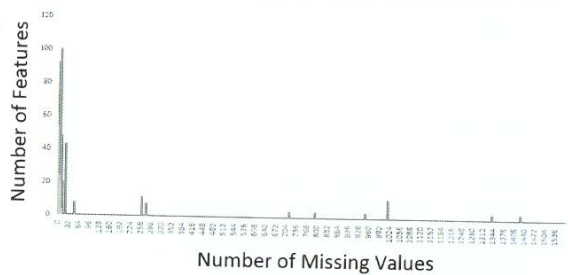
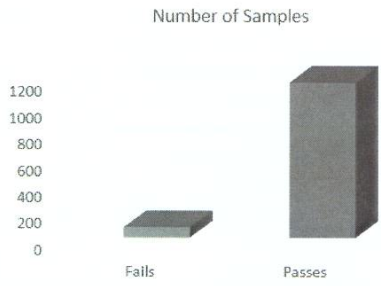


실험제목	센서 기반 제조 공정을 위한 인공지능	
실험목적	반도체 제조 공정의 센서 데이터를 이용한 생산 공정 관리	
<p>배경</p> <p>경쟁이 치열하고, 제품의 변화가 빠르기 때문에 생산 공정의 유연한 조정과 제품 출시 시간 단축은 매우 중요하다.</p> <p>IoT 센서의 발전으로 제조 공정에서 다양한 데이터 수집이 가능해졌으며, 그 데이터를 이용한 공정 관리는 복잡해지고 있음.</p> <p>제조 공정 관리를 위해서 빅데이터와 인공지능 기술의 결합을 통해 유연하고 빠른 공정 관리가 가능할 것으로 예상됨.</p>		
<p>반도체 제조 공정 데이터 (SECOM) 을 이용하여 생산 공정 관리를 위한 기계 학습 기술을 적용할 방법 연구.</p>		
<p>실험 설계</p> <p>1. 데이터 분석을 위한 알고리즘 선택</p> <ul style="list-style-type: none"> <li>- Boruta</li> <li>- Multivariate Adaptive Regression Spline (MARS)</li> <li>- PCA</li> </ul> <p>2. 분류 알고리즘</p> <ul style="list-style-type: none"> <li>- Logistic Regression (LR)</li> <li>- Random Forest (RF)</li> <li>- Gradient Boosted Trees (GBT)</li> </ul> <p>3. 평가</p> <ul style="list-style-type: none"> <li>- FPR</li> <li>- Precision</li> <li>- recall</li> <li>- F-measure</li> <li>- accuracy</li> </ul>		
기록자 Written by 김민태	점검자 Witnessed or Understood by	점검자 Witnessed or Understood by
일자 1/1	일자	일자

실험제목	센서 기반 제조 공정을 위한 인공지능	
실험목적	반도체 제조 공정의 센서 데이터를 이용한 생성 공정 관리	
데이터 분석.		
- 100개 이상의 단계로 구성된 반도체 제조 공정에서 수집된 데이터		
- 1567개의 샘플 데이터		
- 샘플 데이터는 591개의 특성을 가짐.		
- 측정이 되지 않은 값이 'NaN'으로 기록		
- 제품에 문제가 있는 경우 1, 정상인 경우 -1로 나타냄		
실험 설계		
<pre> graph LR     A[데이터 전처리] --&gt; B[feature 추출]     B --&gt; C[교차 검증]     C --&gt; D[학습 데이터 추출]     D --&gt; E[샘플링]     E --&gt; F[알고리즘 선택]     F --&gt; G[평가]     C --&gt; H[테스트 데이터 추출]     H --&gt; G         </pre>		
1. 데이터 전처리	4. 샘플링	
- 결측치 제거	- 불균형 데이터	
- 데이터 [A, B] 정규화	- 오버 샘플링	
2. feature 추출	5. 알고리즘 선택	
- Bruta	- RF	
- MARS	- LR	
- PCA	- GBT	
3. 교차검증		
- k-Fold cross validation		
기록자 Written by	점검자 Witnessed or Understood by	점검자 Witnessed or Understood by
김민태		
일자	일자	일자
1 / 8		

실험제목	센서 기반 제조 공정을 위한 인공지능																
실험목적	반도체 제조 공정의 센서 데이터를 이용한 생산 공정 관리																
데이터 전처리																	
-55% 이상 데이터가 결측된 경우 해당 특성 제거																	
		<ul style="list-style-type: none"><li>- 각 샘플 데이터에서 결측값 수를 시각화</li><li>- 591개의 특성중 4개가 제거됨.</li><li>- 567개의 특성을 고려하여 사용</li></ul>															
Fig. 2. Features Distribution by the Number of Missing Values																	
<ul style="list-style-type: none"><li>- 'NaN' 값은 전, 후 값이 정상치일 경우 전, 후 값의 평균 이용</li><li>- 전, 후 값 중 'NaN'이 존재할 경우 정상치를 그대로 이용</li><li>- 상수 값을 제거 하여 총 116개의 특성이 추가로 제거됨</li><li>- 451개의 특성이 남으며, 다음과 같은 식을 통해 0~1로 정규화</li></ul>																	
$x_{norm} = \frac{x - x_{min}}{x_{max} - x_{min}}$																	
feature 추출																	
-분류에 영향을 미치지 않는 특징 제거																	
<table><tr><th colspan="3">TABLE II COMPARISON OF THE FEATURES SELECTION TECHNIQUES</th></tr><tr><th>Features Selection Algorithm</th><th>Number of Selected Features</th><th>Characteristics</th></tr><tr><td>Boruta</td><td>22</td><td>by default it uses the Random Forest (RF) algorithm</td></tr><tr><td>MARS</td><td>10</td><td>it is a form of regression analysis</td></tr><tr><td>PCA</td><td>111</td><td>it uses orthogonal linear transformations</td></tr></table>		TABLE II COMPARISON OF THE FEATURES SELECTION TECHNIQUES			Features Selection Algorithm	Number of Selected Features	Characteristics	Boruta	22	by default it uses the Random Forest (RF) algorithm	MARS	10	it is a form of regression analysis	PCA	111	it uses orthogonal linear transformations	<ul style="list-style-type: none"><li>- Boruta는 76번의 반복 수행을 통해 22개의 특성 추출</li><li>- MARS는 RSS와 GCV를 기반으로 10개의 특성을 추출</li><li>- PCA는 111개의 특성을 추출</li></ul>
TABLE II COMPARISON OF THE FEATURES SELECTION TECHNIQUES																	
Features Selection Algorithm	Number of Selected Features	Characteristics															
Boruta	22	by default it uses the Random Forest (RF) algorithm															
MARS	10	it is a form of regression analysis															
PCA	111	it uses orthogonal linear transformations															
기록자 Written by 김민태																	
점검자 Witnessed or Understood by																	
일자 1/15																	
점검자 Witnessed or Understood by																	
일자																	

실험제목	센서 기반의 제조 공정을 위한 인공지능	
실험목적	반도체 제조 공정의 센서 데이터를 이용한 생산 공정 관리	
샘플링		
- 소수 클래스와 다수 클래스의 불균형 문제		
- 1567개의 데이터 중 104개는 실패, 1463개는 성공 데이터		
- 1:14의 비율로 불균형 데이터임.		
<div><div></div><div><div>- 불균형을 해결하기 위해 다수 클래스에서 무작위로 데이터를 뽑는 언더 샘플링 이용</div><div>- 소수 클래스를 복사하는 오버 샘플링 이용</div></div></div>		
Fig. 6. Initial Distribution of the Fail Samples and the Pass Samples		
- 언더 샘플링에는 Spread Sub Sample 필터를 이용		
- 오버 샘플링에는 Synthetic Minority Oversampling Technique 이용		
데이터 분할		
- 학습 데이터와 테스트 데이터를 4:1 (80%:20%)로 분할		
- 모델 검증을 위해 교차 검증을 이용		
- k-fold cross validation 이용하며, 5개로 분할하는 5-fold cross validation 이용		
- 모델 평가에는 F-measure, recall, Precision, Accuracy, False Positive Rate (FPR)을 이용		
- 테스트 데이터에서 과적합 문제를 검증하기 위해 샘플링은 적용하지 않음.		
기록자 Written by	점검자 Witnessed or Understood by	점검자 Witnessed or Understood by
김민태		
일자	일자	일자
1 / 22		



실험제목	센서 기반 제조 공정을 위한 인공지능																						
실험목적	반도체 제조 공정의 센서 데이터를 이용한 생산 공정 관리																						
분류 알고리즘																							
- LR은 Broyden - Fletcher - Goldfarb - Shanno (LBFGS) 알고리즘 이용																							
- Random Forest는 다음과 같이 구성																							
<p>TABLE III RANDOM FOREST PARAMETERS</p> <table border="1"> <thead> <tr> <th>Parameter</th> <th>Value</th> <th>Description</th> </tr> </thead> <tbody> <tr> <td>number of classes</td> <td>2</td> <td>the number of classification classes</td> </tr> <tr> <td>number of trees</td> <td>20</td> <td>a parameter that describes how many trees are in the forest</td> </tr> <tr> <td>features subset strategy</td> <td>"auto"</td> <td>how many features to use as candidates for splitting</td> </tr> <tr> <td>impurity</td> <td>"gini"</td> <td>a parameter that describes the homogeneity of the labels</td> </tr> <tr> <td>maximum depth</td> <td>20</td> <td>the maximum depth of the trees from the forest</td> </tr> <tr> <td>maximum number of bins</td> <td>100</td> <td>the maximum number of bins that are used to discretize the features</td> </tr> </tbody> </table>			Parameter	Value	Description	number of classes	2	the number of classification classes	number of trees	20	a parameter that describes how many trees are in the forest	features subset strategy	"auto"	how many features to use as candidates for splitting	impurity	"gini"	a parameter that describes the homogeneity of the labels	maximum depth	20	the maximum depth of the trees from the forest	maximum number of bins	100	the maximum number of bins that are used to discretize the features
Parameter	Value	Description																					
number of classes	2	the number of classification classes																					
number of trees	20	a parameter that describes how many trees are in the forest																					
features subset strategy	"auto"	how many features to use as candidates for splitting																					
impurity	"gini"	a parameter that describes the homogeneity of the labels																					
maximum depth	20	the maximum depth of the trees from the forest																					
maximum number of bins	100	the maximum number of bins that are used to discretize the features																					
- GBT는 다음과 같이 구성																							
<p>TABLE IV GRADIENT BOOSTED TREES PARAMETERS</p> <table border="1"> <thead> <tr> <th>Parameter</th> <th>Value</th> <th>Description</th> </tr> </thead> <tbody> <tr> <td>number of classes</td> <td>2</td> <td>the number of classification classes</td> </tr> <tr> <td>number of iterations</td> <td>3</td> <td>each iteration creates a tree</td> </tr> <tr> <td>maximum depth</td> <td>5</td> <td>the maximum depth of the trees</td> </tr> </tbody> </table>			Parameter	Value	Description	number of classes	2	the number of classification classes	number of iterations	3	each iteration creates a tree	maximum depth	5	the maximum depth of the trees									
Parameter	Value	Description																					
number of classes	2	the number of classification classes																					
number of iterations	3	each iteration creates a tree																					
maximum depth	5	the maximum depth of the trees																					
샘플링 방법																							
<p>TABLE VII SAMPLING METHOD</p> <table border="1"> <thead> <tr> <th>Sampling Method</th> <th>Description</th> </tr> </thead> <tbody> <tr> <td>Unsampled</td> <td>the number of samples is unchanged</td> </tr> <tr> <td>Undersampled Majority Class</td> <td>the number of samples of the majority class is reduced to 78</td> </tr> <tr> <td>Oversampled Minority Class</td> <td>the number of samples of the minority class is increased to 1176</td> </tr> </tbody> </table>			Sampling Method	Description	Unsampled	the number of samples is unchanged	Undersampled Majority Class	the number of samples of the majority class is reduced to 78	Oversampled Minority Class	the number of samples of the minority class is increased to 1176													
Sampling Method	Description																						
Unsampled	the number of samples is unchanged																						
Undersampled Majority Class	the number of samples of the majority class is reduced to 78																						
Oversampled Minority Class	the number of samples of the minority class is increased to 1176																						
기록자 Written by	점검자 Witnessed or Understood by	점검자 Witnessed or Understood by																					
김민태																							
일자	일자	일자																					
1/29																							

실험제목	센서 기반 제조 공정을 위한 인공지능
실험목적	반도체 제조 공정의 센서 데이터를 이용한 생산 공정 관리

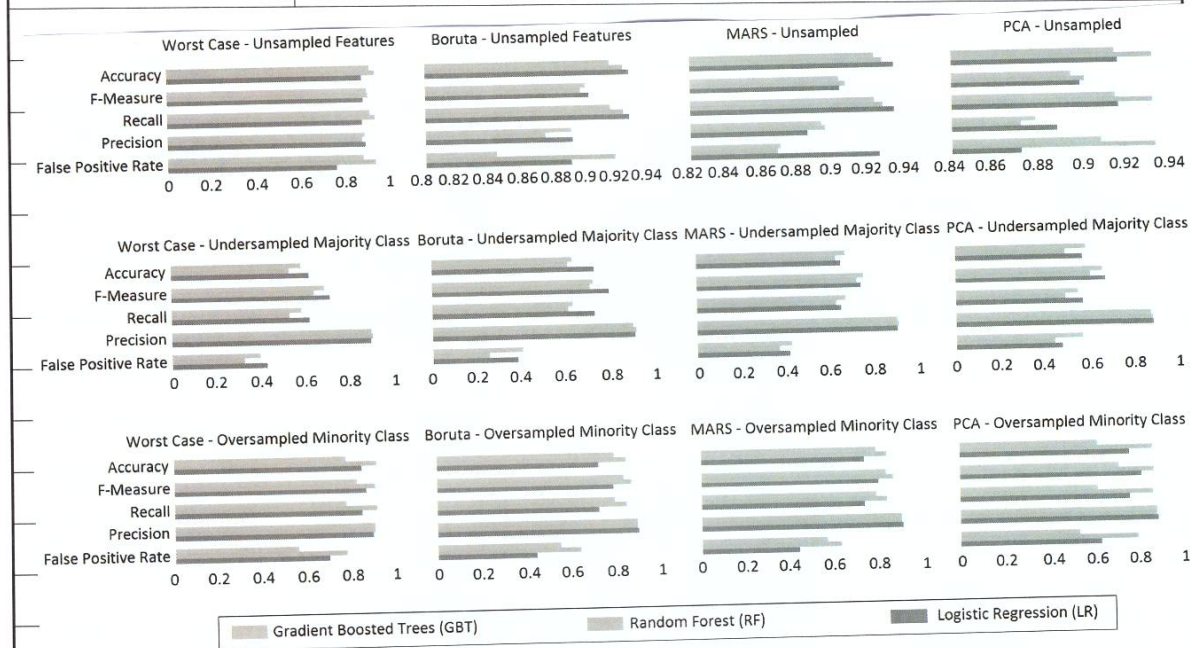


Fig. 8. Experimental Results Summary - Classification Algorithm Perspective

## 실험 결과

- 데이터를 샘플링 하지 않은 경우 Boruta 알고리즘을 통해 특성을 추출 한 뒤 LR을 적용 하였을 때 정밀도가 0.89440으로 가장 높게 나타남
- 다수 클래스 언더 샘플링 시 Boruta 알고리즘과 RF 결합시 정밀도가 0.90506으로 가장 높게 나타났음
- 소수 클래스 오버 샘플링 시 MARS 알고리즘과 LR 이용시 0.91541로 가장 높게 나타났음
- FPR을 기준으로 평가 시에도 정밀도와 동일한 결과가 나타남

ㄷ

## 결론

- 대부분의 경우 Boruta 알고리즘과 RF를 이용한 경우 높은 성능이 나타났음

기록자 Written by 김민태	점검자 Witnessed or Understood by	점검자 Witnessed or Understood by
일자 1/29	일자	일자