

# CAT OR DOG: PREDICTIVE MODELING

## STAT GU4243 - APPLIED DATA SCIENCE

Group 6

Columbia University

March 5, 2018

# 1 OUTLINE

## 2 INTRODUCTION

- Us
- Motivation
- Scope

## 3 METHOD

- Exploratory analysis
- Feature extraction
- Statistical machine learning models
- Tuning and training

## 4 RESULTS

- Feature extraction
- Computational cost
- Prediction accuracy
- Final result

## 5 DISCUSSION

- Confusion matrices
- Further directions

# GROUP MEMBERS

Wanting Cheng, Mingkai Deng, Jiongjiong Li, Kai Li, Daniel Parker

# WHY DO THIS?—MOTIVATION



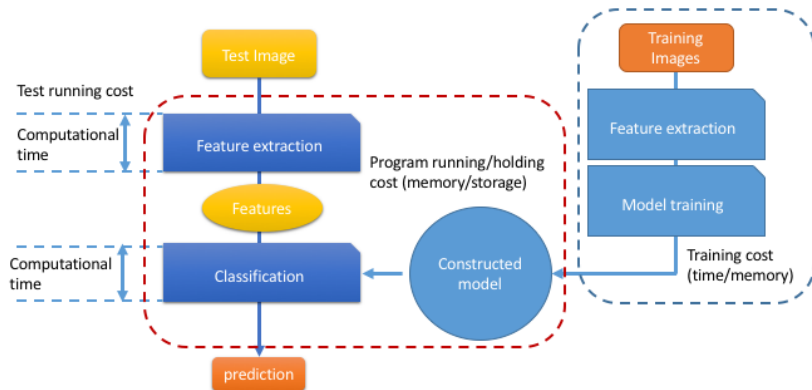
# WHY DO THIS?—MOTIVATION



# SPEC & SCOPE

*[C]arry out model evaluation and selection for predictive analytics on image data ... [using] a set of 4387 labeled images of cats and dogs ... creat[e] a mobile AI program that accurately distinguishes between [them] ... balance between the complexity of variables/features/models used and the predictive performance.*

# SPEC & SCOPE



# EXPLORATORY ANALYSIS

What makes one animal different from another? [Intuition]  
What approaches did previous semesters' groups employ?  
[Research]



# FEATURE EXTRACTION

- ① SIFT = scale-invariant feature transformation.

# FEATURE EXTRACTION

- 1 SIFT = scale-invariant feature transformation.
- 2 HOG = histogram of oriented gradients.

# FEATURE EXTRACTION

- ❶ SIFT = scale-invariant feature transformation.
- ❷ HOG = histogram of oriented gradients.
- ❸ LBP = local binary patterns.

# FEATURE EXTRACTION

- ❶ SIFT = scale-invariant feature transformation.
- ❷ HOG = histogram of oriented gradients.
- ❸ LBP = local binary patterns.
- ❹ HSV = hue, saturation, value.

# FEATURE EXTRACTION

- ❶ SIFT = scale-invariant feature transformation.
- ❷ HOG = histogram of oriented gradients.
- ❸ LBP = local binary patterns.
- ❹ HSV = hue, saturation, value.
- ❺ RGB = red, green, blue.

# STATISTICAL MACHINE LEARNING MODELS

- ➊ Gradient boosting machine—the baseline.

# STATISTICAL MACHINE LEARNING MODELS

- ① Gradient boosting machine—the baseline.
- ② Random forests.

# STATISTICAL MACHINE LEARNING MODELS

- ① Gradient boosting machine—the baseline.
- ② Random forests.
- ③ TensorFlow/Keras neural network.



# STATISTICAL MACHINE LEARNING MODELS

- ➊ Gradient boosting machine—the baseline.
- ➋ Random forests.
- ➌ TensorFlow/Keras neural network.
- ➍ Support vector machine.

# STATISTICAL MACHINE LEARNING MODELS

- ➊ Gradient boosting machine—the baseline.
- ➋ Random forests.
- ➌ TensorFlow/Keras neural network.
- ➍ Support vector machine.
- ➎ Adaptive boosting (“AdaBoost”).

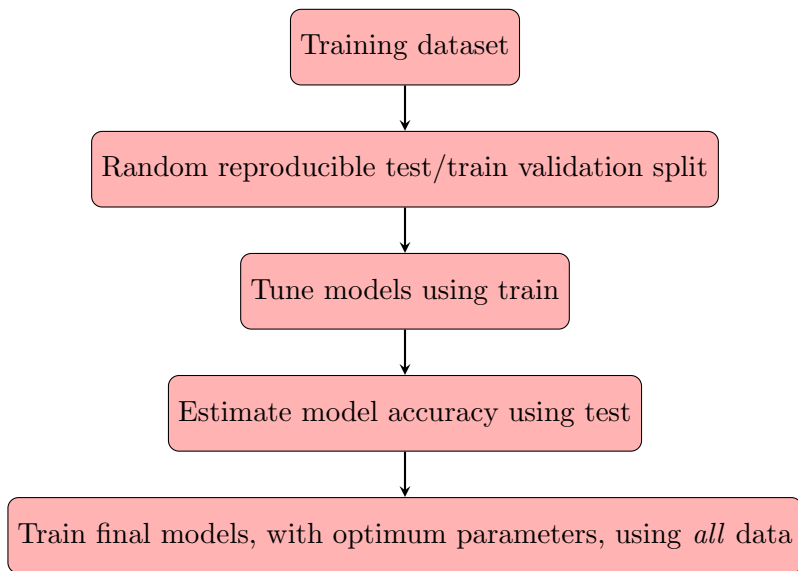
# STATISTICAL MACHINE LEARNING MODELS

- ➊ Gradient boosting machine—the baseline.
- ➋ Random forests.
- ➌ TensorFlow/Keras neural network.
- ➍ Support vector machine.
- ➎ Adaptive boosting (“AdaBoost”).
- ➏ Extreme gradient boosting (“XGBoost”).

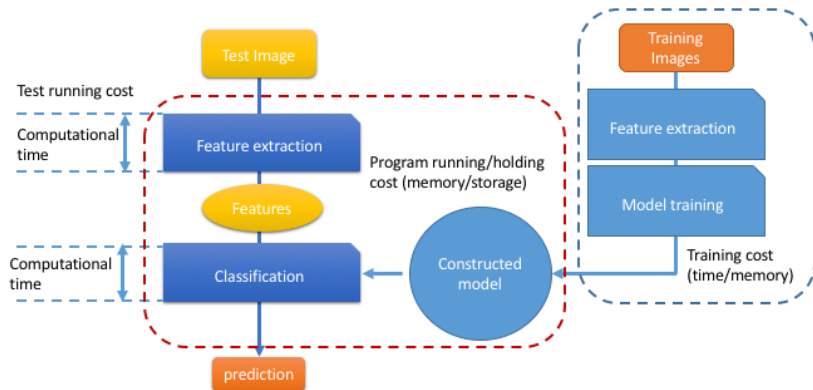
# TUNING AND TRAINING

Simplifying heuristic: use *all* features, rather than subsets.  
Preference for built-in package functions, rather than a generalized syntax.

# HOW WE FINALIZED MODELS—FLOWCHART



# RESULTS



# FEATURE EXTRACTION TIME

Feature type	Color	HOG	LBP
Time (m)	7	5	20

TABLE: Image processing time by feature, in minutes

# TRAINING TIMES—COMPUTATIONAL COST

Model	SIFT	Color	HOG	LBP
GBM	13.452	89.036	116.964	2.74
RF	174.1	905.883	1999.23	24.842
NN	31	65.81	61.58	28.12
SVM	3.484		33.099	0.805
XGBoost	3.829	16.851	27.005	1.67
AdaBoost	16.34	103.61	142.45	3.25

TABLE: Training time per model, in seconds



# PREDICTION ACCURACY

Model	SIFT	Color	HOG	LBP
GBM	73.25	69.5	75.25	69.
RF	72.25	73.	74.25	69.
NN	75.75	64.5	76.75	69.
SVM	77.5		77.5	69.75
XGBoost	72	72	77.25	66.5
AdaBoost	72.75	69.5	71.75	69.75

TABLE: Prediction accuracy by model, in percentage

# FINAL RESULT

Our final feature extraction and model choice is HOG features with XGBoost. This choice optimizes extraction time, training time, and prediction accuracy.

# CONFUSION MATRICES: RANDOM FOREST

RF Prediction	Dog	Cat
Truth		
Dog	54	19
Cat	92	35

TABLE: RF-SIFT Confusion Matrix

RF Prediction	Dog	Cat
Truth		
Dog	247	26
Cat	77	50

TABLE: RF-HOG Confusion Matrix

# CONFUSION MATRICES: GRADIENT BOOSTING MACHINE

GBM Prediction	Dog	Cat
Truth		
Dog	229	44
Cat	63	64

TABLE: GBM-SIFT Confusion Matrix

GBM Prediction	Dog	Cat
Truth		
Dog	232	41
Cat	58	69

TABLE: GBM-HOG Confusion Matrix

# CONFUSION MATRICES: NEURAL NETWORK

NN Prediction	Dog	Cat
Truth		
Dog	233	40
Cat	55	72

TABLE: NN-SIFT Confusion Matrix

NN Prediction	Dog	Cat
Truth		
Dog	263	10
Cat	89	38

TABLE: NN-HOG Confusion Matrix

# FURTHER DIRECTIONS TO EXPLORE

- ❶ Other extractions and combinations thereof.

# FURTHER DIRECTIONS TO EXPLORE

- 1 Other extractions and combinations thereof.
- 2 Dataset manipulation to “grow” more training data for free.

# FURTHER DIRECTIONS TO EXPLORE

- 1 Other extractions and combinations thereof.
- 2 Dataset manipulation to “grow” more training data for free.
- 3 Other models.



# FURTHER DIRECTIONS TO EXPLORE

- ❶ Other extractions and combinations thereof.
- ❷ Dataset manipulation to “grow” more training data for free.
- ❸ Other models.
- ❹ Ensembling.

# FURTHER DIRECTIONS TO EXPLORE

- 1 Other extractions and combinations thereof.
- 2 Dataset manipulation to “grow” more training data for free.
- 3 Other models.
- 4 Ensembling.
- 5 ...

Thank you!