

行列式

$$C(X_1 \dots X_n) \text{ 逆序对个数} \quad \left| \begin{array}{cccc} a_{11} & & & \\ & \ddots & & \\ & & \ddots & \\ & & & a_{nn} \end{array} \right| = \sum (-1)^{\tau(j_1 \dots j_n)} a_{1j_1} \dots a_{nj_n}$$

$$\left| \begin{array}{cccc} a_1+b_1 & \dots & a_n+b_n \\ c_1 & \dots & c_n \\ d_1 & \dots & d_n \end{array} \right| = \left| \begin{array}{ccc} a_1 & \dots & a_n \\ c_1 & \dots & c_n \\ d_1 & \dots & d_n \end{array} \right| + \left| \begin{array}{ccc} b_1 & \dots & b_n \\ c_1 & \dots & c_n \\ d_1 & \dots & d_n \end{array} \right| \quad \begin{matrix} \text{放缩 2 次} \\ \det(AB) = \det(A)\det(B) \end{matrix}$$

把某行 k 倍加到另一行，行列式不变

如果 2 列成比例，那么行列式为 0

余子式 $\left| \begin{array}{ccccc} \dots & \vdots & \dots & \vdots & \dots \\ & i & & j & \\ & \vdots & & \vdots & \vdots \end{array} \right| = M_{ij} : \text{去掉 } i \text{ 行 } j \text{ 列的 } \det$

C_{ij} 代数余子式为 $(-1)^{i+j} M_{ij}$ $|A|$ 为一行一列 $a_{ij} C_{ij}$

$\boxed{\nabla}$ $\boxed{\Delta}$ 对角线的积 $\boxed{\square}$ $\boxed{\triangle}$ $(-1)^{\frac{n(n-1)}{2}}$ 对角线（交换）

$$\left| \begin{array}{cc} A & 0 \\ * & B \end{array} \right| = \left| \begin{array}{cc} A & * \\ 0 & B \end{array} \right| = \det(A)\det(B) \quad \left| \begin{array}{cc} 0 & A \\ B & * \end{array} \right| = \left| \begin{array}{cc} * & A \\ B & 0 \end{array} \right| = (-1)^{mn} |A||B|$$

矩阵

$A - m \times n \quad B - n \times s \Rightarrow AB = 0 \quad r(A) + r(B) \leq n$

Sylvester: $\text{rank}(AB) \geq \text{rank}(A) + \text{rank}(B) - n$

$\text{Null}(AB) \leq (n - r(A)) + (n - r(B))$

把 2 次型写成矩阵乘法: $x_1^2 + 5x_2^2 - 4x_1x_2 = (x_1 x_2) \begin{bmatrix} 1 & -2 \\ -2 & 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$$\text{伴随矩阵 } A^* A = |A| E \quad A^* = \begin{vmatrix} A_{11} & & \\ & \ddots & \\ & & A_{nn} \end{vmatrix} \quad \text{代数余子式}$$

$$(A^{-1})^* = \frac{1}{|A|} A \quad (kA)^* = k^{n-1} A^* \quad (AB)^{-1} = B^{-1} A^{-1}$$

初等矩阵 $\begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix} \rightarrow \begin{array}{l} \text{行列} \times k \\ \text{行列倍加} \\ \text{行列对调} \end{array} \left\{ \begin{array}{l} P \\ PA \text{ 即为对 } A \text{ 相同行变换} \\ AP \text{ 则为列变换} \end{array} \right.$

经过初等变换 rank 不变

$\text{rank}(A) = \text{rank}(B) \Leftrightarrow A, B \text{ 可经过多次初等变换互换}$

所有可逆矩阵都可以看成 $E_n E_{n-1} \cdots E_1 I$, 其中 E 为初等矩阵

矩阵的秩

与矩阵乘法其实在压扁, 压扁 2 次 (AB)

$$\text{rank}(AB) \leq \min(\text{rank}(A), \text{rank}(B))$$

可以认为乘法在互相组合行, 列. 结合出来的不会 rank 更大

如果可以用 $B_1 - B_n$ 表示 $\alpha_1 - \alpha_m$ $\text{rank}(\beta) > \text{rank}(\alpha)$

如果 A 的 k 阶行列式不为 0, 那么 $\text{rank} \geq k$ $\boxed{\square} \Rightarrow \boxed{\square}$

$$\text{rank}(A+B) \leq \text{rank}(A) + \text{rank}(B)$$

$$\text{rank} \begin{bmatrix} A & 0 \\ 0 & B \end{bmatrix} = \text{rank}(A) + \text{rank}(B) \quad \dim(U+V) = \dim(U) + \dim(V) - \dim(U \cap V)$$

特征值与相似对角化

特征向量 \downarrow 特征值

实对称矩阵一定可以正交对角化 $Q^T \Lambda Q$, 特征向量相互正交

$$A+kI \Rightarrow \lambda+k \quad kA \Rightarrow k\lambda \quad A^n \Rightarrow \lambda^n \quad A^{-1} \Rightarrow \frac{1}{\lambda}$$

$$A \sim B, \det(\lambda I - A) = \det(\lambda I - B) \quad \text{rank}(A) = \text{rank}(B)$$

$$\text{trace}(A) = \sum \lambda_i = \text{trace}(B), \quad A \sim \Lambda \sim B, \quad A^n \sim B^n, \quad A^{-1} \sim B^{-1}$$

相似本质上是更换坐标系后的观测

如果一个矩阵 $A = UV^T$ (2个向量张成的)

1. $\text{rank} = 1$ (都是U的倍数) $\text{rank}(AB) \leq \min(\text{rank}(A), \text{rank}(B))$

2. $\lambda = U^T V$, 其余都是0. 3. 代表把空间映射到一个直线上

如果矩阵A对称, 那么 $\lambda_{\max} = \frac{x^T A x}{\|x\|_2} = \|Ax\|_2$ (max)

任意情况下 $\lambda_{\max} = \max \|Ax\|_2$ $\max \|Ax\|_2 = \lambda_{\max} \left(\underbrace{\frac{A+A^T}{2}}_{\text{对称部分}} \right) = 1$

正定矩阵

正定矩阵一定对称, 顺序主子式一定大于0

A 正定 $\Rightarrow A^{-1}$ 正定 $\begin{cases} \text{对称: } (A^{-1})^T = (A^T)^{-1} = A^{-1} \\ \lambda > 0, \frac{1}{\lambda} > 0 \end{cases}$

PCA

对于数据 $X \Rightarrow Cov = \frac{1}{n-1} (X - \mu)^T (X - \mu)$

$$X - \mu = \begin{bmatrix} x_{11} - \mu_1 & x_{1n} - \mu_n \\ \vdots & \vdots \\ x_{ni} - \mu_i & x_{nn} - \mu_n \end{bmatrix}$$

找到几个让 Variance 尽量大的方向



(G 特征值最大的几个特征向量

Cov 矩阵对称, 不同特征值对应的特征向量正交

相同特征值的话可以彼此正交化

$$V = \left[\frac{v_1}{\|v_1\|_2}, \frac{v_2}{\|v_2\|_2}, \dots, \frac{v_n}{\|v_n\|_2} \right]$$

$$X_{PCA} = (X - \mu) \cdot V$$

SVD $AV = \sum \sigma_i U_i V_i^T$, 其中 V 为 $A^T A$ 的特征向量

$A \Rightarrow$ 旋转(U) \Rightarrow 放缩(Σ) \Rightarrow V^T (旋转)

$$A = U \Sigma V^T \quad A^T A = V \Sigma^2 V^T \quad A A^T = U \Sigma^2 U^T$$

U, V 都为正交矩阵(旋转) Σ 是 $A^T A$ / $A A^T$ 特征值的 $\sqrt{\text{sqrt}}$

$$\Sigma = \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_m \end{bmatrix} \quad A_{m \times n} = U_{m \times m} \Sigma_{m \times n} V_{n \times n}^T$$

$$\|A\|_{\max} = \sup_{\|x\|_2=1} \frac{\|Ax\|_2}{\|x\|_2} \quad \|A\|_F = \sqrt{\sum \sigma_i^2}$$

$$\|x\|_1 = \sum |x_i|$$

$$\text{Example 1} \quad \|x\|_2 = \sqrt{\sum x_i^2}$$

Example: The following functions are matrix norms: $\|x\|_\infty = \max |x_i|$

- $\|A\|_1 = \max_j \sum_{i=1}^n |a_{ij}|$ (maximum of the ℓ_1 norms of columns).
Moreover, this is an induced matrix norm:

$$\|A\|_1 = \max_{\|x\|_1=1} \|Ax\|_1.$$

列和

- $\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)}$. Moreover, this is an induced matrix norm:

最大拉伸 $\leq \sigma_{\max}$ $\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2$. 这个向量 x 就是 $A^T A$ 对应 λ_{\max} 的 特征 向量

- $\|A\|_\infty = \max_i \sum_{j=1}^n |a_{ij}|$ (maximum of the ℓ_1 norms of rows).
Moreover, this is an induced matrix norm:

$$\|A\|_\infty = \max_{\|x\|_\infty=1} \|Ax\|_\infty.$$

行和

- $\|A\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^n |a_{ij}|^2}$. This is known as the Frobenius norm.

$$= \sqrt{\sum \sigma_i^2}$$



$$\min_X \|A - X\|_F \text{ s.t. } \text{rank}(X) \leq k \quad X = A_k = \sum_{i=1}^k \sigma_i U_i V_i^T$$

$$\text{Linear Regression } y = \theta_x^T x + \theta_0 = \theta^T \bar{x} \quad \theta^T = [\theta_x \ \theta_0] \quad \bar{x} = [x, 1]^T$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum (y_i - \theta^T \bar{x}_i)^2 = \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \| Y - X\theta \|^2$$

$$= \operatorname{argmin} (Y - X\theta)^T (Y - X\theta)$$

$$\text{导数为 } -2X^T y + 2X^T X \theta = 0 \quad X^T X \theta = X^T y$$

$$\underset{\text{invertible}}{\theta} = X^{-1} y$$

U, V 是保留奇偶值那么
多个

$$\text{Not invertible} \Rightarrow \text{pseudo-inverse} \Rightarrow X = U \Sigma V^T \Rightarrow X^+ = U \Sigma^{-1} V^T$$

Derivative of A Matrix

$$\text{Scalar by vector: } \frac{\partial \|x\|}{\partial n \times 1} = n x$$

$$\frac{\partial a^T x}{\partial x} = a \quad \frac{\partial x^T a}{\partial x} = a \quad \frac{\partial x^T A x}{\partial x} = (A + A^T) \times \frac{\partial \|x\|_2^2}{\partial x} = 2x$$

$$\text{Scalar by Matrix} \quad \frac{\partial \|x\|}{\partial n \times m} = n x m$$

$$\frac{\partial \operatorname{trace}(AB)}{\partial A} = B^T \quad \frac{\partial \operatorname{tr}(AB)}{\partial B} = A^T$$

$$\frac{\partial x^T A x}{\partial A} = x x^T \quad \frac{\partial \operatorname{tr}(A^T A B)}{\partial A} = A (B + B^T) \quad \frac{\partial \det(A)}{\partial A} = |A| \cdot (A^{-1})^T$$

vector by vector

$$\frac{\partial A x}{\partial x} = A \quad \frac{\partial x^T A}{\partial x} = A^T$$

Probability 如果 $x > 0$ $\int_0^\infty P(X > t) dt$

$$EX = \int_{-\infty}^{+\infty} x dF(x) \text{ if } \int_{-\infty}^{+\infty} |x| dF(x) \text{ 存在}$$

$$= \int_{-\infty}^{+\infty} xf(x) dx \text{ 如果 } f(x) \text{ 存在}$$

$$\text{Var}(X) = E(X - EX)^2 = EX^2 - (EX)^2 \quad \text{Corr}(X, Y) = E(X - EX)(Y - EY)$$
$$= E(XY) - EXEY$$

$$\text{Moments: } E(X^k) = \int_{-\infty}^{+\infty} x^k dF(x)$$

$$\text{Normal: } f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$$

$$(\frac{1}{\sqrt{2\pi}})^n \cdot \frac{1}{|\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

在独立情况下是 $\sigma_1 \sigma_2 \dots \sigma_n \begin{vmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{vmatrix}$

$$y \sim N(\mu, \Sigma) \quad a^T y \sim N(a^T \mu, a^T \Sigma a) \quad Ay + d \sim N(A\mu + d, A\Sigma A^T)$$

$$y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \quad y_1 \sim N(\mu_1, \Sigma_{11})$$
$$y_2 \sim N(\mu_2, \Sigma_{22})$$

$$x \sim N(\mu_x, \Sigma_x) \quad y \sim N(\mu_y, \Sigma_y) \quad \begin{pmatrix} x \\ y \end{pmatrix} \sim N\left(\begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \Sigma_x & \\ & \Sigma_y \end{pmatrix}\right) \text{ 左乘} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$
$$x+y \sim N(\mu_x + \mu_y, \Sigma_x + \Sigma_y)$$

$$\text{归一化: } (\Sigma^{\frac{1}{2}})^{-1}(y - \mu) \quad \Sigma = \sum \lambda_i e_i e_i^T \quad \Sigma^{\frac{1}{2}} = \sum \sqrt{\lambda_i} e_i e_i^T$$

$$\text{MGF: } E(e^{tx}) \quad \frac{d^n}{dt^n} \Big|_{t=0} = E(X^n) \quad \text{PGF:}$$

$$\text{泰勒展开: } f(x+d) = f(x) + f'(x)d + \frac{1}{2}f''(x)d^2 + \dots + \frac{1}{n}f^{(n)}(x+\theta d)d^n$$

$$f(x) + \nabla f(x)^T d + \frac{1}{2} d^T \nabla^2 f(x) d \dots$$

$$f(y) = f(x) + \nabla f(x)^T (y - x) + \frac{L}{2} \|y - x\|^2 \quad L-\text{Smooth}$$

$$+ \frac{1}{2} (y - x)^T \nabla^2 f(x) (y - x) + O(\|y - x\|^2)$$

$$\text{or } + \frac{1}{2} (y - x)^T \nabla f(\xi) (y - x) \quad x < \xi < y$$

Optimization

目标：做 constrained optimization

Local minimum x^* 下 $\nabla f(x^*)^T d \geq 0$ feasible d: $x + ad \in S$

Let $f \in C^1$, if x^* is a local minimum and interior point $\nabla f(x^*) = 0$

Convex 问题

只有1个 local min
唯一解

Convex set: $\alpha \in [0, 1]$

$\forall x, y, \alpha x + (1-\alpha)y \in S$ (连线都在)

Convex function:

$$f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$$

① $f+g$ convex ② $cf, c>0$

③ $f: x \in S, f(x) \leq c$ convex set. 小于的都连续区间

$$\text{等价于 } f(y) \geq f(x) + \nabla f(x)^T (y - x)$$

$$\nabla^2 f(x) \geq 0$$

在KKT这个人墙条件下才会有(KKT必要)

$$\min f(x) \quad \text{势能}$$

Active set.

$$\text{s.t. } \begin{cases} C_i(x) = 0 & i \in E \text{ 墙上} \\ C_i(x) \geq 0 & i \in I \text{ 墙内} \end{cases} \quad A(x) = \{i \in I \mid C_i(x) = 0\}$$

这些在边界上，要管用

LICQ: $\{C_i(x), i \in A(x)\}$ 线性无关 [边界上各个方向无关]

如果最优解 x^* 处 LICQ, 则 x^* 为

Constrained
unconstrained

KKT 理想情况

$$\textcircled{1} \quad \nabla f(x^*) - \sum_{i \in A(x)} \lambda_i^* \nabla C_i(x) = 0 \quad \lambda_i^* \geq 0 \quad \lambda_i^* C_i(x) = 0$$

人对墙力 满足作用力 不靠墙就退力

对偶
问题

我们把墙换成风, 工 $\lambda_i C_i(x)$ 动能, $f(x)$ 是势能

$$L(x, \lambda) = f(x) - \lambda^T C(x) \quad \text{能量为力平衡态}$$

$$q(\lambda) = \min_x L(x, \lambda) \leq P^* \quad (\text{有墙}) \quad C(x) \text{ 为距离墙的深度}$$

$\max_\lambda q(\lambda)$, 调整风力, 模拟有墙的情况

这个入是不同的力量, 入越大越难动摇(入为打破这个约束可以收益多少)

Slater 条件: 如果有一个可行点 \hat{x} 满足所有 $C_i(\hat{x}) > 0 \quad \forall i \in I$, $C_i(\hat{x}) = 0 \quad \forall i \in E$

\Rightarrow 强对偶 + 凸问题 KKT 充要, 非凸必要 墙内有空间, 波浪死

普通问题在 x^* 如果有 CQ ($LICQ \Rightarrow MFCQ$), KKT 必要

凸问题下, KKT 的解就是全局最优 \Rightarrow

\hookrightarrow 配合 Slater 可得 KKT 必要 \Leftarrow

MFCQ: 对于 $i \in E$ $\nabla C_i(x)^T d = 0 \quad \forall i \in E \quad \nabla C_i(x)^T d > 0$

有一个方向 d 可以走 防止 \checkmark 被卡在这里

$$\text{凸共轭 } f^*(y) = \sup_x \langle y, x \rangle - f(x)$$

LICQ: \hookrightarrow 线性无关

在 LP 问题中强对偶, KKT 充要, 可用来验证答案

Fenchel 对偶 $f^*(k) = \sup_x (x^T k - f(x))$ 给定余料率 R, kx 到 f(x) max d

如果 f 是 convex + closed, $f^*(x) = f(x) \quad f(x) + f^*(y) \geq \langle x, y \rangle$

对于 $\min_x f(x) + g(Ax)$ 对偶问题是 $\sup_y -f^*(-A^T y) - g^*(y)$

一般都是对 $\|x\|$, 这种不光滑不可求导的

Linear Programming

$$\min c^T x$$

$$Ax = b \quad x \geq 0$$

$$\Rightarrow \left\{ \begin{array}{l} \text{若 max 的话是原的 min} \\ \sum a_{ij} x_j \leq b_i \Rightarrow \sum a_{ij} x_j + x_{n+1} = b_i \quad x_{n+1} \geq 0 \\ (\text{没有说 } x_i \geq 0) \\ \text{有 } x_i \text{ 无限制的话, } x'_i \geq 0 \quad x''_i \geq 0 \quad x_i = x'_i - x''_i \end{array} \right.$$

Gradient descent

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla L(\theta^{(t)})$$

条件数 $K = \frac{L}{\mu}$, L 是最大曲, μ 是最小曲

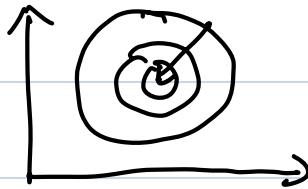


4

往垂直等高线方向走会 $2\eta - 2\log$

最优 $\frac{2}{L+\mu}$ 下

$$\theta^t - \theta^* \leq \left(\frac{1-K}{1+K} \right)^t (\theta^0 - \theta^*)$$



$$\text{收缩: } L(\theta^{(t)}) - L(\theta^*) \leq \underbrace{(1-\frac{1}{K})^t}_{\text{保 } \frac{1}{L} \text{ 下}} (L(\theta^0) - L(\theta^*))$$

optimal step size $\frac{2}{L+\mu}$

$\frac{1}{L}$

K 越大, 收缩越慢

$O(K \log \frac{1}{\epsilon})$ iterations

在 x 很小的时候 $1-x \approx e^{-x}$

$$(1-\frac{1}{K})^t \approx e^{-\frac{t}{K}}$$

给定误差 ϵ , $e^{-\frac{t}{K}} = \epsilon$ $t = -K \log \epsilon$

误差随时间指数 \downarrow $= K \log \frac{1}{\epsilon}$

SGD:

SGD 在使用 Batch 下会有

$$E(L\theta^{(t)}) \leq L\theta^t - \left(\eta - \frac{\eta n^2}{2} \right) \| \nabla L\theta^t \|_2^2 + \frac{\eta n^2 \sigma^2}{2}$$

$\eta \rightarrow 0$ Common choice for SGD

Variance Term

are $O(\frac{1}{t})$ or $O(\frac{1}{\sqrt{t}})$

误差随时间

$$\frac{K-1}{K+t} \quad \frac{\mu-L}{\mu+L}$$

对于 convergence rate, 去算 $\theta_{t+1} - \theta_t$

$$\theta_{t+1} = \theta_t - \eta \nabla f(\theta_t) \quad \nabla f(\theta_t) \approx \nabla f(\theta^*) + (\theta_t - \theta^*) \nabla^2 f(\theta^*)$$

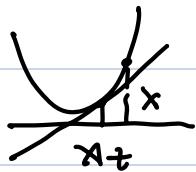
$$= \theta_t - \eta (\theta_t - \theta^*) \nabla^2 f(\theta^*) \quad e_{t+1} = (1 - \eta \nabla^2 f(\theta^*)) e_t$$

$$e_{t+1} = (1 - \eta M) e_t \quad \eta = \frac{1}{L}, (1 - \frac{1}{K})^t, \eta = \frac{2}{L+\mu}, \left(\frac{K-1}{K+t} \right)^t$$

$$\begin{aligned}
 \text{对于 } \max_{\theta} f(\theta), \text{ 算 } f(\theta_{t+1}) &= f(\theta_t) + \nabla f(\theta_t)^T (\theta_{t+1} - \theta_t) + \frac{\gamma}{2} \|\theta_{t+1}\|^2 \\
 &= f(\theta_t) - \nabla f(\theta_t)^T \cdot \eta \nabla f(\theta_t) + \frac{\gamma}{2} \eta^2 \|\nabla f(\theta_t)\|^2 \\
 \eta - \frac{\gamma}{2} \eta^2 > 0, \quad \eta < \frac{2}{\gamma}
 \end{aligned}$$

牛顿法

求零点: $g(x) \approx g(x_t) + \nabla g(x_t)(x - x_t)$



$$g(x) = 0 \quad x = x_t - \frac{g(x_t)}{\nabla g(x_t)}$$

如果要求极值, 相当于求 $\nabla g(x) = 0$

$$x = x_t - \frac{\nabla g(x_t)}{\nabla^2 g(x_t)} \leftarrow \text{不适用于多元情况 (无除法)}$$

$$L(\theta) = L(\theta_t) + \nabla L(\theta_t)(\theta - \theta_t) + (\theta - \theta_t)^T \nabla^2 L(\theta_t)(\theta - \theta_t)$$

$$\nabla L(\theta) = \nabla L(\theta_t) + \nabla^2 L(\theta_t)(\theta - \theta_t) = 0$$

$$\theta = \theta_t - (\nabla^2 L(\theta_t))^{-1} \nabla L(\theta_t)$$

$$= \theta_t - [H(\theta_t)]^{-1} \nabla L(\theta_t) \quad \text{需要严格凸才有逆}$$

这是一种收敛极快的方法, $\varepsilon_{t+1} \leq C \cdot \underline{\varepsilon_t^2}$

但是要在离的很近的时候才收敛快, 可以 GDI + Newton

算 $\nabla^2 L(\theta)$ 比较费时间, 也经常不用

$$\nabla f(x_k) = \underbrace{\nabla f(x^*)}_{0} + H(x^*)(x_k - x^*) + \mathcal{O}(\|x_k - x^*\|^2)$$

$$x_{k+1} - x^* = x_k - \underbrace{H^T(x_k)}_{\nabla f(x_k)} \underbrace{(H(x^*)(x_k - x^*) + \mathcal{O})}_{0} - x^*$$

$$e_{k+1} = \underbrace{(I - H^T(x_k) H(x^*))}_{0} e_k - \underbrace{\frac{H^T(x_k) \mathcal{O}}{C \cdot \underline{\varepsilon_k^2}}}_{\frac{H^T(x_k) \mathcal{O}}{C \cdot \underline{\varepsilon_k^2}}}$$

$$\text{Logistic Regression} \quad \sigma(x) = \frac{1}{1+e^{-x}} \quad \sigma'(x) = \sigma(x)(1-\sigma(x))$$

$$P(y=1/x) = \hat{y} = \sigma(XW) = \frac{1}{1+e^{-xW}}$$

$$L(w) = - \sum [y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)] \\ = -y^T \log \hat{y} + (1-y) \log (1-\hat{y})$$

† 定义 $P(y|w,x) = (\hat{y}_i)^{y_i} (1-\hat{y}_i)^{1-y_i}$

$$\begin{cases} y_i = 0 & 1 - \hat{y}_i \\ y_i = 1 & \hat{y}_i \end{cases}$$

$$\frac{\partial L}{\partial w} = X^T(\hat{y} - y) = X^T(\sigma(XW) - y)$$

$$\nabla^2 L = X^T S X \quad X = \text{diag}(\hat{y}_i(1-\hat{y}_i))$$

$$u^T X^T S X u = v^T S v \quad S_{ii} = \hat{y}_i(1-\hat{y}_i) > 0 \quad v^T S v > 0$$

$\nabla^2 L > 0$ convex

另一种情况 $y = \pm 1$

那么可以写成 $\sigma(YXW)$, $Y = \text{diag}(y_1 - y_n)$

$$L(w) = Y^T \log(1 + \exp(-YXW))$$

$$\nabla L(w) = \frac{Y e(-YXW)}{1 + e(-YXW)}$$

Generalization Error

$$\hat{R}(\theta) = \frac{1}{n} \sum L_i(x_i, y_i, \theta) \quad \text{是你能算出的 Empirical}$$

$$R(\theta) := E[L(x, y, \theta)] \quad \text{population risk, generalization loss}$$

$$R(\hat{\theta}) = \underbrace{\hat{R}(\hat{\theta})}_{\text{能有的}} + \underbrace{(R(\hat{\theta}) - \hat{R}(\hat{\theta}))}_{\text{Gap}}$$

Bias - Variance Trade-off

$$\text{Error} = \text{Bias}^2 + \text{Variance} + \text{Noise}$$

想要 Low variance & Low Bias

Uniform Convergence

$$n \rightarrow \infty \quad \hat{R}(\theta) \rightarrow R(\theta) \quad \sup_{\theta \in H} |\hat{R}(\theta) - R(\theta)| < \varepsilon$$