ax^{k+1} - b)$$

ADMM:

problem: $\begin{bmatrix} \underset{x, z}{\text{minimize}} & f(x) + g(z) \\ \text{s.t.} & Ax + Bz = C \end{bmatrix}$

Augmented Lagrangian:

$$L_p(x, z, y) = f(x) + g(z) + \langle y, Ax + Bz - C \rangle + \underbrace{\frac{\rho}{2} \|Ax + Bz - C\|_F^2}$$

ADMM:

$$x^{k+1} := \underset{x}{\text{argmin}} L_p(x, z^k, y^k)$$

$$z^{k+1} := \underset{z}{\text{argmin}} L_p(x^{k+1}, z, y^k)$$

$$y^{k+1} := y^k + \rho(Ax^{k+1} + Bz^{k+1} - C)$$

add regularization to
dual decomposition
(faster convergence)
控制往可行集靠

no need for projection

Convergence: iterates approach feasibility: $Ax^k + Bz^k - C \rightarrow 0$

(assume convex f,g) objective approaches optimal value: $f(x^k) + g(z^k) \rightarrow P^*$

OK for non-differentiable speed: slow → high accuracy;

↓ efficient → moderate accuracy

but can extend to nonconvex

L21 cont'd.

Example 1 : Robust PCA

$$\begin{bmatrix} \underset{\text{L}, \text{S}}{\text{minimize}} & \| \text{L} \|_* + \lambda \| \text{S} \|_1 \\ \text{s.t.} & \text{L} + \text{S} = M \end{bmatrix} \xrightarrow{\text{convex but not differentiable}}$$

Example 2 : Graphical Lasso

Precision matrix estimation from Gaussian samples

$$\begin{bmatrix} \underset{\Theta}{\text{minimize}} & -\log \det \Theta + \underbrace{\langle \Theta, S \rangle}_{\text{neg. log likelihood}} + \lambda \| \Theta \|_1 \\ \text{s.t.} & \Theta > 0 \end{bmatrix}$$

引入slack variable
方便构造ADMM

$$\begin{bmatrix} \underset{\Theta}{\text{minimize}} & -\log \det \Theta + \langle \Theta, S \rangle + \lambda \| \Psi \|_1 \\ \text{s.t.} & \Theta > 0, \quad \Theta = \Psi \end{bmatrix}$$

Summary : Algorithms

Unconstrained Optimization:

- (1) Descent Method ; Gradient Descent Method
- (2) Newton's Method.

Constrained Optimization:

- (1) Equality Constrained Optimization

general { (2) Gradient Projection Method.

- (3) Interior - Point Method. (IPM)

block { (4) Block Coordinate Descent (BCD) sequential

- (5) Jacobi Algorithm parallel

difficult problem { (6) Majorization - Minimization (MM)

- (7) Successive Convex Approximation (SCA) parallel

block { (8) Block Majorization - Minimization (BMM) sequential

- (9) Primal Decomposition

- (10) Dual Decomposition

coupling { (11) Alternating Direction Method of Multipliers (ADMM)

Extension: PCA and eigenvalue

What is PCA?

PCA tries to find the direction such that the projection of the data on it has the highest possible variance.

Sample covariance matrix: $S = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top = \mathbf{X}^\top \mathbf{X} / n$

which is a PSD matrix with an eigenvalue decomposition: $S = Q \Lambda Q^\top$
 $(\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0)$

PCA seeks to solve a sequence of optimization problem,

the first is : $\underset{\mathbf{u} \in \mathbb{R}^P}{\text{maximize}} \frac{\mathbf{u}^\top S \mathbf{u}}{\mathbf{u}^\top \mathbf{u}}$

常见形式

(见于 L7)

$$\underset{\substack{\mathbf{u} \\ \text{s.t.}}}{\text{maximize}} \quad \mathbf{u}^\top S \mathbf{u}, \quad \mathbf{u}^\top \mathbf{u} = 1$$

$$\begin{aligned} \mathbf{u}^\top S \mathbf{u} &= \mathbf{u}^\top Q \Lambda Q^\top \mathbf{u} = \mathbf{w}^\top \Lambda \mathbf{w} = \sum_{i=1}^P \lambda_i w_i^2 \\ \|\mathbf{w}\|_2 &= \|Q^\top \mathbf{u}\|_2 = \|\mathbf{u}\|_2 = 1 \end{aligned} \rightarrow \underset{\substack{\mathbf{w} \\ \text{s.t.}}}{\text{maximize}} \quad \sum_{i=1}^P \lambda_i w_i^2 \rightarrow \mathbf{w} = (1, 0, \dots, 0)^\top = \mathbf{e}_1$$

rotation

\therefore the first principal component is the eigenvector wr.t. the largest eigenvalue.

the remaining components

are found by solving the sequence

$$\underset{\substack{\mathbf{u} \\ \text{s.t.} \\ \mathbf{u}^\top \mathbf{u}_j = 0, \forall 1 \leq j \leq i}}{\text{maximize}} \quad \mathbf{u}^\top S \mathbf{u}$$

$$\rightarrow \mathbf{u}_i = \mathbf{q}_i$$

note: rank and # eigenvalue if A is an $n \times n$ matrix, $\text{rank}(A) + \text{null}(A) = n$

$\therefore \text{null}(A)$ is the dimension of the kernel of the matrix (V)

$$Av = 0 = 0V$$

The kernel of A is precisely the eigenspace corresponding to eigenvalue 0

$\therefore \text{rank}(A) = n - \text{eigenspace of } A \text{ corresponding to eigenvalue 0}$