

Code Clone Detector

ระบบตรวจสอบความเหมือนของโค้ดและตรวจสอบการคัดลอกโค้ด

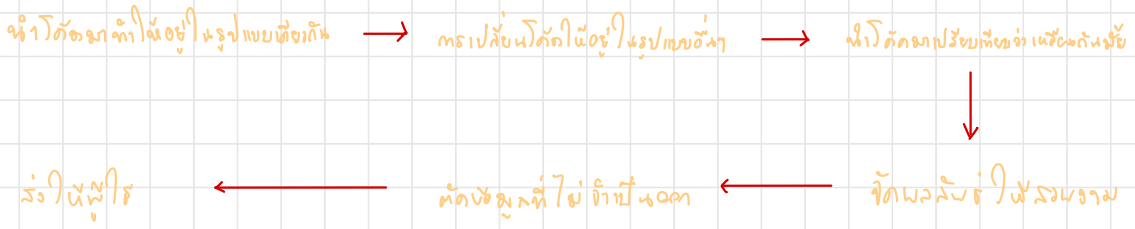
โค้ดที่เหมือนกันกี่เปอร์เซ็นต์?

มีโค้ดสองตัวที่เหมือนกันเพียง 10%

โค้ดที่เหมือนกันหลายบรรทัด

- 1. มีจำนวนเหมือนกัน 100%
- 2. จำนวนเหมือนกันเกือบทั้งหมด แต่มีบางส่วนของมันคือ มีบางส่วนของมัน
- 3. จำนวนที่เหมือนกันที่น้อยกว่า มีจำนวนน้อย โค้ดที่เหมือนกัน 100% หรือ 100%
- 4. ให้ที่โปรแกรมเมอร์ได้รู้ว่ามีกี่เปอร์เซ็นต์

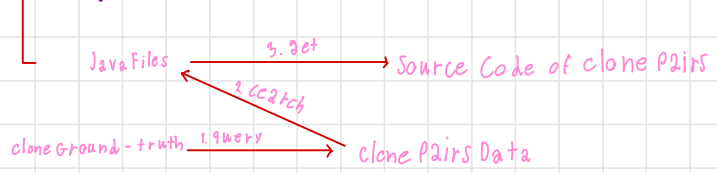
วิธีที่ตรวจสอบว่าโค้ดที่เหมือนกัน



วิธีที่ตรวจสอบ

1. เก็บข้อมูล

2. ใช้ฐานข้อมูล Big Clone Bench
นำข้อมูลมาเปรียบเทียบ ?



แบบ Data ที่	training set	test set
	↓	↓
	แบบเป็น type	แบบเป็น type
จำนวน Data	113,750 33,104 33,120 33,462	12,393 21,557 33,097 33,137

L 2. Metrics

• Syntactic metrics

L ၂။ ဂရမ်မာတစ်ခုကိုသုံးသပ်ခြင်း
 token, uni, identifier, operator...

code 1 → ကိုယ်တိုင်သုံးသပ်ခြင်း
 code 2 → ကိုယ်တိုင်သုံးသပ်ခြင်း

No.	Metric
1	Token No [9]
2	Unique Token No [10]
3	Identifier No [10]
4	Unique Identifier No [10]
5	Operator No [10]
6	Unique Operator No [10]
7	Token Types Diversity [9]
8	Diff File Name Score
9	Diff Method Name Score
10	Similar Return Type
11	DIFFLOC

• Semantic metrics

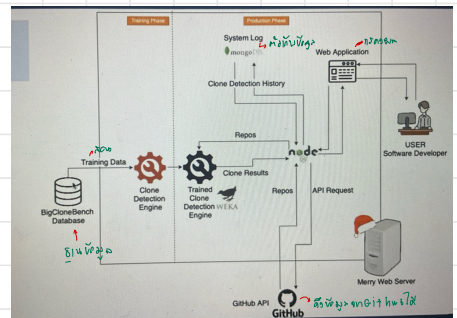
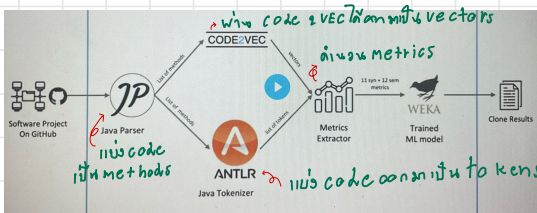
L ၂။ ဂရမ်မာတစ်ခုကိုသုံးသပ်ခြင်း

code 1 → code 2 vec → vector → ချက်ပြုတ်မှုတူခြင်း
 code 2 → code 2 vec → vector → ချက်ပြုတ်မှုတူခြင်း

L 3. ပြု machine Learning models ပြုလုပ်ရန်အတွက်အသုံးပြုသော code 4 မျိုး

1. Decision Tree
2. Random Forest
3. Support Vector Machine
4. SVM using SMO

၂။ မြန်မာနိုင်ငံရှိ ML အသုံးပြုသော code



การวัดประสิทธิภาพ | ชีวิตจริงแบบง่าย

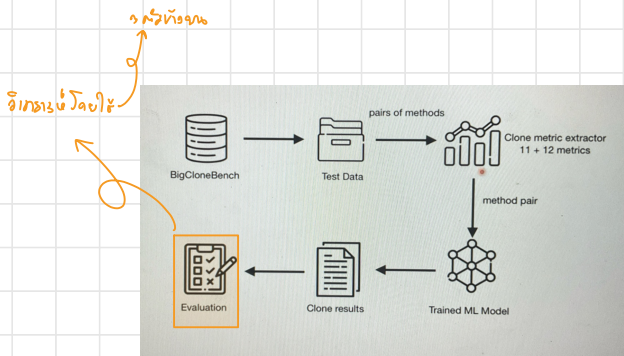
ตัวอย่าง 3 ตัว

$$1. \text{ Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

$$2. \text{ Recall} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

$$3. \text{ F1-Score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

หาค่าเฉลี่ย



ผลการทดลอง

Model	Metrics	Precision	Recall	F1-Score
Randomization (baseline)		0.20	0.49	0.28
Decision Tree	Syntactic + Semantic	0.89	0.86	0.87
	Syntactic	0.95	0.72	0.86
	Semantic	0.68	0.87	0.76
Random Forest	Syntactic + Semantic	0.97	0.86	0.91
	Syntactic	0.97	0.80	0.87
	Semantic	0.70	0.87	0.78
SVM	Syntactic + Semantic	0.97	0.85	0.91
	Syntactic	0.97	0.79	0.87
	Semantic	0.62	0.90	0.73
SVM using SMO	Syntactic + Semantic	0.98	0.89	0.93
	Syntactic	0.97	0.69	0.81
	Semantic	0.63	0.90	0.74

→ ผลลัพธ์ที่ดีที่สุด

ใช้วิธีนี้