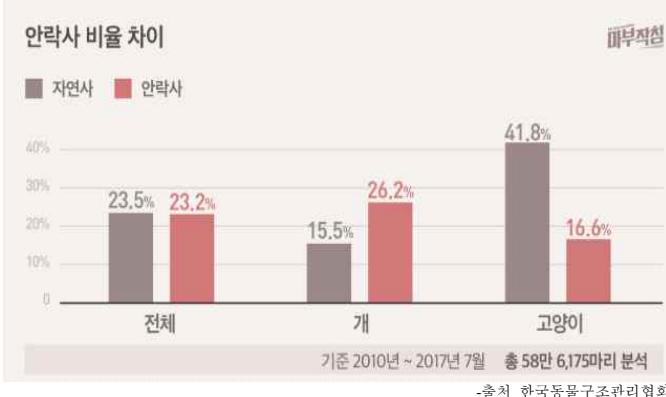
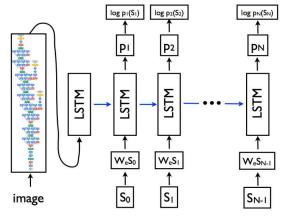


## 캡스톤 디자인 결과보고서

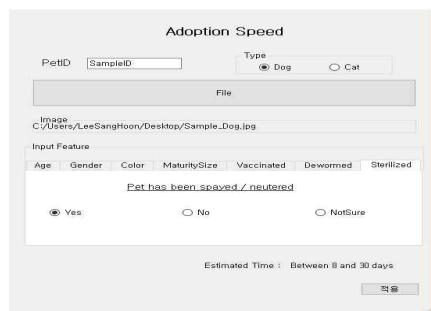
주제		유기동물 입양 기간 예측																																																		
팀명		팀원명 (총 3명)	담당교수																																																	
집•돌		이상훈(팀장), 김윤하, 안중현	황두성 교수님																																																	
과제 개요	설계 과제 요약/목표	<p>Deep Learning 과 Machine Learning 을 통해 메타데이터 및 이미지를 학습하여 동물 형태에 따른 입양기간을 예측한다. 그로 인해 유기동물보호소에서 학습된 모델을 이용하여 유기동물의 입양기간을 예측하고 입양기간이 길게 예측된 동물에 홍보를 집중하는 차등적 홍보등의 방법을 통한 유기동물의 입양이 순환이 잘 되도록 하여 더 많은 동물이 입양되어 안락사를 줄이고자 한다.</p>																																																		
현황 분석	기준 방법의 특징/한 계	<ul style="list-style-type: none"> <li>수백만의 길 잃은 애완동물들이 유기되고 보호소로 보내진다. 유기동물들은 입양자를 기다리다 보호소 수용 최대가 되면 안락사를 당하게 된다.</li> </ul>  <p><b>유기견보호센터 유기동물 보호중입니다!</b> 동물보호관리시스템 유기동물 공고!</p> <p>유기동물을 보호 중입니다! 등록하기</p> <p>실증동물 선택 [로드 보호동물] [포스트보]</p> <table border="1"> <thead> <tr> <th>사진</th> <th>동물종류</th> <th>상세설명</th> <th>발견장소</th> <th>발견날짜</th> <th> 조회</th> </tr> </thead> <tbody> <tr> <td></td> <td>인천 구월동 스코티어풀드 보호중...</td> <td>새벽3시쯤 구월동 만의점 앞에서 어린 고양이가 혼자 암아있...</td> <td>인천 남동구 구월동</td> <td>05-26</td> <td>15</td> </tr> <tr> <td></td> <td>끼양 닉스풀트 어마</td> <td>온라인 커뮤니티에서 실근자를 배회하는 아이입니다.</td> <td>경기 파주시 마창초수 근처</td> <td>05-26</td> <td>20</td> </tr> <tr> <td></td> <td>강북구 미스건 주인 찾습니다.</td> <td>강북노인종합복지관 앞 골목에서 자 릴의 숨 어있는 것을 발견...</td> <td>서울 강북노인종합복지관</td> <td>05-26</td> <td>15</td> </tr> <tr> <td></td> <td>말티즈 애미 보호중</td> <td>말티즈 애미 2살가량 내곡동 주민센터 옥고</td> <td>서울 서초구 내곡동</td> <td>05-23</td> <td>21</td> </tr> <tr> <td></td> <td>자와와 일보중입니다</td> <td>제주대구역 세종로와본길 #포인핸드 #대구강 대구 동구 현대백화점 아출...</td> <td>05-25</td> <td>21</td> <td></td> </tr> </tbody> </table> <p><b>안락사 비율 차이</b></p>  <table border="1"> <thead> <tr> <th>Category</th> <th>자연사 (%)</th> <th>안락사 (%)</th> </tr> </thead> <tbody> <tr> <td>전체</td> <td>23.5%</td> <td>23.2%</td> </tr> <tr> <td>개</td> <td>15.5%</td> <td>26.2%</td> </tr> <tr> <td>고양이</td> <td>41.8%</td> <td>16.6%</td> </tr> </tbody> </table> <p>기준 2010년 ~ 2017년 7월 총 58만 6,175마리 분석</p> <p>-출처 한국동물구조관리협회</p>			사진	동물종류	상세설명	발견장소	발견날짜	조회		인천 구월동 스코티어풀드 보호중...	새벽3시쯤 구월동 만의점 앞에서 어린 고양이가 혼자 암아있...	인천 남동구 구월동	05-26	15		끼양 닉스풀트 어마	온라인 커뮤니티에서 실근자를 배회하는 아이입니다.	경기 파주시 마창초수 근처	05-26	20		강북구 미스건 주인 찾습니다.	강북노인종합복지관 앞 골목에서 자 릴의 숨 어있는 것을 발견...	서울 강북노인종합복지관	05-26	15		말티즈 애미 보호중	말티즈 애미 2살가량 내곡동 주민센터 옥고	서울 서초구 내곡동	05-23	21		자와와 일보중입니다	제주대구역 세종로와본길 #포인핸드 #대구강 대구 동구 현대백화점 아출...	05-25	21		Category	자연사 (%)	안락사 (%)	전체	23.5%	23.2%	개	15.5%	26.2%	고양이	41.8%	16.6%
사진	동물종류	상세설명	발견장소	발견날짜	조회																																															
	인천 구월동 스코티어풀드 보호중...	새벽3시쯤 구월동 만의점 앞에서 어린 고양이가 혼자 암아있...	인천 남동구 구월동	05-26	15																																															
	끼양 닉스풀트 어마	온라인 커뮤니티에서 실근자를 배회하는 아이입니다.	경기 파주시 마창초수 근처	05-26	20																																															
	강북구 미스건 주인 찾습니다.	강북노인종합복지관 앞 골목에서 자 릴의 숨 어있는 것을 발견...	서울 강북노인종합복지관	05-26	15																																															
	말티즈 애미 보호중	말티즈 애미 2살가량 내곡동 주민센터 옥고	서울 서초구 내곡동	05-23	21																																															
	자와와 일보중입니다	제주대구역 세종로와본길 #포인핸드 #대구강 대구 동구 현대백화점 아출...	05-25	21																																																
Category	자연사 (%)	안락사 (%)																																																		
전체	23.5%	23.2%																																																		
개	15.5%	26.2%																																																		
고양이	41.8%	16.6%																																																		

과제 내용	개발 프로세스 구현환경 진행일정	인도적 처리(안락사) 대상 동물 순위				비부직침
		동물보호센터 운영지침 20조 2항				
1순위		전염성과 치사율이 높은 질환에 감염되고 상해로 건강회복이 불가능할 것으로 판단되는 개체				
2순위		치료비용이나 치료기간 등 고려시 추가적인 보호가 불가능할 것으로 판단되는 개체				
3순위		심장질환, 백내장, 호르몬 질환 등에 감염돼 분양 후에도 지속적인 치료가 필요한 개체				
4순위		교정이 어려운 행동 장애 등으로 분양이 어려울 것으로 판단되는 개체				
5순위		센터 수용능력, 분양가능성 등을 고려해 보호 및 관리가 어려울 것으로 판단되는 개체				
		실제로 건강상에 문제가 있어서 동물들을 안락사하는 경우도 많지만 유기동물보호센터의 수용능력에 문제와 효율적인 분양시스템이 이루어지지 않아서 안락사하는 경우도 많다.				
		<ul style="list-style-type: none"> <li>Tensorflow, Keras : image 와 metadata 를 분석 및 학습하여 유기동물들의 입양기간을 예측하는 모델 제작</li> <li>PyQT : 학습이 완료된 머신러닝 모델을 적용하여 유기동물보호소에서 사용할 수 있는 간단한 GUI 모델 제작</li> <li>Pycharm, Anaconda : Python 개발환경, 머신러닝 개발환경</li> <li>Colab : 텍스트, 코드 출력을 하나의 공동작업 문서로 통합해 주는 머신러닝 교육과 연구를 위한 데이터 분석 도구</li> <li>Deep learning &amp; Machine learning: 이미지 학습을 위한 imageNet 과 데이터 분류를 위한 Randomforest, K-NN, Logistic Regression 활용</li> </ul>				
세부내용		수행기간(월)				비고
		3	4	5	1 2 3 4 1 2 3 4 1 2 3 4	
아이디어 기획 및 설계						
스터디 및 모델 탐색 데이터 분석						cs231/tensor
데이터 프로세싱						image, meta
모델 제작						vggNet, RF
Model Traning&Testing						
보고서 작성 및 제출						

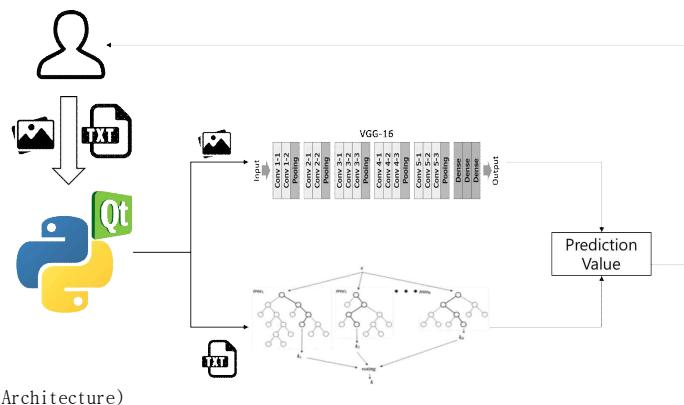
내용	<p>1. Data Processing</p> <ul style="list-style-type: none"> <li>Metadata 중에서 학습에 방해가 되는 feature들을 제외하였고 Image data 중에서는 동물 식별이 힘든 image를 제거하였다.</li> </ul> <p>2. Data Training</p> <ul style="list-style-type: none"> <li>2-1 Metadata Machine Learning           <ul style="list-style-type: none"> <li>-KNN Algorithm</li> <li>-Random Forest Algorithm</li> <li>-Logistic Regression Algorithm</li> </ul>           3가지의 알고리즘을 이용하여 입양기간 정확도가 가장 높은 Model을 선정         </li> <li>2-2 Image Classification(Deep learning)           <ul style="list-style-type: none"> <li>CNN model(VGG, Resnet)을 이용해 image를 학습시켜 입양기간을 예측</li> </ul> </li> </ul> <p>3. GUI Program을 만들어 사용자가 data를 입력하여 모델을 통해서 예측된 입양기간을 출력해준다.</p>
내용	<p>개발은 Machine learning 부분에서는 Random forest와 K-NN 모델들의 구현이 끝나고 Matadata들을 학습을 시켜보았고 Deeplearning CNN모델 중 Resnet을 구현하였고 VGG는 구현 중 생긴 오류를 수정하고 있었다.</p> <p>완성된 모델들에 데이터를 학습시켜 보았는데 학습을 통한 예측의 정확도가 유의미한 결과가 나오지 못하였다. 추후 VGG와 logistic regression까지 완성이 된 후에 정확도를 측정하였지만 마찬가지로 높은 정확도를 보여주지는 못하였다. 데이터 전처리에 충분한 시간을 투자하였지만 더 이상 정확도가 올라 가지는 않았다. 그래서 정확도가 낮은 다른 이유를 찾고자 하였다. 우선 ML 모델은 학습에서 오버피팅이 심하게 일어나 Depth를 줄이고 하이퍼파라미터를 계속 조절하며 데이터 불균형을 줄이고자 하였다. 그 방법으로 오버 샘플링 방법을 찾았고 그 중 SMOTE 기법을 사용하였다. SMOTE 기법을 사용하여 ML 모델들의 정확도를 많이 끌어 올렸으며 특히 Random forest가 0.39 <math>\rightarrow</math> 0.49로 ML 모델들 중에 그나마 가장 높은 정확도를 보였다. DL CNN 모델에서 또한 정확도를 높이고자 하는 노력을 많이 하였는데, 그 중 개와 고양이를 따로 분리하여 학습을 시키는 것은 큰 효과가 없었다. 오히려 Activation이나 데이터 프로세싱 기법들과 같은 기본적인 변화들을 주었더니 VGG가 20% <math>\rightarrow</math> 37% 큰 폭으로 상승하였다.</p> <p>분명 정확도는 많이 상승시켰지만 정형화되지 않은 데이터를 이용해서 인지</p>

	<p>실제로 사용하기는 힘든 유의미하지 않은 결과가 나왔다. 그렇기에 다른 방식을 찾아보았고 LSTM 기법을 도입해 보았다.</p>  <p>CNN 모델에서 사용하는 이미지 데이터와 ML에서 사용하던 Metadata를 염이 Many-to-one 형식의 CNN-LSTM 구조를 만들었다. 이미지를 CNN에 통과시켜서 LSTM의 초기 데이터로 사용할 수 있도록 데이터를 편 다음 LSTM을 통화시키며 Metadata의 각 feature들을 적용하는 방식을 취하였으나 오히려 VGG와 Randomforest보다 정확도가 많이 낮게 나왔다. 그래서 다음 방법으로 이미지만 LSTM을 통과시키는 방법과 Meta data만 LSTM을 통과시키는 방법을 취해 보았지만 전과 큰 차이가 없어 결국 LSTM 도입은 실패하였다.</p> <p>과제수행 기간이 막바지에 다다라 다른 방법을 찾고자 하다 단순하게 DL과 ML의 결과를 염는 방식을 택하였다. 데이터들을 사용자가 입력할 시에 DL과 ML이 각각 결과를 낸 후 그 결과치를 종합하는 방식이었다. 우선 각 모델 선택은 DL과 ML에서 정확도가 가장 높은 모델로 Randomforest와 VGG를 선택하였다. 입양기간의 예측을 범위 단위로 진행하게 되는데 (0~7일 내 입양은 0, 7~21일은 1 등...) 두 모델이 입양기간을 다르게 예측할 경우 어떤 방식으로 범위를 선택해야 할지를 선택해야 했다. 그런데 입양기간 데이터 특성상 숫자가 클수록 점점 넓은 범위를 갖는 특성을 가지고 있어서 두 모델의 예측 기간 중 Max값을 가져오는 방법을 취하였다. 그리고 최종 정확도는 52%를 기록하였다.</p> <p>처음 캡스톤의 결과물로 단순하게 Accuracy를 보여주는 형태를 취하려고 하였다. 그러나 멘토와의 만남에서 결과물이란 Visualize가 되어야 한다는 말씀을 해 주셨고 단순하게 숫자만 보여주는 것이 아닌 사람들이 보고 반응할 수 있는 결과물이 어떤 것이 있을 것인가를 같이 고민해 주셨다. 입양기간에 따라 우선순위가 자동으로 바뀌는 웹 어플리케이션 등 다양한 아이디어가 나왔으나 2차 멘토링을 진행하는 시기가 결과물 제출일까지 많은 시간이 남지 않았고 그 때는 아직 정확도를 조정하는 단계여서 많은 시간을 투자할 수 없</p>
--	--

있기 때문에 PythonQT로 GUI를 만들어 사용자가 실습을 하고자 할 때 직접 Image와 Metadata 입력 시에 예측한 입양기간을 보여주는 Application을 만들었다.



Process를 도식화하여 하단과 같이 나타내었다.

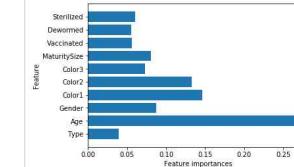


ML과 DL 모델 결과 각각의 예측 결과를 단순 Max 취하는 것은 정확도가 오르더라도 오히려 설득력이 떨어지기에 실제로 쓰기에는 무리가 있다는 것을 알고 있다. 하지만 50%도 안되는 정확도로 프로젝트를 마무리하기보다 다양한 모델들을 시도해 보았으며 정확도를 올리기 위해 다양한 방법을 사용해 보았다라는 메시지를 전달하고 싶었다. 마지막 발표 때에 교수님께서 데이터 자체가 문제일 수도 있다라는 말씀을 해 주셨다. 캡스톤 기간 동안은 유의미한 결과가 나오지 못했지만 데이터부터 다시 확인하며 유의미한 결과를 만들어내고자 한다.

## ● Machine Learning model

■ 정확도: Random Forest > Logistic Regression > KNN

총련 세트 정확도 : 0.715  
테스트 세트 정확도 : 0.361  
특성 중요도 :  
[0.03971239 0.26862534 0.08711554 0.14670956 0.13232278 0.07287098  
0.08042326 0.05617344 0.05521242 0.0601444]



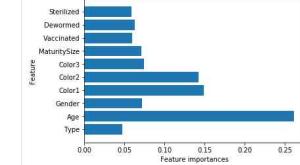
<Randomforest>

하지만 현실 데이터인 경우 불규칙 데이터가 많아 정확도가 낮게 나온다.

따라서 SMOTE(Synthetic Minority Over-Sampling Technique)라는 Over-Sampling 기법을 사용하여 정확도를 올릴 수 있었다.

- SMOTE 사용 후 36% -> 48%

총련 세트 정확도 : 0.763  
테스트 세트 정확도 : 0.480  
특성 중요도 :  
[0.04713301 0.26149447 0.07214365 0.14899595 0.14264649 0.07415923  
0.07101826 0.05997955 0.06324295 0.05917196]

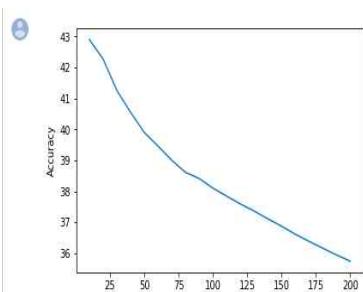


- 다른 Machine Learning Algorithm(K-NN, Logistic Regression)들도 SMOTE 기법 적용하여 다시 결과를 재비교

testing dataset: 44,761/90476/90476%

Confusion matrix						
		True label		Predicted label		
		0	1	2	3	4
0	361	83	21	46	79	350
1	174	141	78	89	118	350
2	133	113	99	100	121	350
3	173	75	51	135	125	350
4	141	92	57	95	193	350

<Logistic Regression>

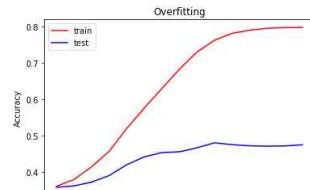


<K-NN>

## 결과 및 문제점

### ■ Overfitting

SMOTE를 통해서 Randomforest의 정확도가 올라갔지만 여전히 overfitting의 문제가 있었다.



(Random Forest의 overfitting)

따라서 parameter의 값을 변경해본 결과 트리의 깊이 때문에 overfitting이 이루어짐을 알 수 있었다. 트리의 깊이를 얕게 하면 overfitting을 해결할 수 있었지만 동시에 정확도도 줄어들기에 parameter를 비교하여 최적의 parameter의 값을 설정하였다.

### • Deep Learning model

#### VGG16 Model Architecture

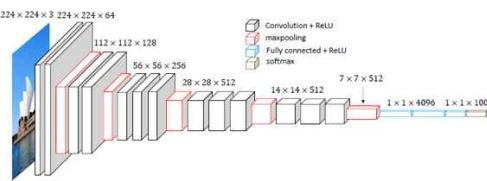
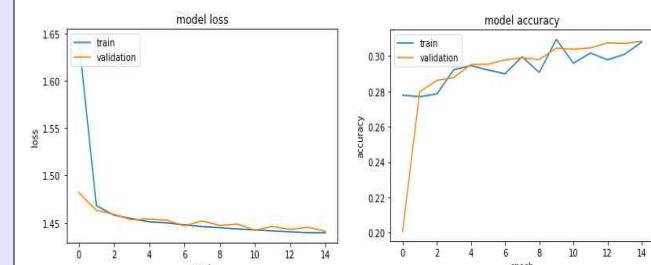


Figure 2: The architecture of VGG16 model .

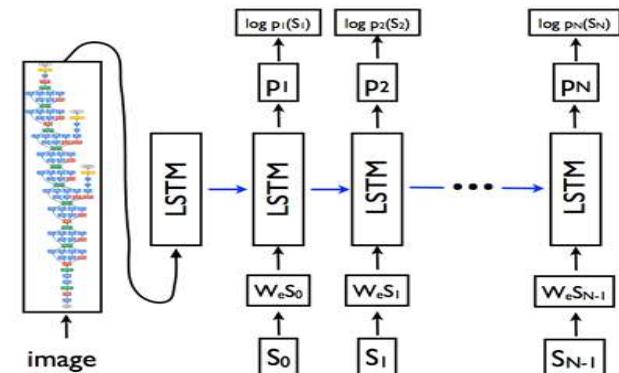
VGGNet의 모델을 선택한 이유는 이해가 쉽고 변형을 시켜가면서 테스트하기에 용이하다는 것 이었다. 일단 layer층이 현재 다른 ResNet이나 GoogelNet들에 비해서 얕았기 때문에 단시간에 이해하기가 쉬웠던 것 같다.

뿐만 아니라 VGGNet의 가장 큰 특징은 작은 filter 크기의 convolution 연산이기 때문에 그 만큼 parameter의 개수가 많아져서 연산량이 증가하여 학습시간이 오래 걸렸지만 높은 정확도를 얻을 수 있었다.

### ■ Image Training 결과



### • LSTM&CNN



정확도 문제를 해결하기 위해서 CNN과 LSTM을 이용하여 이미지 데이터와 Metadata를 함께 사용하고자 하였고 이미지를 CNN을 통과시킨 후 LSTM을 이용하여 Many-one 예측 구조를 이용하려고 하였다. 하지만 오히려 다른 ML,DL 모델보다 성능이 더 떨어졌다.

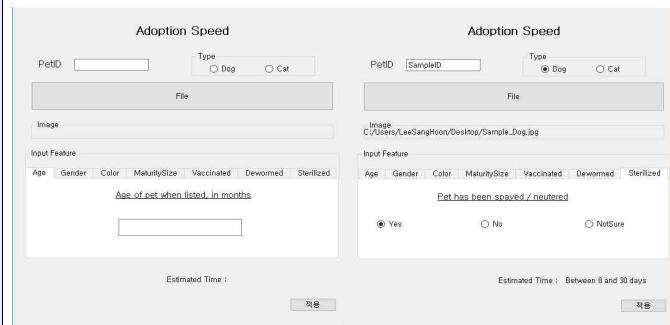
### • 결과 예측

각 모델들의 예측 값들과 실제 값이, 예측 값의 평균, 작은 값, 큰 값 중 어떤 값을 선택했을 때 정확도가 향상되는지 비교하였다.

	<b>Prevoius acc</b>	<b>Last acc</b>
VGG	20%	37%
Randomforest	36%	49%
<b>VGG+Randomforest</b>	-	<b>52%</b>

다른 두 방법은 떨어지거나 변화가 없는 데에 비해 높은 값으로 예측한 값을 선택하는 방법을 택했을 때에 정확도가 상승하였다. 그래서 VGG와 Randomforest의 각 결과값 중 큰 값을 기준으로 택하는 방식을 취하였고 정답으로 예측한 모델 선정하여 최종 정확도는 52%를 기록하였다.

- GUI



QT designer Tool을 이용하여 Window의 틀을 만들고 PYQT5 Module을 이용, 각 button, text, label에 대한 event 처리하였고 model을 가져와 사용자들이 입력한 값을 적용하여 입양기간의 예측값을 출력하였다.

<b>결과물에 대한 기 대 효과 및 활용 방안</b>	<input type="checkbox"/> 유기동물 입양 기간 예측 <ul style="list-style-type: none"> <li>■ 유기된 동물의 입양 기간 예측 가능</li> <li>■ 예측 입양기간에 따라 효율적 홍보 방식을 택하여 입양 순서의 원활한 순환구조를 만든다.</li> <li>■ 예측에 따른 선별적 집중 홍보 가능</li> <li>■ 선별적 집중에 따른 입양 실패 동물 감소</li> </ul>
---	---

		<ul style="list-style-type: none"> <li>■ 입양 실패 감소에 따른 안락사 감소</li> </ul>
		<ul style="list-style-type: none"> <li>□ 동물 보호소               <ul style="list-style-type: none"> <li>■ 전체 동물 홍보에서 선별적 홍보방식으로 홍보 효율을 높임</li> <li>■ 안락사 감소에 의한 동물보호단체의 이미지 상승</li> </ul> </li> </ul>