

Outline

1. Open set classification

2. EVT

① "from mean to extremes"

② limit distribution

③ MDA and GEV

④ GPD

3. EVM algorithm.

① algorithm

② limitation

4. The GPD classifier

① Settings:

② Th 3

- Th 4 (con) ~~?~~

- Hill estimator

- Th 5

③ algorithm

5. GEV classifier

6. Application

1. Open set classification (开集)

Supervised learning { regression,
classification.

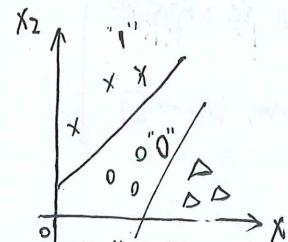
• Data: $\{(x_i, y_i)\}_{i=1}^n$, $x_i \in \mathbb{R}^P$,
 $y_i \in \{c_1, \dots, c_M\}$

• classifier $f: X \rightarrow Y$, $f \in F$

{ training set : $\{\cdot\}$
validation set : $\{\cdot\}$ hyperparameters
testing set : Recall, Precision, AUC, ...

⇒ samples are drawn from the same distribution.

eg 1. $P = 2$, $y_i \in \{"0", "1"\}$



training set: $\{x, 0\}$

testing set

closed set $\{x, 0\}$ — closed set

$\{x, 0, \Delta, \dots\}$ open set

unknown classes

all possible classes are present
during the training phase.

eg 2. face recognition.

I identification: Figure → ID 肯定是, 是谁?

V verification: 是否是员工?

• Open set classification

① distinguish between known and unknown classes.

— novelty detection: rare classes
— outlier detection

work

② in an incremental way

— life-long learning: new classes eg. KNN

③ few hyper parameters.

— cross-validation is not available.

"We know nothing of the unknown classes."

"why extreme value theory?"

"extreme features", rather than the average ones, are the most important

CV: for discriminating between different objects."

2.EVT

e.g. A tire of a car can fail in two ways: x_1, x_2, \dots, x_n i.i.d. f

- partial sums exceed some threshold

$$\frac{\frac{1}{n} \sum_{i=1}^n x_i - E x_i}{\sigma/\sqrt{n}} \xrightarrow{d} N(0,1)$$

- partial maxima ... [Block Maxima]

$M_n = \max(x_1, \dots, x_n) \rightarrow$ Unusually large (or small) levels.

$$\frac{M_n - a_n}{b_n} \xrightarrow{d} G(x) \quad (n \rightarrow \infty)$$

the behavior of x_i is usually unknown.

exact calculation on M_n is impossible

suitable assumptions, $n \rightarrow \infty, M_n \rightarrow \square$.

a_n, b_n ?

$$\text{let } x^* := \sup \{x : F(x) < 1\}$$

$$\text{then } \max(x_1, \dots, x_n) \xrightarrow{P} x^*, n \rightarrow \infty$$

$$\begin{aligned} \text{since } P\{M_n \leq x\} &= P\{x_1 \leq x, \dots, x_n \leq x\} \\ &= F^n(x) \rightarrow 0, \text{ for } x < x^* \\ &\rightarrow 1, \text{ for } x \geq x^* \end{aligned}$$

we need a normalization.

Suppose $\exists a_n > 0, b_n \in \mathbb{R}, n=1,2,\dots$

such that $\frac{M_n - b_n}{a_n}$ has a nondegenerate limit distribution

as $n \rightarrow \infty$,

$$\text{i.e., } \lim_{n \rightarrow \infty} F^n(a_n x + b_n) = G(x)$$

The class of F satisfying this condition is called the

Maximum domain of attraction of G

MDA(G)

$$F \in \text{MDA}(G)$$

$$G(x) = \exp \left\{ - \left[1 + \frac{x-\mu}{\theta} \right]^{-\frac{1}{\gamma}} \right\}$$

$$G(z) = \exp \left\{ - \left[1 + \frac{z-\mu}{\theta} \right]^{-\frac{1}{\gamma}} \right\}, z \in \mathbb{R}; 1 + \frac{z-\mu}{\theta} > 0$$

$\beta, \mu \in \mathbb{R}, \theta > 0$. (shape, location, scale)

① $\gamma = 0$, Gumbel

$$G(z) = \exp \left\{ - \exp \left[- \frac{z-\mu}{\theta} \right] \right\} \text{ for } z < \mu$$

② $\gamma > 0$, Frechet

$$G(z) = \begin{cases} \exp \left\{ - \left(\frac{z-\mu}{\theta} \right)^{-\frac{1}{\gamma}} \right\}, & z > \mu \\ 0, & z \leq \mu \end{cases}$$

③ $\gamma < 0$, reversed Weibull

$$G(z) = \begin{cases} \exp \left\{ - \left[\frac{z-\mu}{\theta} \right]^{\frac{1}{\gamma}} \right\}, & z > \mu \\ 1, & z \leq \mu \end{cases}$$

EVT-II.

Th1. Block Maxima

$$G(x) = \exp \left\{ - \left[1 + \frac{x-\mu}{\theta} \right]^{-\frac{1}{\gamma}} \right\}, \quad x \in \mathbb{R}; \quad 1 + \frac{x-\mu}{\theta} > 0$$

$\gamma \in \mathbb{R}$: shape, $\mu \in \mathbb{R}$, $\theta > 0$ | γ : different representations of extreme value behavior.

① $\gamma = 0$, Gumbel, Type 1

$$G_1(x) = \exp \left\{ - \exp \left[- \frac{x-\mu}{\theta} \right] \right\}$$

Exponential decay in the tail of $F(x)$

$$F \in \text{MDA}(G_1) \text{ iff } \lim_{x \rightarrow \infty} \frac{1 - F(t+x)}{1 - F(x)} = e^{-t}, \text{ for all } t > 0.$$

② $\gamma > 0$, Fréchet, Type 2

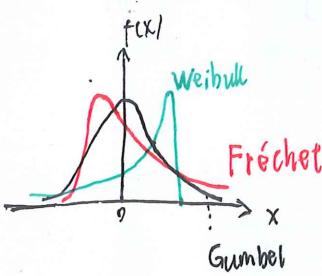
$$F \in \text{MDA}(G_2) \text{ iff } \lim_{x \rightarrow \infty} \frac{1 - F(tx)}{1 - F(x)} = t^{-\frac{1}{\gamma}}, \text{ for all } t > 0.$$

Polynomial decay in the tail of $F(x)$. "fat tail"

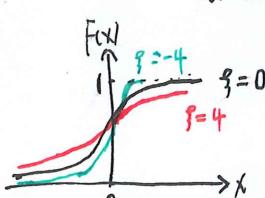
③ $\gamma < 0$, reversed Weibull, Type 3. (Lemma 1)

$$F \in \text{MDA}(G_3) \text{ iff } \lim_{x \rightarrow \infty} \frac{1 - F(x^* - \frac{t}{x})}{1 - F(x^* - \frac{1}{x})} = t^{-\frac{1}{\gamma}}, \text{ for all } t > 0$$

$x^* = \sup \{x : F(x) \leq 1\}$ $F(x)$ has a finite upper endpoint of x^* .



$P\{X > x_p\} = p$, $p = 0.01$: 百年一遇, 大坝.



Th2. Peak over threshold

$Y := X - u \mid X > u$: more efficient use of data.

$$H(y) = 1 - \left(1 + \frac{y}{\theta} \right)^{-\frac{1}{\gamma}}, \quad y > 0; \quad 1 + \frac{y}{\theta} > 0$$

$$\bar{G} = \frac{\gamma}{\bar{\gamma}} (u - M)$$

need to choose u .

4. The extreme value machine (Rudd et al. 2016)

$$x_i \in \mathbb{R}^p, y_i \in \{c_1, \dots, c_k\}, i=1, \dots, n$$

x_0 : a new point (known or unknown)

since $\bar{M}^{(i)}$ is bounded above by 200.

$$M^{(i)} = \min_{j: y_j \neq y_i} D_j^{(i)} = \min_{j: y_j \neq y_i} \frac{1}{2} \|x_i - x_j\|.$$

↑
that is, assuming $z^* = 0$. (lemma)

$$\bar{M}^{(i)} = -M^{(i)} \xrightarrow{d} W^{(i)}(z) = \begin{cases} \exp\left\{-\left(\frac{0-z}{b_i}\right)^{\frac{1}{\theta}}\right\}, & \text{if } z < 0 \\ 1, & \text{if } z \geq 0 \end{cases} \quad (\text{fit it to the } k \text{ largest } -D_j^{(i)} \text{ for each point } x_i)$$

$(\hat{c}_i, \hat{w}_i), \hat{W}^{(i)}$

" x_0 落入任一训练样本点的 Margin distance, 即 known"

then, label x_0 as known if $\max_{i=1, \dots, n} \hat{W}^{(i)}(-\|x_0 - x_i\|) \geq \delta$

If x_0 is likely enough to fall in at least one of the margins, of the n training data, then label it as known.

limitation.

① $z^* = 0$: all classes in the training set share the same support.

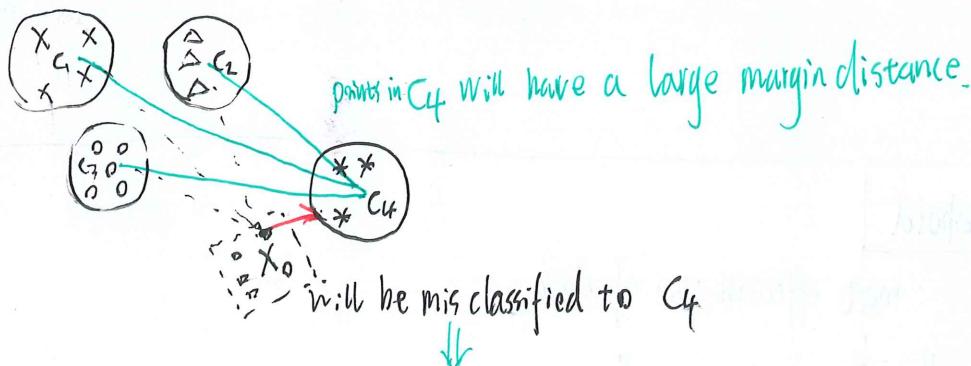
(for each i) \Rightarrow It's impossible two classes are perfectly separated.

② Estimating the upper tail of $\bar{M}^{(i)}$ based on the k largest observed $-D_j^{(i)}$ for each x_i should rely on Th2. GP instead of GEV. peaks over threshold ✓

Block Maxima.

③ δ should be chosen by fixing the type-I error.

④ A non-justified premium: classes far from all the others.



We can not use the geometry given by the known classes to do open set classification.

↓
at unknowns 假设

对 unknown - 无所知, 故不可假设距离关系!

5. GPDC

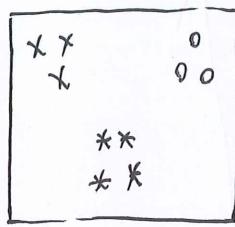
(5)

Intuition: Test points that are extremely far from the training classes are more likely to be **unknown**.

$$\{(x_i, y_i)\}_{i=1}^n, x_i \in \mathbb{R}^p, y_i \in \{c_1, \dots, c_J\}$$

$$x_0, D_i = \|x_i - x_0\| \quad (\text{consider Euclidean distance for simplicity})$$

then the **lower endpoint** of the distribution D of the D_i $\left\{ \begin{array}{ll} = 0 & : x_0 \text{ is known} \\ \neq 0 & : \text{unknown.} \end{array} \right.$



training set

$\int_{x \in A} f_g(x) dx$: the probability of a point in C_j falls in the set $A \subseteq \mathbb{R}^p$

training data generating process: $f(x) = \sum_{j=1}^J w_j f_g(x), w_j \in [0, 1], \sum_{j=1}^J w_j = 1$



Ball $B_\delta(x_0)$

↓
approximate directly f .

$$P\{B_\delta(x_0)\} = P(D < \delta) = P(-D > -\delta)$$

Th3. Consider D_1, \dots, D_n and assume $x_0 \in \text{Supp}(f)$, let $R = -D$.

↓ that is, (the upper endpoint of $-D$ is zero.) $r^* = \sup\{r : F(r) < 1\}$

$\hat{\xi}_n$ (MLE) let $Y := R_1 - u | R > u, H(y) = 1 - (1 + \frac{y}{\theta})^{-\frac{1}{\beta}}, y > 0, 1 + \frac{y}{\theta} > 0$

$$h(y) = \frac{1}{\theta} (1 + \frac{y}{\theta})^{-\frac{1}{\beta}-1}$$

suppose y_1, \dots, y_k are k excesses of the threshold u ,

$$\text{since } \beta \neq 0, \log L(y_1, \dots, y_k; \bar{\theta}, \beta) = -k \log \bar{\theta} - (1 + \frac{1}{\beta}) \sum_{i=1}^k \log (1 + \frac{y_i}{\bar{\theta}})$$

$$\text{as } r^* = u - \frac{\bar{\theta}}{\beta} = 0 \Rightarrow \bar{\theta} = u\beta = \frac{-k \log \bar{\theta}}{1 + \frac{1}{\beta}}$$

$$\Rightarrow \log L(R_1, \dots, R_n; \bar{\theta}, \beta) = -k \log u\beta - \frac{1}{\beta} \sum_{i=1}^k \log \frac{R_{(n+1-i)} - u}{u}, \text{ where } R_{(n)} \geq R_{(n-1)} \geq \dots \geq R_{(1)}$$

$$\hat{\xi}_n = \arg \max_{\beta} -k \log \beta - \frac{1}{\beta} \sum_{i=1}^k \log \frac{R_{(n+1-i)}}{u}$$

$$= \frac{1}{k} \sum_{i=1}^k \log \frac{R_{(n+1-i)}}{u}$$

Th 4. let $u = R(n-k)$, assume $k = k(n) \rightarrow \infty$, $\frac{k(n)}{n} \rightarrow 0$, as $n \rightarrow \infty$

⑥

the shape parameter of the distribution of $-D$ is $\xi = -\frac{1}{p}$,

and $\hat{\xi}_n \xrightarrow{P} -\frac{1}{p}$, where $P \in N$.

$$\hat{\xi}_n = \frac{1}{k} \sum_{i=1}^k \log \frac{R(n+1-i)}{u}$$

proof

① the distribution of $-D$ is in the max-domain of a GEV distribution with $\xi = -\frac{1}{p}$.

$R = -D$, $F_R \in \text{MDA}(GEV_3)$

$$\begin{aligned} P\{R > -s\} &= P\{D < s\} = P\left(\frac{1}{D} > \frac{1}{s}\right) \\ &= P\{B_s(x_0)\} = \int_{B_s(x_0)} f(x) dx \\ &= f(x_0) V\{B_s(x_0)\} + o(1) \\ &= f(x_0) \pi^{\frac{p}{2}} \cdot s^p / \Gamma\left(\frac{p}{2} + 1\right) + o(1) \end{aligned}$$

$$r^* := \sup\{r : F(r) < 1\} = 0$$

$$\begin{aligned} \lim_{r \rightarrow \infty} \frac{1 - F(r^* - \frac{t}{r})}{1 - F(r^* - \frac{1}{r})} &= \lim_{r \rightarrow \infty} \frac{1 - F(-\frac{t}{r})}{1 - F(-\frac{1}{r})} \\ &= \lim_{r \rightarrow \infty} \frac{f(x_0) \pi^{\frac{p}{2}} \cdot (-\frac{t}{r})^p + o(1)}{f(x_0) \pi^{\frac{p}{2}} \cdot (-\frac{1}{r})^p + o(1)} \\ &= t^p, \text{ for all } t > 0 \end{aligned}$$

$$\Rightarrow p = -\frac{1}{\xi}, \quad \xi = -\frac{1}{p}$$

$$\gamma = \frac{1}{p} > 0$$

② The upper tail of $\frac{1}{D}$ is in the max-domain of attraction of a GEV distribution with

$$z = \frac{1}{D}, \quad z^* := \sup\{z : F(z) < 1\} = \infty$$

$F_z \in \text{MDA}(G_z)$

$$\begin{aligned} \lim_{z \rightarrow \infty} \frac{1 - F(tz)}{1 - F(z)} &= \lim_{z \rightarrow \infty} \frac{f(x_0) \pi^{\frac{p}{2}} (tz)^{-p} / \Gamma\left(\frac{p}{2} + 1\right) + o(1)}{f(x_0) \pi^{\frac{p}{2}} (z)^{-p} / \Gamma\left(\frac{p}{2} + 1\right) + o(1)} \\ &= t^{-p} \end{aligned}$$

$$\Rightarrow -p = -\frac{1}{\xi}, \quad \xi = \frac{1}{p} > 0$$

A positive function C , on $[0, +\infty)$ is

① slowly varying at ∞ , if $\lim_{x \rightarrow \infty} \frac{C(tx)}{C(x)} = 1$, for all $t > 0$.

② regularly varying at ∞ with index $p \in \mathbb{R}$.

$$\text{if } \lim_{x \rightarrow \infty} \frac{C(tx)}{C(x)} = t^p, \quad t > 0.$$

$$\begin{aligned} \hat{\beta}_n &= \frac{1}{k} \sum_{i=1}^k \log \frac{R(n+i)}{n} \\ &= -\frac{1}{k} \sum_{i=1}^k \log \frac{-1/R(n+i)}{-1/R(n+k)} \rightarrow \frac{1}{D}, \quad \hat{\beta} = \frac{1}{P} \end{aligned}$$

↓
Hill estimator

(2006)

De Haan L. Extreme value theory : an introduction

$$\hat{\beta}_{\text{Hill}} \xrightarrow{P} \hat{\beta} = \frac{1}{P}$$

k: the number of the exceedances above U. 超過參數.

$$H_0: \beta_n = -\frac{1}{P} \quad (r^* = 0)$$

$$H_1: \beta_n = 0 \quad (r^* < 0), \quad x_0 \notin \text{supp}(f)$$

Th5. Suppose $r^* < 0$, $U = R(n-k)$, $k = k(n) \rightarrow \infty$, $\frac{k(n)}{n} \rightarrow 0$, as $n \rightarrow \infty$

$$\text{then } \hat{\beta}_n = \frac{1}{k} \sum_{i=1}^k \log \frac{R(n+i)}{R(n-k)} \rightarrow 0.$$

$$\begin{aligned} \text{proof [page 10].} \quad 0 &\leq -\hat{\beta}_n = \frac{1}{k} \sum_{i=1}^k \log \frac{-R(n-k)}{-R(n+i)} && R(n+i) < r^* \\ &\leq k \cdot \frac{1}{k} \log \frac{-R(n-k)}{-r^*} \rightarrow 0. && -R(n+i) > -r^* \end{aligned}$$

$$\Rightarrow \text{the first test} \quad \begin{cases} H_0: P\hat{\beta}_n = -1 \\ H_1: P\hat{\beta}_n = 0 \end{cases}$$

choose $s > -1$, if $P\hat{\beta}_n < s$, mark x_0 as possibly known, go to next step.

$$s = -0.8, \quad P\hat{\beta}_n < s, \quad \text{靠近} -1.$$

else if $P\hat{\beta}_n > s$, unknown.

(更靠近 0)

Why $P\hat{\beta}_n$ instead of $\hat{\beta}_n$?

$$-\frac{1}{P} \rightarrow 0, \text{ as } P \rightarrow \infty$$

$$P\hat{\beta}_n \sim N(-1, 1)$$

What if $P_{\mathcal{B}_n}$ is close to 1? (Can't reject H_0)

⑧

From Th3, let $u = R_{(n-k)}$

$$P(-D > x) \approx \frac{k}{n} [1 - \hat{H}(x - R_{(n-k)})] \\ \Leftrightarrow P < x = \frac{k}{n} (x/R_{(n-k)})^{-\frac{1}{k}}$$

$$Y \approx \frac{k}{n} [1 - \hat{H}(x - R_{(n-k)})]$$

$$1 - \frac{nY}{k} = \hat{H}(x - R_{(n-k)})$$

possibility that x_0
was generated from α_k class.

Proof:

$$P\{-D > x, x > R_{(n-k)}\} = \\ \cancel{P\{-D > x\}} P(-D > x | x > R_{(n-k)}) P(x > R_{(n-k)})$$

$$\approx \frac{k}{n} [1 - P(-D \leq x | x > R_{(n-k)})]$$

$$= \frac{k}{n} [1 - \hat{H}(x - R_{(n-k)})]$$

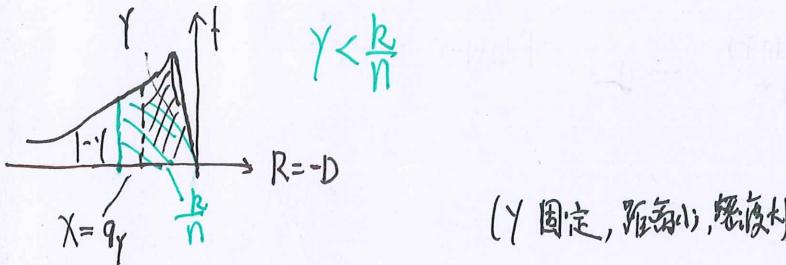
$$= \frac{k}{n} \left(1 + \frac{x - R_{(n-k)}}{\bar{G}}\right)^{-\frac{1}{k}}$$

$$\Rightarrow \bar{G} = u, u = R_{(n-k)}$$

$$= \frac{k}{n} \left(\frac{x}{R_{(n-k)}}\right)^{-\frac{1}{k}}$$

$$\begin{aligned} q_Y &= R_{(n-k)} + \hat{H}^{-1}\left(1 - \frac{ny}{k}\right) \\ &= R_{(n-k)}\left(\frac{ny}{k}\right)^{-\frac{1}{k}} \end{aligned}$$

$q_Y < 0$ 为 $(1-\gamma)$ -quantile of $-D$



$-q_Y \downarrow, q_Y \nearrow, q_Y$ 靠近 0, $f(x_0) \uparrow$

$-q_Y \uparrow, q_Y \downarrow, q_Y$ 远离 0, unknown

$$\hat{P}\{B_{-q_Y}(x_0)\} = 1 - \gamma ?$$

So, if $-q_Y > t > 0$, mark x_0 as unknown.

otherwise, known.

$-q_Y$ 越小, $f(x_0)$ 处密度越大。
(尾部)

$$\begin{aligned} f &= H(y) = 1 - \left(1 + \frac{y}{\bar{G}}\right)^{-\frac{1}{k}} \\ y &= \hat{H}^{-1}(f) \\ &= \frac{\bar{G}}{k} \left[\left(1-f\right)^{-\frac{1}{k}} - 1\right] \\ &= u \left[\left(1-f\right)^{-\frac{1}{k}} - 1\right] \end{aligned}$$

5.3. GPDC Algorithm.

1. D_1, \dots, D_n

input: ~~$\{x_i\}_{i=1}^n$~~ , x_0 .

1. $-D_1, \dots, -D_n$

$$2. \hat{\gamma}_n = -\frac{1}{k} \sum_{i=1}^k \log \left(\frac{R(n+i)}{R(n-k)} \right)$$

$$3. \text{first test: } \begin{cases} H_0: p_{\hat{\gamma}_n} = 1 \\ H_1: p_{\hat{\gamma}_n} = 0. \end{cases}$$

if $p_{\hat{\gamma}_n} < s$, where $s \geq 1$,

if $p_{\hat{\gamma}_n} \geq s$, where $s > 1$, x_0 is unknown.

else go next

$$4. q_y = R_{(n-k)} + \hat{\gamma}_n^{-1} \left(1 - \frac{1}{k} \right)$$

($1 - \frac{1}{k}$)-quantile of $-D$

$$5. \text{if } -q_y > t, \text{ where } t > 0, \quad x_0 \text{ is unknown,} \\ \text{else } x_0 \text{ is known.}$$

①

k : the number of upper order statistics
bias-variance trade-off

② s.t.: control type-I error

a jackknife fashion (instead of using test statistic under H_0)

$\hat{\gamma}_n^{(1)}, \dots, \hat{\gamma}_n^{(n)}$, $1 - \frac{2}{2}$ quantiles $\rightarrow s$

$-q_y^{(1)}, \dots, -q_y^{(n)}$, $1 - \frac{2}{2}$ quantiles $\rightarrow t$.

(等) 两重检验, Bonferroni's correction. ↗

设 H_1, \dots, H_m

let H_1, \dots, H_m be a family of hypotheses,
the familywise error rate (FWER) 族错误率

$$\overline{FWER} = P \left\{ \bigcup_{i=1}^{m_0} (P_i \leq \frac{\alpha}{m}) \right\}$$

$$\begin{aligned} \text{出现至少一个I类错} &\leq \sum_{i=1}^{m_0} P \left\{ P_i \leq \frac{\alpha}{m} \right\} \\ \text{错误的概率} &\leq m_0 \cdot \frac{\alpha}{m} \\ &\leq m \cdot \frac{\alpha}{m} = \alpha. \end{aligned}$$

P_i : 每一假设的P值.

m_0 : 实际为真的零假设总数.

6. GEV classifier

① for each x_i , calculate

$$D_i^{\min} = \min_{j \neq i} \|x_i - x_j\|$$

$\mathcal{O}(n \log n)$

$$-D_i^{\min} = \max_{j \neq i} \|x_i - x_j\| \leq 0.$$

② fit \hat{w} to $-D_1^{\min}, \dots, -D_n^{\min}$

$$③ -d_0^{\min} = \max_{i=1, \dots, n} -\|x_i - x_0\|,$$

H_0 : x_0 is known

H_1 : x_0 is unknown.

④ $\hat{w}(-d_0^{\min}) < \lambda$: reject $H_0 \Rightarrow x_0$ is unknown.

$$P(-D_1^{\min} < -d_0^{\min}),$$

K-D本对.

对于给定的精度, 查找时间 $\mathcal{O}(\log n)$