

支持向量机

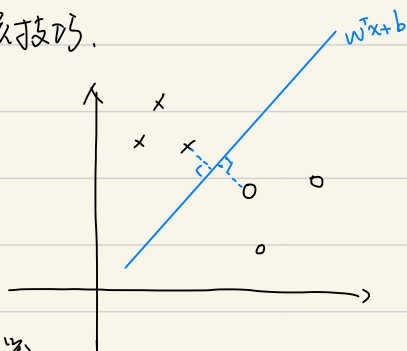
Support Vector Machine

- ① 硬间隔 SVM
- ② 软间隔 SVM
- ③ 约束优化问题

硬间隔SVM (hard margin SVM)

SVM所涉及到的东西：间隔、对偶、核技巧。

SVM $\begin{cases} \text{hard-margin SVM} \\ \text{soft-margin SVM} \\ \text{kernel SVM} \end{cases}$



其中 hard-margin SVM 又叫作最大间隔分类器

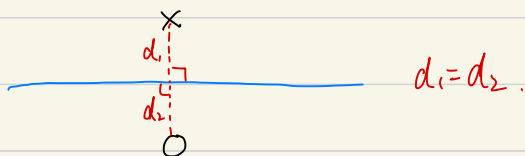
目的是寻找一个超平面，该超平面能够将不同类数据分开，且所有数据点到该超平面的最小距离最大

定义超平面为 $w^T x + b = 0$

则数据 x_i 到该超平面的距离为 $\frac{|w^T x_i + b|}{\|w\|}$

则 SVM 目标：
$$\begin{cases} \max_{w, b} \min_{x_i, y_i} \frac{1}{\|w\|} |w^T x_i + b| \\ \text{s.t. } y_i (w^T x_i + b) > 0 \end{cases}$$
 条件是超平面能够将数据分开。

稍微想一想就能知道，两类数据中，到满足上述目标的超平面的距离最小值是相等的。见下图。



若 $d_1 < d_2$ ，那么最小距离为 d_1 ，但此时的 d_1 明显不是最大值。

$$\text{优化问题} \begin{cases} \max_{w, b} \min_{x_i, y_i} \frac{1}{\|w\|} |w^T x_i + b| & \textcircled{1} \\ \text{s.t. } y_i (w^T x_i + b) > 0 & \textcircled{2} \end{cases}$$

对于①, 可以写成 $\exists r > 0, \min_{x_i, y_i} y_i (w^T x_i + b) = r$.

则①也可以写成 $\max_{w, b} \frac{1}{\|w\|} \min_{x_i, y_i} y_i (w^T x_i + b) = \max_{w, b} \frac{1}{\|w\|} \min r$.

令 $r=1$, 则② = s.t. $y_i (w^T x_i + b) \geq 1$

① = $\max_{w, b} \frac{1}{\|w\|} = \min_{w, b} \frac{1}{2} w^T w$

→ 为了后面求导时约去系数

$w^T x_i + b = 0$ 与 $2w^T x_i + 2b = 0$ 代表的是同一个超平面.

所以不论 r 为多少, 我们都可以通过放缩来使得 $r=1$.

则目标为
$$\begin{cases} \min_{w, b} \frac{1}{2} w^T w & \text{二次} \\ \text{s.t. } y_i (w^T x_i + b) \geq 1, \quad \forall i=1, 2, \dots, N. & \text{N个约束, 线性.} \end{cases} \quad (1)$$

凸优化问题.

关于凸优化问题, 见如下链接

https://blog.csdn.net/xu_fengyu/article/details/84727096

针对于 w, b 而言

很自然地, 想到用拉格朗日乘子法 将有约束问题转化为无约束问题

$$\mathcal{L}(w, b, \lambda) = \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i [1 - y_i (w^T x_i + b)], \quad \lambda_i \geq 0.$$

$$\lambda = \begin{pmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_N \end{pmatrix}$$

因为上面有 N 个约束.

$$\text{则目标转化为} \begin{cases} \min_{w, b} \max_{\lambda} L(w, b, \lambda) \\ \text{s.t. } \lambda_i \geq 0 \end{cases} \quad (2)$$

新目标与原目标等价的探讨

由于 $\lambda_i \geq 0$. 则 当 $1 - y_i(w^T x_i + b) > 0$ 时, $\max_{\lambda} L(w, b, \lambda) = \infty$
 当 $1 - y_i(w^T x_i + b) \leq 0$ 时, $\max_{\lambda} L(w, b, \lambda) = 0$.

$$\therefore \min_{w, b} \max_{\lambda} L(w, b, \lambda) = \min_{w, b} (\infty, \frac{1}{2} w^T w) = \min_{w, b} \frac{1}{2} w^T w$$

隐式地, 此时 $1 - y_i(w^T x_i + b)$ 必定小于等于 0

也就是说, 加入 \max_{λ} 和 $\lambda_i \geq 0$ 的条件是为了在求解的过程中,
 将 $1 - y_i(w^T x_i + b) > 0$ 这种情况给屏蔽掉

\therefore 两者是等价的.

利用对偶性质, 得到对偶问题为

$$\begin{cases} \max_{\lambda} \min_{w, b} L(w, b, \lambda) \\ \text{s.t. } \lambda_i \geq 0 \end{cases} \quad (3)$$

一般地, $\min \max L \geq \max \min L$ (弱对偶)

但是, 原问题目标函数是二次的, 约束是线性的 (凸优化二次问题)

则此时, $\min \max L = \max \min L$ (强对偶).

则对其进行求解:
$$\begin{cases} \max_{\lambda} \min_{w, b} \mathcal{L}(w, b, \lambda) \\ \text{s.t. } \lambda_i \geq 0 \end{cases}$$

先解 $\min_{w, b} \mathcal{L}(w, b, \lambda)$,

$$\frac{\partial \mathcal{L}(w, b, \lambda)}{\partial b} = \frac{\partial}{\partial b} \left[\sum_{i=1}^N \lambda_i (1 - y_i (w^T x_i + b)) \right]$$

$$= \frac{\partial}{\partial b} \left[\sum_{i=1}^N \lambda_i y_i b \right]$$

$$= \sum_{i=1}^N \lambda_i y_i$$

令 $\frac{\partial \mathcal{L}(w, b, \lambda)}{\partial b} = 0$, $\{ \frac{0}{0} \}$, $\sum_{i=1}^N \lambda_i y_i = 0$. (a)

将(a)代入 $\mathcal{L}(w, b, \lambda)$. $\{ \frac{0}{0} \}$:

$$\mathcal{L}(w, b, \lambda) = \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \lambda_i y_i w^T x_i$$

令 $\frac{\partial \mathcal{L}(w, b, \lambda)}{\partial w} = w - \sum_{i=1}^N \lambda_i y_i x_i = 0$.

$\{ \frac{0}{0} \}$ $w = \sum_{i=1}^N \lambda_i y_i x_i$ (b)

将(b)代入 $\mathcal{L}(w, b, \lambda)$. $\{ \frac{0}{0} \}$:

$$\begin{aligned} \min_{w, b} \mathcal{L}(w, b, \lambda) &= \frac{1}{2} \left(\sum_{i=1}^N \lambda_i y_i x_i \right)^T \sum_{j=1}^N \lambda_j y_j x_j + \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \lambda_i y_i \left(\sum_{j=1}^N \lambda_j y_j x_j \right)^T x_i \\ &= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j - \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_j^T x_i + \sum_{i=1}^N \lambda_i \\ &= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_{i=1}^N \lambda_i \end{aligned}$$

则问题变为
$$\begin{cases} \max_{\lambda} -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_{i=1}^N \lambda_i \end{cases} \quad (4)$$

带约束 (1) $\xLeftrightarrow{\text{无约束}}$ (2) $\xLeftrightarrow{\text{强对偶}}$ (3) $\xLeftrightarrow{\quad}$ (4)

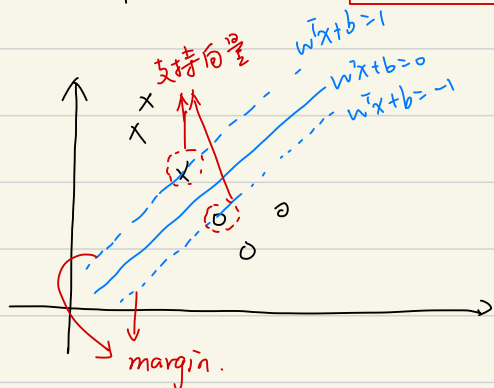
原问题与对偶问题是强对偶关系 \iff 满足KKT条件.

则该问题的KKT条件为.

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial w} = 0 \\ \frac{\partial \mathcal{L}}{\partial b} = 0 \\ \lambda_i [1 - y_i (w^T x_i + b)] = 0 \quad \text{松弛互补条件} \\ \lambda_i \geq 0 \\ 1 - y_i (w^T x_i + b) \leq 0 \end{cases}$$

当 $1 - y_i (w^T x_i + b) < 0$ 时, $\lambda_i = 0$
 当 $1 - y_i (w^T x_i + b) = 0$ 时, $\lambda_i < 0$.

由上述计算, 已经得出 $W^* = \sum_{i=1}^N \lambda_i y_i x_i$, 下面需要求解 b^* .



由图可知, 将支持向量代入 $w^T x + b$, 结果为1
 即当 x_k 为支持向量时.

$$w^T x_k + b = 1 \Rightarrow \begin{cases} b^* = 1 - w^T x_k \\ = 1 - \sum_{i=1}^N \lambda_i y_i x_i^T x_k. \end{cases}$$

从KKT条件中还可以看出, 对于不在margin上的数据, 其对应的 $\lambda_i = 0$.

对(4)进行求解可以得到 λ , 这里不进行求解, 可以用SMO算法.

或者用现成的求解QP问题的工具就可以得到

支持向量机就是一个QP问题 (Quadratic Programming, 二次规划)

软间隔SVM (Soft margin SVM)

软间隔SVM允许SVM出现错误。通过在优化条件中加入损失项来达到该目的。即 $\min \frac{1}{2} W^T W + \text{loss}$ $\rightarrow y_i (W^T x_i + b) < 1$

① 一般地，会用 0/1 损失 来统计出错的个数。

$$\text{loss} = \sum_{i=1}^N \mathbb{I} \{ y_i (W^T x_i + b) < 1 \}$$

$\mathbb{I}\{\}$ 为指示函数， z 为真时为 1， z 为假时为 0。

$$\text{令 } z = y_i (W^T x_i + b)$$

图像如下所示。



可见，其不连续。
不利于求解。

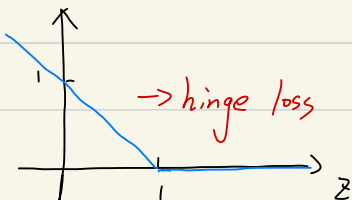
② 当统计个数行不通时，使用距离来代替。

当 $y_i (W^T x_i + b) \geq 1$ 时， $\text{loss} = 0$

当 $y_i (W^T x_i + b) < 1$ 时， $\text{loss} = 1 - y_i (W^T x_i + b)$ 。 $\frac{1}{\|W\|} y_i (W^T x_i + b)$ 是距离，但 $\frac{1}{\|W\|}$

$$\text{则 } \text{loss} = \sum_{i=1}^N \max(0, 1 - y_i (W^T x_i + b))$$

图像如下图所示。



对所有 x_i 都一样，所以可以只考虑 $y_i (W^T x_i + b)$ ， $1 - y_i (W^T x_i + b)$ 就等价于当前点与 margin 之间的距离。

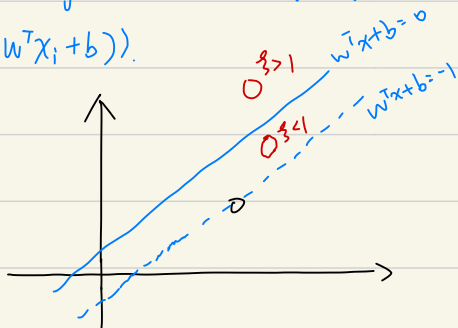
引入松弛变量, $\xi_i \geq 0$

$\xi_i = 1 - y_i(w^T x_i + b)$, 由于 $\xi_i \geq 0$, 故当 $y_i(w^T x_i + b) \geq 1$ 时, $\xi_i = 0$.

或者可以直接看作 $\xi_i = \max(0, 1 - y_i(w^T x_i + b))$.

此时软间隔的目标问题为

$$\begin{cases} \min_{w, b} \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i \\ \text{s.t. } y_i(w^T x_i + b) \geq 1 - \xi_i \end{cases}$$



显然, 对于每一个 x_i 都有一个对应的 ξ_i , 来反映其不满足约束的程度。 C 是常数, 当 $C \rightarrow \infty$ 时, 会迫使所有 $\xi_i = 0$, 此时约束条件又会变为硬间隔的约束条件; 当 C 取有限常数时, 允许有一定的误差。

约束优化问题

(Constraint Optimization)

弱对偶性证明

约束优化问题 (原问题, primal problem)

$$\text{一般表达} \begin{cases} \min_x f(x) \\ \text{s.t. } m_i(x) \leq 0, i=1, 2, \dots, M \\ n_j(x) = 0, j=1, 2, \dots, N \end{cases}$$

原问题是关于 x 的函数.

拉格朗日函数 (Lagrange)

$$\mathcal{L}(x, \lambda, \eta) = \underbrace{f(x)}_{\text{target function}} + \sum_{i=1}^M \underbrace{\lambda_i}_{\lambda_i \geq 0} m_i(x) + \sum_{j=1}^N \underbrace{\eta_j}_{\text{无要求}} n_j(x)$$

则使用拉格朗日乘子法将约束优化问题转化为 (依旧是原问题)

$$\begin{cases} \min_x \max_{\lambda, \eta} \mathcal{L}(x, \lambda, \eta) \\ \text{s.t. } \lambda_i \geq 0 \end{cases}$$

原问题对 x 的无约束的开放式
仍然是关于 x 的函数.

在逻辑上考察两者的等价性:

$$\begin{cases} \text{对于 } m_i(x), \text{ 若 } m_i(x) > 0, \text{ 则 } \max_{\lambda, \eta} \mathcal{L} = \infty. \text{ 若 } m_i(x) \leq 0, \text{ 则 } \max_{\lambda, \eta} \mathcal{L} \text{ 存在} \\ \text{对于 } n_j(x), \text{ 若 } n_j(x) \neq 0, \text{ 则 } \max_{\lambda, \eta} \mathcal{L} = \infty, \text{ 若 } n_j(x) = 0, \text{ 则 } \max_{\lambda, \eta} \mathcal{L} \text{ 存在} \end{cases}$$

$$\text{即 } \min_x \max_{\lambda, \eta} \mathcal{L}(x, \lambda, \eta) = \min_x (\infty, \infty, \max_{\lambda, \eta} \mathcal{L}) = \min_x \max_{\lambda, \eta} \mathcal{L}.$$

则只在约束优化的条件下, 才成立, 故两者等价

2. 对偶问题 (dual problem) 为

$$\begin{cases} \max_{\lambda, \eta} \min_x \mathcal{L}(x, \lambda, \eta) \\ \text{s.t. } \lambda_i \geq 0. \end{cases} \quad \text{对偶问题为关于 } \lambda, \eta \text{ 的函数}$$

弱对偶性: 对偶问题 \leq 原问题.

$$\max_{\lambda, \eta} \min_x \mathcal{L}(x, \lambda, \eta) \leq \min_x \max_{\lambda, \eta} \mathcal{L}(x, \lambda, \eta).$$

$$\text{证明: } \underbrace{\min_x \mathcal{L}(x, \lambda, \eta)}_{A(\lambda, \eta)} \leq \mathcal{L}(x, \lambda, \eta) \leq \underbrace{\max_{\lambda, \eta} \mathcal{L}(x, \lambda, \eta)}_{B(x)}.$$

$$\text{由于 } A(\lambda, \eta) \leq B(x)$$

$$\therefore \max A(\lambda, \eta) \leq \min B(x), \text{ 即对偶问题} \leq \text{原问题}$$

弱对偶性的几何解释

$$\text{原问题} \begin{cases} \min_x f(x) \\ \text{s.t. } m_1(x) \leq 0 \end{cases} \quad \text{定义域: } D = \text{dom } f \cap \text{dom } m_1$$

忽略等式约束

拉格朗日函数: $\mathcal{L}(x, \lambda) = f(x) + \lambda m_1(x)$, $\lambda \geq 0$.

$$\text{对偶问题} \begin{cases} \max_{\lambda} \min_x \mathcal{L}(x, \lambda) \\ \text{s.t. } \lambda \geq 0 \end{cases}$$

$$\begin{aligned} p^* &= \min_x f(x) && \text{原问题的最优解} \\ d^* &= \max_{\lambda} \min_x \mathcal{L}(x, \lambda) && \text{对偶问题的最优解} \end{aligned}$$

假设一个集合 $G = \{(m_1(x), f(x)) \mid x \in D\}$.

$$\text{令 } u = m_1(x), t = f(x) \Rightarrow G = \{(u, t) \mid x \in D\}.$$

$$p^* = \inf \{t \mid (u, t) \in G, u \leq 0\}$$

集合的下确界

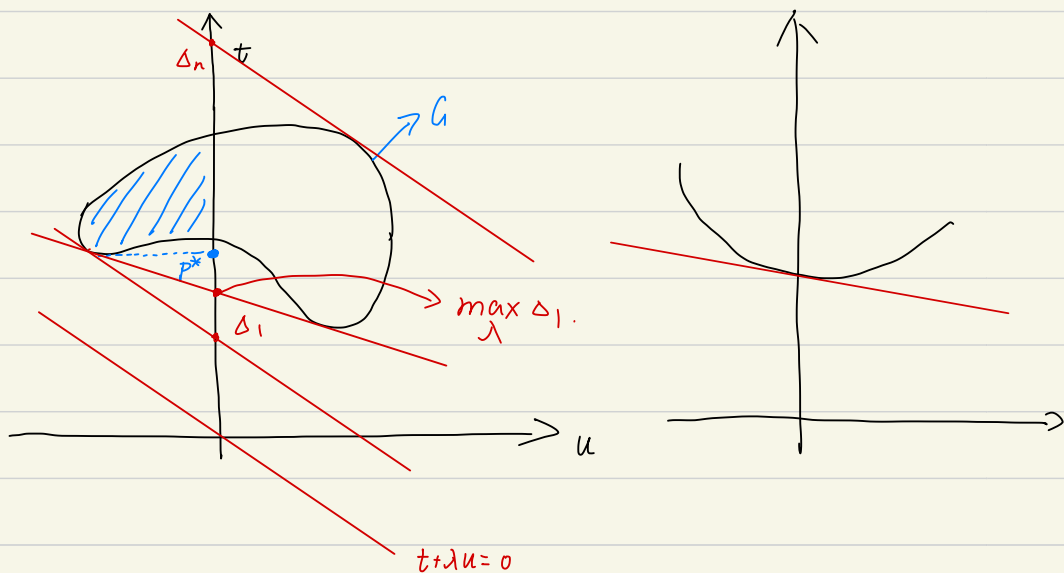
$$\begin{aligned} d^* &= \max_{\lambda} \min_x \mathcal{L}(x, \lambda) \\ &= \max_{\lambda} \min_x (t + \lambda u) \end{aligned}$$

$$\text{令 } g(\lambda) = \min_x (t + \lambda u)$$

$$d^* = \max_{\lambda} (t + \lambda u)$$

$$g(\lambda) = \inf \{t + \lambda u \mid (u, t) \in G\}.$$

对上一页的定义进行了可视化. (G 定义为非凸集, 更具一般性).



由左图可知: $\{t + \lambda u \mid (u, t) \in G\} = \{\Delta_1, \dots, \Delta_n\}$.

$$g(\lambda) = \inf \{t + \lambda u \mid (u, t) \in G\} = \Delta_1$$

$$\text{则 } d^* = \max_{\lambda} \Delta_1$$

如何取 λ , 使得 Δ_1 最大? 就等同于, 通过改变 λ , 使得直线进行旋转, 直到与 G 同时出现两个切点, (这是针对于所画的图像而言, 不同集合使得 Δ_1 最大的条件会有些许不同).

但此时, $d^* \leq p^*$.

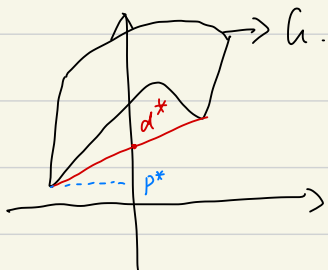
何时 $d^* = p^*$. 通过图解释的话, 就是当集合下界为凸时 (右图) 也就是说, 在坐标轴之交点处存在集合的切线.

实际上, 凸优化 + Slater 条件 $\implies d^* = p^*$.

但 Slater 条件只是充分不必要条件, 有些强对偶性不满足 Slater 条件

SVM 是一个二次规划问题, 满足上述强对偶性成立的条件

但是我有一个问题. 如果集合如下所示, 是不是就错了?



Slater 条件

定义: $\exists \hat{x} \in \text{relint } D$,

s.t. $\forall i=1, \dots, M, m_i(\hat{x}) < 0$.

relative interior 意思是不考虑集合的边界

\Rightarrow 几何意义: 对于上页图例, Slater 规定了必须有一个点 (a, b) 在 t 轴之左侧.

① 对于大多数凸优化问题, Slater 条件成立

② 放松的 Slater: 若 n 中有 K 个仿射函数, 则只需校验另外 $M-K$ 个函数是否满足条件.
最高次数为 1 的多项式函数.

则 SVM 的目标函数 $f(x)$ 是凸的, 不等式约束是仿射函数, 则 SVM 与其对偶函数是强对偶关系.

偶函数是强对偶关系.

KKT 条件

$$\text{原问题} \begin{cases} \min_x f(x) \\ \text{s.t.} \quad m_i(x) \leq 0, i=1, \dots, M \\ \quad \quad n_j(x) = 0, j=1, \dots, N \end{cases}$$

$$\text{拉格朗日: } \mathcal{L}(x, \lambda, \eta) = f(x) + \sum_{i=1}^M \lambda_i m_i(x) + \sum_{j=1}^N \eta_j n_j(x), \lambda_i \geq 0.$$

$$\text{最优解: } p^* \rightarrow x^* \quad d^* \rightarrow \lambda^*, \eta^*$$

$$\text{Convex} + \text{slater} \Rightarrow \text{strong duality} \Leftrightarrow \text{KKT 条件}.$$

$$\text{KKT:} \begin{cases} \text{① 满足可行条件} \\ \text{② 互补松弛} \\ \text{③ 梯度为0.} \end{cases}$$

$$(\lambda^* = p^*)$$

$$\text{① 可行条件: } m_i(x^*) \leq 0, \eta_j(x^*) = 0, \lambda^* \geq 0$$

$$\begin{aligned} \text{② 互补松弛: } d^* &= \max_{\lambda, \eta} g(\lambda, \eta) = g(\lambda^*, \eta^*) = \min_x \mathcal{L}(x, \lambda^*, \eta^*) \\ &\leq \mathcal{L}(x^*, \lambda^*, \eta^*) \\ &= f(x^*) + \sum_{i=1}^M \lambda_i^* m_i(x^*) + \sum_{j=1}^N \eta_j^* n_j(x^*) \\ &\leq \underbrace{f(x^*)}_{p^*} \quad \underbrace{\sum_{i=1}^M \lambda_i^* m_i(x^*)}_{\leq 0} \quad \underbrace{\sum_{j=1}^N \eta_j^* n_j(x^*)}_{=0} \end{aligned}$$

$$x: d^* = p^*.$$

∴ 上述推导中的两个“ \leq ”号都需要取“ $=$ ”

$$\text{对于第二个“}\leq\text{”号, 如果其为“}=\text{”, 那么 } \lambda_i^* m_i(x^*) = 0$$

即为互补松弛的条件

③ 梯度为0. 还是看上述②中的推导.

当第一个" \leq "号取" $=$ "时.

$$\min_x \mathcal{L}(x, \lambda^*, \eta^*) = \mathcal{L}(x^*, \lambda^*, \eta^*)$$

且 $m_i(x)$, $\eta_j(x)$ 都收敛 (声明原问题时忘说了)

$$\text{故} \left| \frac{\partial \mathcal{L}(x, \lambda^*, \eta^*)}{\partial x} \right|_{x=x^*} = 0$$

即为梯度为0的条件.

综上: KKT:

$$\begin{cases} m_i(x^*) \leq 0 \\ \eta_j(x^*) = 0 \\ \lambda \geq 0 \\ \lambda_i^* m_i(x^*) = 0 \\ \frac{\partial \mathcal{L}(x, \lambda^*, \eta^*)}{\partial x} \Big|_{x=x^*} = 0 \end{cases}$$