1. 概念

<span style="color:red">正定矩阵 P 对称，且 x^TPx>0 特征值大于 0</span>

- **affine combination**

仿射和凸组合

$$\theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_k x_k \quad \text{with} \quad \theta_1 + \theta_2 + \cdots + \theta_k = 1$$

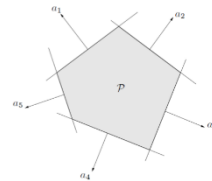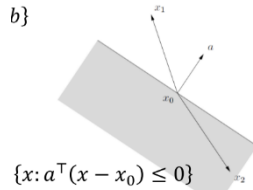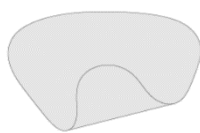$\theta = 0 \quad x_2$
$\theta = -0.2$

- **convex combination**

$$\theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_k x_k \quad \text{with} \quad \theta_1 + \theta_2 + \cdots + \theta_k = 1 \quad \text{and } \theta_i \geq 0, \forall i$$

凸集：集合内的点的凸组合也属于集合　　内点（领域也在集合内）**空集和全集**<span style="color:red">即开又闭</span>
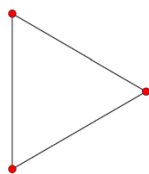
凸包（最小凸集）　　　　**超平面，半空间（a 是法向量）**　　多面体由超平面和半空间围成

$b\}$

$\{x: a^\top(x - x_0) \leq 0\}$

$\{x: a_i^\top x \leq b_i, c_j^\top x = d_j\}$

仿射独立（一个点无法由其他点的仿射变换表示，**2D 最多 <span style="color:red">3</span> 个**）　　　　锥：<span style="color:red">过原点</span>的射线

单纯形　（一系列<span style="color:red">仿射独立</span>的点的凸包）　　　　　　　　凸锥:锥变换在集合内

两点一线段，三点一图形

$$\theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_k x_k \quad \text{with} \quad \theta_i \geq 0$$

椭球　**P 正定**

$$\{x: (x - x_c)^\top P^{-1}(x - x_c) \leq 1\}$$

Proper Cone
Generalized Inequality
- $K$ is a proper cone:
  - closed
  - solid: non-empty interior
  - <mark>pointed: contains no line</mark>

perspective function $P: \mathbf{R}^{n+1} \to \mathbf{R}^n$

$$P(x, t) = x/t$$

<span style="color:red">保凸运算</span>　　1.交集 2.仿射变换（凸集的和，放大）3.透视降维

4.线性分式函数　仿射+透视　　　　　　　　透视（KL 散度）　g(x,t)

$$f(x) = \frac{Ax + b}{c^\top x + d}$$

5. 非负加权和　　$\mathbf{dom}\, f = \{x | c^\top x + d > 0\}$

- $f(x) = -\log x$ is convex and its perspective function is
$$g(x, t) = tf(x/t) = -t\log\frac{x}{t} = t\log t - t\log x$$
- Kullback-Leibler divergence
$$D_{kl}(u, v) = \sum_{i=1}^{n}(u_i \log(u_i/v_i) - u_i + v_i)$$

- if $f(x, y)$ is convex w.r.t. $(x, y)$ and $C$ is a convex set, then
  the minimization $g(x) = \inf_{y \in C} f(x, y)$ is convex
- distance to a convex set $S$
$$\text{dist}(x, S) \triangleq g(x) = \inf_{y \in S} \|x - y\|$$
difference to
$$f(x) = \sup_{y \in C}\|x - y\|$$

$$f = w_1 f_1 + \cdots + w_m f_m, \ w_m \geq 0$$

7 最小值和舒尔补

- Shur complement
- $f(x, y) = x^\top Ax + 2x^\top By + y^\top Cy$ with $\begin{bmatrix} A & B \\ B^\top & C \end{bmatrix} \geq 0, C > 0$

6.逐点上确界　$f(x) = \max\{f_1(x), \ldots, f_m(x)\}$ is convex

分离超平面（必须凸集，不可逆，<span style="color:red">严格的话不可取等号</span>）　　　　　　支撑超平面（必须凸集）
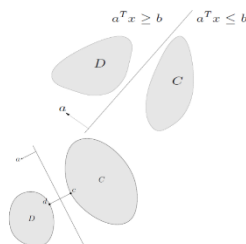
**Separating Hyperplane Theorem**

- if $C$ and $D$ are disjoint convex sets, i.e., $C \cap D = \emptyset$, then there exists at least one
  **separating hyperplane** such that
  $$a^\top x \leq b, \forall x \in C$$
  $$a^\top x \geq b, \forall x \in D$$

  $a^\top x \geq b \quad a^\top x \leq b$
  $D \quad C$

- Proof sketch for the case, the distance
  between $C$ and $D$ are positive
  - find the closest points
  - define the hyperplane by the middle
    point and the orthogonal direction

**Supporting Hyperplane Theorem**

For a set $C \subseteq \mathbf{R}^n$ and a point in its boundary $x_0$,
if $a \neq 0$ satisfies $a^\top x \leq a^\top x_0, \forall x \in C$, then we
call the corresponding hyperplane a **supporting
hyperplane** to $C$ at $x_0$

- **Supporting Hyperplane Theorem**: for a convex set, there exists at least one
  supporting hyperplane at every boundary point. <span style="color:red">how about non-convex sets?</span>
  <span style="color:red">can we have more?</span>
- a convex set could be represented by (maybe infinite) linear inequalities
- (partial converse): if a set is close, has nonempty interior, and has a supporting
  hyperplane at every point in its boundary, then it is convex.

vertex

**凸集可以由线性不等式的解集表示（通过分离超平面定理）**

2. 凸函数（定义域，可不光滑，是否严格） 一阶和二阶

**Definition**

- $f: R^n \to R$ is convex, if **dom** $f$ is a convex set and
  $f(\theta x + (1-\theta)y) \le \theta f(x) + (1-\theta)f(y)$ $\forall x, y \in \mathbf{dom}\, f, \theta \in [0,1]$

定义域必须是凸集

- "line is above the curve":

- **strictly convex** if strict inequality holds for $x \ne y$ and $\theta \in (0,1)$

- **First-order condition**

$$f(y) \ge f(x) + \nabla f(x)^\top (y - x) \quad \forall x, y \in \mathbf{dom}\, f$$

"tangent plane is below the surface"

$\nabla^2 f(x) \ge 0, \forall x \in \mathbf{dom}\, f \iff f$ is convex

$\nabla^2 f(x) > 0, \forall x \in \mathbf{dom}\, f \implies f$ is strictly convex

水平集（**针对定义域的**，保凸）

**sublevel set** of $f: R^n \to R$ $C_\alpha = \{x \in \mathbf{dom}\, f : f(x) \le \alpha\}$

sublevel sets of convex functions are convex

上镜图（升维，保凸）

- **epigraph** of $f: R^n \to R$

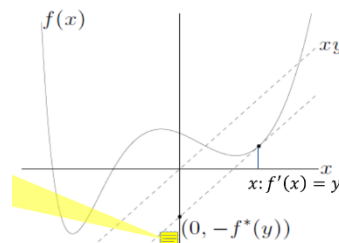  $\mathbf{epi}\, f = \{(x,t) \in R^{n+1} \, x \in \mathbf{dom}\, f : f(x) \le t\}$

- $f$ is convex iff **epi** $f$ is a convex set

**共轭函数**
**切线和函数差值最大值，y 为斜率**
在 f(x)'=y 的点找 永远是凸函数

for a given $y$, $f^*(y)$ is the largest gap between the affine function $xy$ and the function $f(x)$

when $f$ is differentiable, the optimality condition tells

$$\nabla(y^\top x - f(x)) = y - f'(x) = 0$$

3. 凸优化问题 在凸集上找凸函数的最小值
可行解属于可行集（凸函数和约束的定义域交集）

$$\begin{aligned} \min \quad & f_0(x) \\ \text{s.t.} \quad & f_i(x) \le 0, i = 1, \dots, m \\ & h_i(x) = 0, i = 1, \dots, p \end{aligned}$$

全局最优解 f0(x)最小
局部最优解：距离 x 范围 d 内，f0(x+d)>=f0(x)
线性规划（relaxation 法）：例如整数约束变成非整数
线性分式（规划换变量法）
二次规划（梯度法）向量优化
反正拉格朗日对偶最好用

**Linear-fractional Programming**

$d = f = 0$

$$\begin{aligned} \min \quad & \frac{c^\top x + d}{e^\top x + f} \\ \text{s.t.} \quad & Gx \le h \\ & Ax = b \\ & e^\top x + f > 0 \end{aligned}$$
$y = \dfrac{x}{e^\top x + f}$
$z = \dfrac{1}{e^\top x + f}$

- **Linear-fractional Programming** is equivalent to the following LP

$$\begin{aligned} \min \quad & c^\top y + dz \\ \text{s.t.} \quad & Gy \le hz \\ & Ay = bz \\ & e^\top y + fz = 1 \\ & z \ge 0 \end{aligned}$$
$x, y \in R^n$
$z \in R$
$x = y/z$

4. 常用模型
监督学习
**线性拟合**用最小二乘法 缺点：偏离值影响太大
 解决方法：**huber loss**（连续和光滑问题）

**分位拟合：拟合曲线上下点成固定比例**

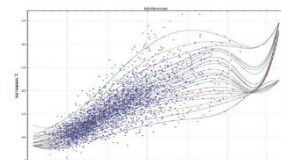**正则化项： 额外 min 所有参数**
**Lasso: l1norm，使系数稀疏**
Ridge: l2 norm

## Huber loss

$$L_{\text{huber}}(r) = \begin{cases} r^2 & \text{if } |r| < c \\ |r| & \text{if } |r| \ge c \end{cases}$$

result is the median value

regression with asymmetric loss

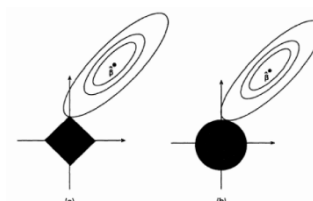$$L_\tau(r) = \begin{cases} -(1-\tau)r & \text{if } r < 0 \\ \tau r & \text{if } r \ge 0 \end{cases}$$
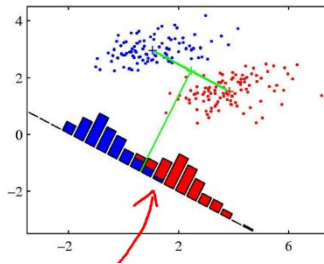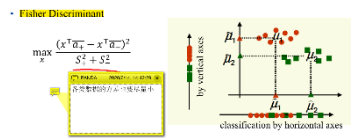
**LDA: 在直线上投影均值最大**
方差越小越好

**Fisher Discriminant**

$$\max_x \frac{(x^\top \overline{a_+} - x^\top \overline{a_-})^2}{S_+^2 + S_-^2}$$

**CCA 找到两条线,使得数据上面的投影相关性最大**
**数据 A,B 描述同个 object**

$$\max_{x,y} \frac{x^\top \text{cov}(A,B)y}{\sqrt{x^\top \text{cov}(A,A)x}\sqrt{y^\top \text{cov}(A,A)y}}$$

**SVM 最大边距 硬边距**

$$\min_{x,z} \|x\|_2^2$$
$$\text{s.t.} \quad b_i(x^\top b_i + z) \geq 1$$

**软边距**

$$\min_{x,z} \frac{1}{2}\|x\|_2^2 + C\sum s_i$$
$$\text{s.t.} \quad b_i(x^\top b_i + z) \geq 1 - s_i$$
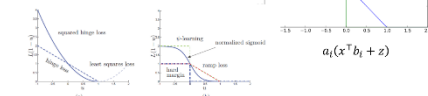$$s_i \geq 0$$

**hinge loss（和分类结果<=1 数据有关）**

- loss function + regularization term

$$\min_{x,z} \frac{1}{2}\|x\|_2^2 + C\sum_i L(b_i(x^\top a_i + z))$$

- hinge loss

$$L(u) = \max\{0, 1 - u\}$$

- the convex approximation for misclassification loss

$$a_i(x^\top b_i + z)$$

**非监督学习**
降维分析
**主成分分析 PCA（线性分有用）**
A 是数据，x 和 x^T 是主成分和无效成分，
x 就是 C 的特征向量

- find a direction that maximizes the data's variance

$$\max_{x^T x = 1} x^T C x$$

$$C = \frac{1}{n-1} A^T A$$

it should be zero-mean, otherwise, the covariance matrix will be ?

**局部线性嵌入 LLE（流形学习）**
1 用一个点周围的点的加权和表示它
（用最小二乘法得到权重系数）
2 通过低维的线性变换矩阵 b 尽量保
留所有点和周围点的加权关系

$$M = (I - W)^T (I - W)$$

W 是权重系数矩阵，特征向量为
低维特征向量

- find the relationship $a_0 = \sum x_i a_i$

$$\min_x \|a_0 - \sum x_i a_i\|_2$$

- keep the relationship in reduction

$$\min_b \|b_0 - \sum x_i b_i\|_2$$

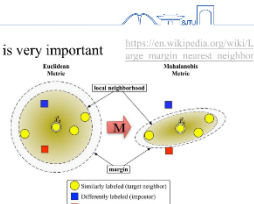**聚类 K-means 质心迭代**
**度量学习**

**Metric Learning**

- distance/similarity/dissimilarity/metric is very important
  - non-negative
  - identity
  - symmetry
  - subadditivity (triangle inequality)
- Mahalanobis distance

$$d_M(a,b) = \sqrt{(La - Lb)^\top(La - Lb)} = \sqrt{(a-b)^\top L^\top L(a-b)} = \sqrt{(a-b)^\top M(a-b)}$$

$$\min_M \sum_{i,j \in N_i} d(a_i, a_j) + \lambda \sum_{i,j,l} \xi_{ijl} \quad \text{s.t.} \begin{cases} d(a_i, a_j) + 1 \leq d(a_i, a_l) + \xi_{ijl}, \forall j \in N_i, l \notin N_i \\ \xi_{ijl} \geq 0 \\ M \succcurlyeq 0 \end{cases}$$

**自编码器（高维数据编码至低维神经元再解码）**
目标：**神经网络输出和输入一样**

**5．无约束问题解法**
1.OLS 直接求梯度为 0
Pseudo inverse 计算量大

$$\min_x \sum_{i=1}^m (a_i^\top x - y_i)^2$$

$$\min_x \|Ax - Y\|_2^2 = \min_x (Ax - Y)^\top(Ax - Y)$$

$$\nabla \|Ax - Y\|_2^2 = \nabla(Ax - Y)^\top(Ax - Y) = 2A^\top(Ax - Y) = 0$$

$$x^* = (A^\top A)^{-1} A^\top Y$$

2 梯度下降　　梯度为 0 需要 strong convexity

复杂度和 m/M 成线性关系

General descent method.
**given** a starting point $x \in \text{dom } f$.
**repeat**
1. Determine a descent direction $\Delta x$.
2. Line search. Choose a step size $t > 0$.
3. Update. $x := x + t\Delta x$.
**until** stopping criterion is satisfied.

- staring point $x^{(0)}$
- descent direction $\Delta x^{(k)}$
- step size /step length/ learning rate $t^{(k)}$

$$mI \preceq \nabla^2 f(x) \preceq MI,$$

line search   exact line search:找一条直线上 f(x)的最小值,计算量大
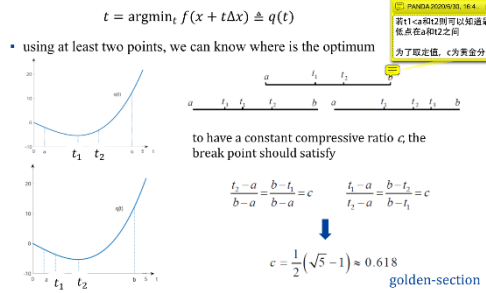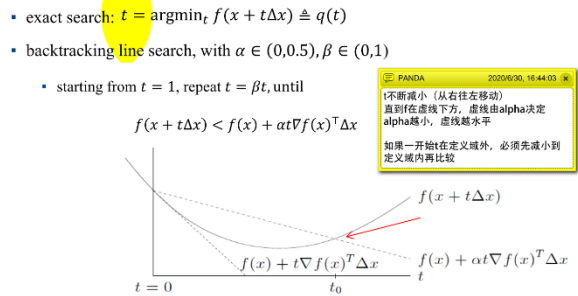


**Line Search: by Function Value**

$$t = \arg\min_t f(x + t\Delta x) \triangleq q(t)$$

- using at least two points, we can know where is the optimum

to have a constant compressive ratio $c$, the break point should satisfy

$$\frac{t_2 - a}{b - a} = \frac{b - t_1}{b - a} = c \qquad \frac{t_1 - a}{t_2 - a} = \frac{b - t_2}{b - t_1} = c$$

$$c = \frac{1}{2}(\sqrt{5} - 1) \approx 0.618$$

golden-section

**Line Search: Backtracking**

- exact search: $t = \arg\min_t f(x + t\Delta x) \triangleq q(t)$
- backtracking line search, with $\alpha \in (0, 0.5), \beta \in (0,1)$
  - starting from $t = 1$, repeat $t = \beta t$, until

$$f(x + t\Delta x) < f(x) + \alpha t \nabla f(x)^\top \Delta x$$

$f(x + t\Delta x)$

$f(x) + t\nabla f(x)^T \Delta x$   $f(x) + \alpha t \nabla f(x)^T \Delta x$

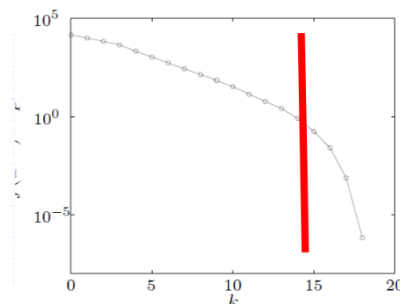$t = 0$   $t_0$   $t$

最速下降法 先假设模长和梯度，找负梯度上投影最长点（不同 norm 结果不同）

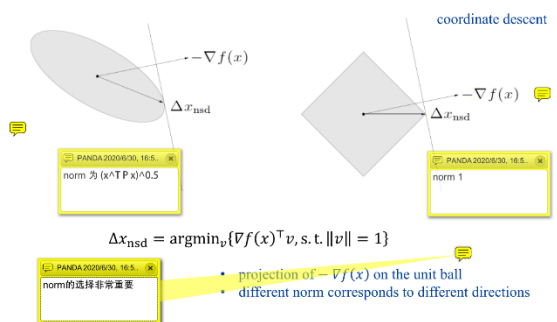**牛顿法 仿射不变** 最优解的仿射变换
是仿射变换后最优解，  $t^k$ 用 line search

$$x^{k+1} = x^k + t^k d^k \qquad d^k = \nabla^2 f(x^k)^{-1} \nabla f(x^k)$$

Damp phase 下降速度为 c 线性
Quadratic phase 速度很快，6 步到位



**Steepest Descent Method**

coordinate descent

$-\nabla f(x)$   $\Delta x_{\mathrm{nsd}}$

norm 为 (x^T P x)^0.5

$-\nabla f(x)$   $\Delta x_{\mathrm{nsd}}$

norm 1

$$\Delta x_{\mathrm{nsd}} = \arg\min_v \{\nabla f(x)^\top v, \text{ s.t. } \|v\| = 1\}$$

norm的选择非常重要

- projection of $-\nabla f(x)$ on the unit ball
- different norm corresponds to different directions

**拟牛顿法**

$$q_k = \nabla f(x_{k+1}) - \nabla f(x_k) \approx \nabla^2 f(x_k)(x_{k+1} - x_k)$$
$$\text{for quadratic case}$$
$$q_k = \nabla f(x_{k+1}) - \nabla f(x_k) = \nabla^2 f(x_k)(x_{k+1} - x_k)$$

- let $H = (\nabla^2 f(x_k))^{-1}$, there should be $Hq_k = x_{k+1} - x_k$
- we want to use rank one correction update

$$H_{k+1} = H_k + a_k z_k z_k^\top \qquad H_{k+1} q_k = x_{k+1} - x_k$$

$$H_{k+1} = H_k + \frac{(x_{k+1} - x_k - H_k q_k)(x_{k+1} - x_k - H_k q_k)^\top}{q_k^\top(x_{k+1} - x_k - H_k q_k)}$$

$$\Delta x^k = -H_{k+1}\nabla f(x_k) \quad \text{modified Newton's method with rank 1 correction}$$

**DFP 和 BFGS(都是求二阶黑塞矩阵)**

- Davidon-Fletcher-Powell uses Rank two correction

- $H_{k+1} = H_k + a_k z_k z_k^\top + \beta_k y_k y_k^\top \qquad H_{k+1} q_k = x_{k+1} - x_k$

$$H_{k+1} = H_k + \frac{(x_{k+1} - x_k)(x_{k+1} - x_k)^\top}{(x_{k+1} - x_k)^\top q_k} - \frac{H_k q_k q_k^\top H_k}{q_k^\top H_k q_k}$$

**BFGS**

**L-BFGS 直接算** $B_{k+1} \nabla f(x_{k+1})$
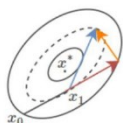
$$B_{k+1} = B_k - \frac{B_k(x_{k+1} - x_k)^\top B_k}{(x_{k+1} - x_k)^\top B_k(x_{k+1} - x_k)} + \frac{q_k q_k^\top}{q_k^\top(x_{k+1} - x_k)}$$
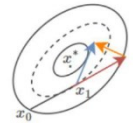
**动量法：Polyaks, Nestrov**

- Polyak's momentum algorithm

$$x_{k+1} = x_k - t_k \nabla f(x_k) + \delta_k(x_k - x_{k-1}) \qquad x_{k+1} \approx x_k - t_k \nabla f(x_k + \delta_k(x_k - x_{k-1})) + \delta_k(x^k - x^{k-1})$$



Polyak's Momentum     Nesterov Momentum

实用的 SGD 实用少数样本更新梯度 计算量小
能收敛 需要 normalization

## Adagrad(学习率一直减小)　　RMSprop（梯度大学习率小）　　Adam

$$G_{ii} = \Sigma_{t=1}^{k} \big(\nabla f(x_t)(i)\big)^2$$

$$x_{k+1} = x_k - \frac{\eta}{\sqrt{G_{ii} + \varepsilon}} \nabla f(x_k)$$

$$z_{k+1} = \rho z_k + (1-\rho)\big(\nabla f(x_k) \cdot \nabla f(x_k)\big)$$

$$x_{k+1} = x_k - \frac{\eta}{\sqrt{z_{k+1} + \varepsilon}} \nabla f(x_k) \quad \text{also for stoch}$$

**Require:** $\alpha$: Stepsize
**Require:** $\beta_1, \beta_2 \in [0,1)$: Exponential decay rates for the moment estimates
**Require:** $f(\theta)$: Stochastic objective function with parameters $\theta$
**Require:** $\theta_0$: Initial parameter vector
$\quad m_0 \leftarrow 0$ (Initialize 1st moment vector)
$\quad v_0 \leftarrow 0$ (Initialize 2nd moment vector)
$\quad t \leftarrow 0$ (Initialize timestep)
$\quad$**while** $\theta_t$ not converged **do**
$\quad\quad t \leftarrow t+1$
$\quad\quad g_t \leftarrow \nabla_\theta f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep $t$)
$\quad\quad m_t \leftarrow \beta_1 \cdot m_{t-1} + (1-\beta_1) \cdot g_t$ (Update biased first moment estimate)
$\quad\quad v_t \leftarrow \beta_2 \cdot v_{t-1} + (1-\beta_2) \cdot g_t^2$ (Update biased second raw moment estimate)
$\quad\quad \hat{m}_t \leftarrow m_t/(1-\beta_1^t)$ (Compute bias-corrected first moment estimate)
$\quad\quad \hat{v}_t \leftarrow v_t/(1-\beta_2^t)$ (Compute bias-corrected second raw moment estimate)
$\quad\quad \theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t/(\sqrt{\hat{v}_t} + \epsilon)$ (Update parameters)
$\quad$**end while**
$\quad$**return** $\theta_t$ (Resulting parameters)

## ISTA FISTA

$$\min_x \quad \lambda\|x\|_1 + \|Ax - B\|_2^2$$

$$x^{k+1} = S_\lambda\big(x^k - A^\top(Ax^k - B)\big)$$

$$S_\lambda = \arg\min_{\mathbf{x}} \left\{ g(\mathbf{x}) + \frac{L}{2}\left\| \mathbf{x} - \left( \mathbf{y} - \frac{1}{L}\nabla f(\mathbf{y}) \right) \right\|^2 \right\}$$

$$y^{k+1} = S_\lambda\big(x^k - A^\top(B - Ax^k)\big)$$
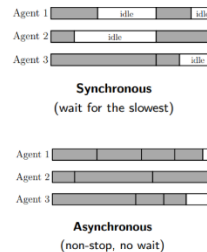
$$t^{k+1} = \frac{1 + \sqrt{1 + 4t_k}}{2}$$

$$x^{k+1} = y^{k+1} + \left(\frac{t^k - 1}{t^{k+1}}\right)(y^{k+1} - y^k)$$

$$x^{k+1} = x^k + \lambda_k A^\top(B - Ax^k)$$
$$= (I - \lambda_k A^\top A)x^k + \lambda_k A^\top B$$
$$\triangleq T_k x^k + b_k$$

## Parallel computing （同时收敛，非同时效率高）
x 没有高项就可以拆分 x 并行算

Agent 1 | idle | idle
Agent 2 | idle |
Agent 3 | | idle
**Synchronous**
(wait for the slowest)

Agent 1
Agent 2
Agent 3
**Asynchronous**
(non-stop, no wait)

## 6. 拉格朗日对偶

$$\min \quad f_0(x)$$
$$\text{s.t.} \quad f_i(x) \le 0, i = 1, \dots, m$$
$$\quad\quad h_i(x) = 0, i = 1, \dots, p$$

$$g(\lambda, \nu) = \inf_{x \in D} L(x, \lambda, \nu)$$

$$= \inf_{x \in D} \left( f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x) \right)$$

• **Primal**
$$\min \quad f_0(x)$$
$$\text{s.t.} \quad f_i(x) \le 0, i = 1, \dots, m$$
$$\quad\quad h_i(x) = 0, i = 1, \dots, p$$
$$f^* = f_0(x^*)$$

**Dual**
$$\max \quad g(\lambda, \nu) = \inf_{x \in D}(f_0(x) + v^\top f(x) + v^T h(x))$$
$$\text{s.t.} \quad \lambda \ge 0$$
$$g^* = g(\lambda^*, \nu^*)$$

## Slater 条件

Slater's constraint qualification requires the problem is strictly feasible:

$$\exists x \in \text{int } D, f_i(x) < 0, h_i(x) = 0$$

• if the problem is convex and the Slater's qualification satisfied, then there is

**strong duality**

## KKT 条件（必要，slater 满足则充要）

primal feasible
$$f_i(x^*) \le 0$$
$$h_i(x^*) = 0$$

dual feasible
$$\lambda_i^* \ge 0$$

complementary slackness
$$\lambda_i^* f_i(x^*) = 0$$

$$\nabla f_0(x^*) + \sum_i \lambda_i^* \nabla f_i(x^*) + \sum_i \nu_i^* h(x^*) = 0$$

# 7 大规模数据
## 约束（投影）梯度下降法

$$x^{k+1} = P_D(x^k - t_k \Delta x_{nt})$$

## ADMM（轮流更新两种参数）
## 增广项影响并行运算

### Alternating direction method of multipliers

▪ Dual decomposition: Lagrangian

$$L(x, z; \alpha) = \lambda \|x\|_1 + \|b - Az\|_2^2 + \alpha^T(x - z)$$

增广项，也是惩罚项

▪ Method of multipliers: augmented Lagrangian

$$L_\rho(x, z; \alpha) = \lambda \|x\|_1 + \|b - Az\|_2^2 + \alpha^T(x - z) + \frac{\rho}{2}\|x - z\|_2^2$$

strong duality
condition number adjustment

▪ ADMM

$$x^{k+1} = \operatorname{argmin}_x \ L_\rho(x, v^k; \alpha^k)$$

$$z^{k+1} = \operatorname{argmin}_z \ L_\rho(x^{k+1}, z; \alpha^k)$$

$$\alpha^{k+1} = \alpha^k + \rho(x^{k+1} - z^{k+1})$$

固定一个参数更新另一个

这个rou由KKT中的导数为
0得到

### Alternating direction method of multipliers

$$\min \quad f(x) + g(z)$$
$$\text{s.t.} \quad Ax + Bz = c$$

▪ Method of multipliers: augmented Lagrangian

$$L_\rho(x, z; \alpha) = \lambda f(x) + g(z) + \alpha^T(Ax + Bz - c) + \frac{\rho}{2}\|Ax + Bz - c\|_2^2$$

▪ ADMM

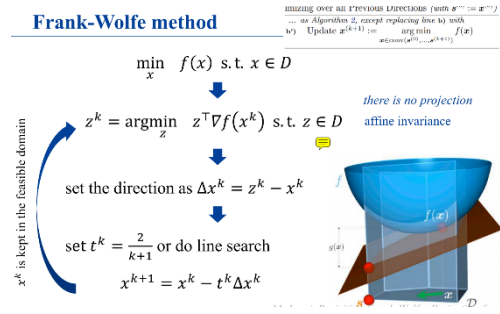$$x^{k+1} = \operatorname{argmin}_x \ L_\rho(x, z^k; \alpha^k)$$

$$z^{k+1} = \operatorname{argmin}_z \ L_\rho(x^{k+1}, z; \alpha^k)$$

$$\alpha^{k+1} = \alpha^k + \rho(Ax^{k+1} + Bz^{k+1} - c)$$

### Frank-Wolfe method

$$\min_x \ f(x) \ \text{s.t.} \ x \in D$$

$$z^k = \operatorname{argmin}_z \ z^\top \nabla f(x^k) \ \text{s.t.} \ z \in D$$

there is no projection
affine invariance

set the direction as $\Delta x^k = z^k - x^k$

set $t^k = \frac{2}{k+1}$ or do line search

$$x^{k+1} = x^k - t^k \Delta x^k$$

$x^k$ is kept in the feasible domain

x-update

$$x^{k+1} = \operatorname{argmin}_x \ L_\rho(x, z^k; \alpha^k)$$

$$= \operatorname{argmin}_x \ \lambda\|x\|_1 + x^\top \alpha^k + \frac{\rho}{2}\|x - z^k\|_2^2$$

$$= S_{\lambda/\rho}(z^k + \alpha^k/\rho)$$

▪ z-update

$$z^{k+1} = \operatorname{argmin}_z \ L_\rho(x^{k+1}, z; \alpha^k)$$

$$= \operatorname{argmin}_z \ \|b - Az\|_2^2 + z^\top \alpha^k + \frac{\rho}{2}\|x^{k+1} - z\|_2^2$$

$$= (A^\top A + \rho I)^{-1}(A^\top b + \rho x^{k+1} - \alpha^k)$$