

Lappangó programhibák megkeresése gépi tanulással

Bolgár Bence

Méréstechnika és Információs Rendszerek Tanszék

Szoftver verifikáció és validáció, 2013. december 12.

Lappangó programhibák

Brun Y. et al.: Finding latent code errors via machine learning over program executions.

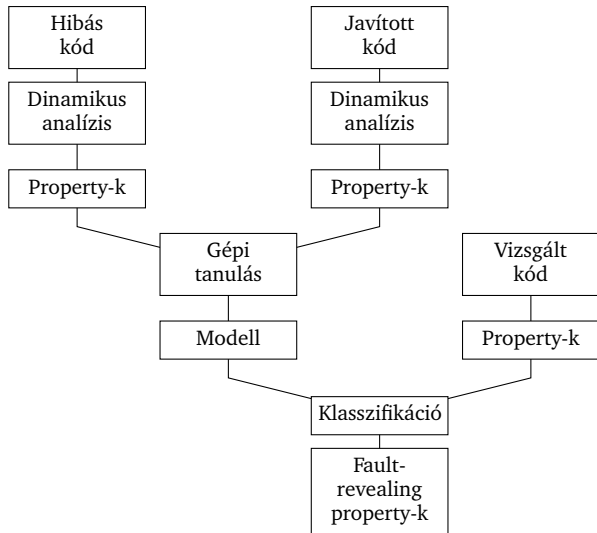
Cél: olyan programhibák megtalálása, amelyekre nincsen tesztünk. Hasznos lehet pl.:

- ▷ Ha új tesztesetek generálása nem megoldható vagy túl költséges.
- ▷ Ha az elvárt viselkedés reprezentációja nehézkes.
- ▷ Nehezen megtalálható programhibák azonosítása.

Eljárás: felügyelt tanulás

1. Tanulás: ismert hibás programkódok és javított verziók alapján
 - ▷ Pl. korábbi projektek kritikus hibáit felhasználva
 - ▷ Property-k generálása dinamikus analízissel (Daikon)
 - ▷ Modellépítés „fault-revealing” property-k megtalálására (SVM, döntési fa)
2. Klasszifikáció: vizsgálandó kód
 - ▷ Property-k klasszifikációja a modell felhasználásával
 - ▷ „Pozitív” property-k ellenőrzése manuálisan

Felügyelt tanulás



Példa

```
// Return a sorted copy of the arg.
double [] bubble_sort(double [] in)
{
    double [] out=array_copy(in);
    for (int x=out.length-1; x>=1; x--)
    {
        // lower bound should be 0, not 1
        for (int y=1; y<x; y++)
        {
            if (out[y]>out[y+1])
                swap(out[y],out[y+1]);
        }
    }
    return out;
}
```

Property	FR
$out[1] \leq in[1]$	Yes
$\forall i: in[i] \leq 100$	No
$in[0] = out[0]$	Yes
$size(out) = size(in)$	No
$in \subseteq out$	No
$out \subseteq in$	No
$in \neq null$	No
$out \neq null$	No

Dinamikus analízis: Daikon

A módszer szempontjából lehetne statikus analízist is használni! Így viszont:

- ▷ Sokkal inkább a program viselkedését vizsgáljuk a szintaxis helyett.
- ▷ A futási idejű property-k alkalmasabbak a helyes és helytelen viselkedés elkülönítésére.

Daikon (dinamikus invariáns detektor):

- ▷ Property-k detekciója specifikus pontokon: pl. eljárások belépési és kilépési pontjai.
- ▷ Detektált property-k pl.:
 - ▷ Egyenlőség, értékvétel adott tartományban.
 - ▷ Nullitás.
 - ▷ Modulus.
 - ▷ Lineáris ill. egyéb függvénykapcsolatok.
 - ▷ Tömbökre: a fentiek elemenként, sorrendezés (lexografikus is), szélsőértékek, stb.
 - ▷ Implikációk, diszjunkciók, pl.: `if p ≠ null then p.value > x`

Feature vektorok előállítás

Property	Equation				Variable type			Number of vars
	\leq	$=$	\neq	\subseteq	int	double	array	
<code>out[1] ≤ in[1]</code>	1	0	0	0	0	1	0	2
<code>∀ i: in[i] ≤ 100</code>	1	0	0	0	0	1	0	1
<code>in[0] = out[0]</code>	0	1	0	0	0	1	0	2
<code>size(out) = size(in)</code>	0	1	0	0	1	0	0	2
<code>in ⊆ out</code>	0	0	0	1	0	0	1	2
<code>out ⊆ in</code>	0	0	0	1	0	0	1	2
<code>in ≠ null</code>	0	0	1	0	0	0	1	1
<code>out ≠ null</code>	0	0	1	0	0	0	1	1

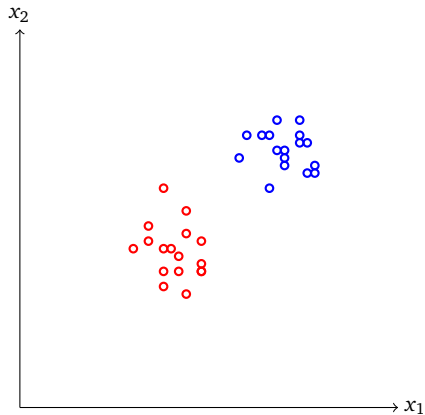
Gyakorlatban: $\mathbf{x}_i \in \mathbb{R}^{388}$, $y_i \in \{-1, +1\}$, $T := \{\mathbf{x}_i\}_i^P$,

$F := \{\text{Hibás program property-jei}\}$, $S := \{\text{Javított program property-jei}\}$.

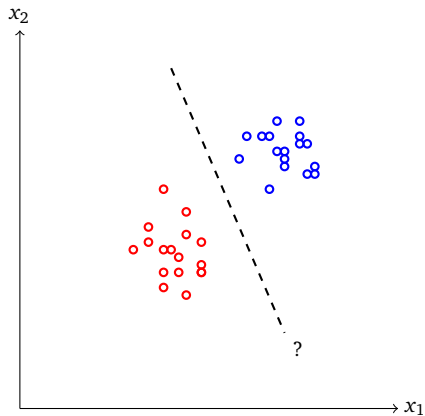
Címkék:

- ▷ $y_i = +1 \Leftrightarrow \mathbf{x}_i \in F, \mathbf{x}_i \notin S$
- ▷ $y_i = -1 \Leftrightarrow \mathbf{x}_i \in S \cup F$
- ▷ $\mathbf{x}_i \notin F \Rightarrow \mathbf{x}_i \notin T$

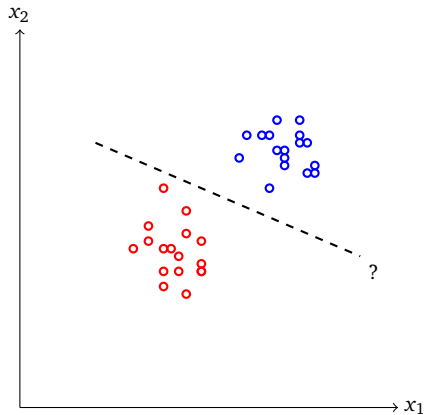
Szupportvektor-gépek



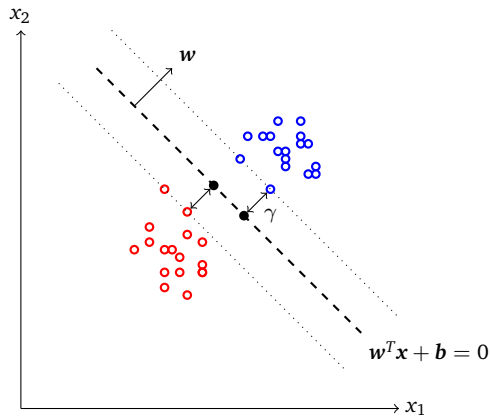
Szupportvektor-gépek



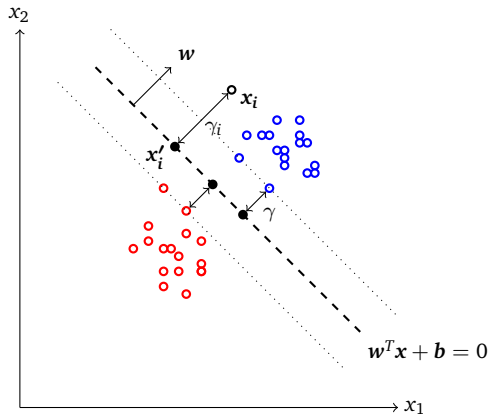
Szupportvektor-gépek



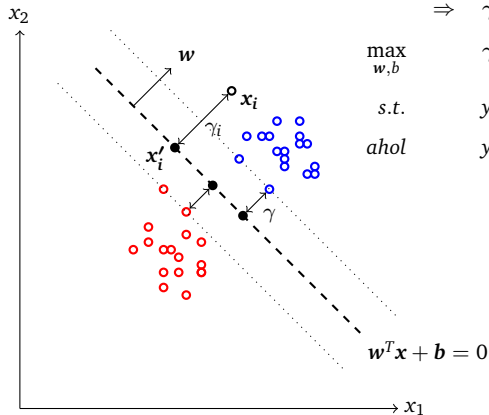
Szupportvektor-gépek



Szupportvektor-gépek



Szupportvektor-gépek



$$w^T x'_i + b = w^T \left(x_i - \gamma_i \frac{w}{\|w\|} \right) + b = 0$$

$$\Rightarrow \gamma_i = \left(\frac{w}{\|w\|} \right)^T x_i + \frac{b}{\|w\|}.$$

max
 w, b

γ

s.t.

$$y_i (w^T x_i + b) \geq \gamma, \quad \|w\| = 1,$$

ahol

$$y_i \in \{-1, +1\}.$$

Szupportvektor-gépek

Ekvivalens feladat:

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t.} \quad & y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1. \end{aligned}$$

„Soft-margin” verzió:

$$\begin{aligned} \min_{w,b,\xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \xi_i \\ \text{s.t.} \quad & y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0. \end{aligned}$$

Bázisfüggvények + „soft-margin”:

$$\begin{aligned} \min_{w,b,\xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \xi_i \\ \text{s.t.} \quad & y_i (\mathbf{w}^T \varphi(\mathbf{x}_i) + b) \geq 1 - \xi_i, \quad \xi_i \geq 0. \end{aligned}$$

Duál feladat

$$\max_{\alpha \geq 0, \beta \geq 0} \min_{\mathbf{w}, b, \xi} \mathcal{L} := \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_i (C - \beta_i) \xi_i - \sum_i \alpha_i \left[y_i \left(\mathbf{w}^T \varphi(\mathbf{x}_i) + b \right) - 1 + \xi_i \right].$$

$$\nabla_{\mathbf{w}} \mathcal{L} = \mathbf{w} - \sum_i \alpha_i y_i \varphi(\mathbf{x}_i) = 0$$

$$\Rightarrow \mathbf{w} = \sum_i \alpha_i y_i \varphi(\mathbf{x}_i).$$

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_i \alpha_i y_i = 0$$

$$\Rightarrow \sum_i \alpha_i y_i = 0.$$

$$\frac{\partial \mathcal{L}}{\partial \xi_i} = C - \beta_i - \alpha_i = 0$$

$$\Rightarrow 0 \leq \alpha_i \leq C.$$

$$\max_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$$

$$\text{s.t.} \quad \sum_i \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C.$$

A duál közelebbről

Mátrixos formában:

$$\begin{aligned} \max_{\boldsymbol{\alpha}} \quad & \mathbf{1}^T \boldsymbol{\alpha} - \frac{1}{2} \boldsymbol{\alpha}^T Q \boldsymbol{\alpha} \\ \text{s.t.} \quad & \mathbf{y}^T \boldsymbol{\alpha} = 0, \quad 0 \leq \alpha_i \leq C \\ \text{where} \quad & Q_{ij} = y_i y_j \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j). \end{aligned}$$

a predikált címke pedig:

$$y(\mathbf{x}) = \text{sgn}(\mathbf{w}^T \varphi(\mathbf{x}) + b) = \text{sgn}\left(\sum_i \alpha_i y_i \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}) + b\right) = \text{sgn}\left(\sum_i \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) + b\right),$$

vagy rangsorolás esetén:

$$\text{score}(\mathbf{x}) = \frac{(\sum_i \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) + b)}{\sqrt{\boldsymbol{\alpha}^T Q \boldsymbol{\alpha}}}.$$

Döntési fák

- ▶ Iteratívan, mohó módon mindig a legnagyobb információ-nyereséggel járó particionálás kiválasztása.
- ▶ Ember számára is könnyen értelmezhető ha-akkor szabályrendszer tanulása.

Entrópia:

$$H(X) = - \sum_i \log(p(x_i))p(x_i).$$

Információ-nyereség: a particionálás által kiváltott várható csökkenés az entrópiában:

$$IG(X, a) = H(X) - \sum_{v \in \text{values}(a)} \frac{|\{\mathbf{x} \in X | x_a = v\}|}{|X|} H(\{\mathbf{x} \in X | x_a = v\}).$$

Mérési környezet

Program	# fgv	Átlagos NCNB	LOC	Hibás verziók
print_tokens	18	452	539	7
print_tokens2	19	379	489	10
replace	21	456	507	32
schedule	18	276	397	9
schedule2	16	280	299	10
space	137	9568	9826	34
tcas	9	136	174	41
tot_info	7	334	398	23
C összesen	245	11881	12629	166
Geo	49	825	1923	95
Pathfinder	18	430	910	41
Streets	19	1720	4459	60
FDAnalysis	277	5770	8864	11
Java összesen	363	7145	16156	207

Eredmények

Mérések:

- ▷ Pozitívnak ítélt minták megítélése javított verziók felhasználásával.
- ▷ Fault-revealing property-k gyakorlati felhasználhatóságának megítélésre emberi szakértők segítségével.

Eredmények:

- ▷ Rangsorolás jobb a klasszifikációnál.
- ▷ C: az első 80 property 45%-a fault-revealing, Java: 59%.
- ▷ A legnagyobb méretű programon (space) működött a legjobban.

Szakértői validáció:

- ▷ 410 pozitívnak ítélt property-t szakértő ellenőrzött: 65% valós programhiba felfedezéséhez vezetett volna.

Eredmények – példa

Procedure	Program	Description of error	Fault-revealing property
addstr	replace	maxString is initialized to 100 but maxPattern is 50	maxPattern > 50
upgrade_process_prio	schedule	prio is incorrectly set to 2 instead of 1	$(prio \geq 2) \Rightarrow return \leq 0$