

Mathematical Foundations of ML (CS 4785/5783)

Lecture 1

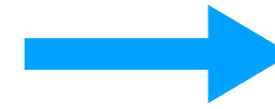
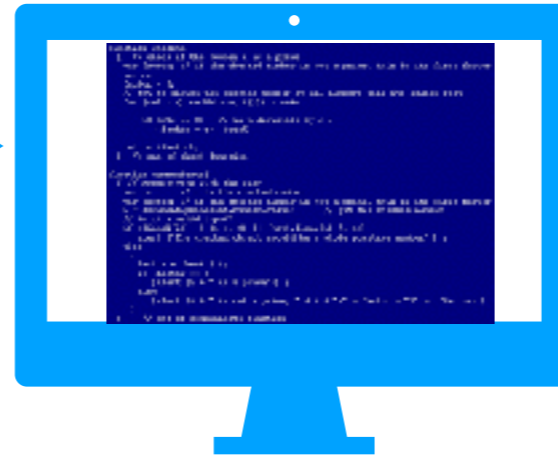
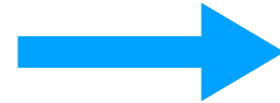
Setting up the Learning Problem

<http://www.cs.cornell.edu/Courses/cs4783/2022sp/notes01.pdf>

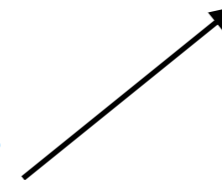
TRADITIONAL COMPUTER SCIENCE

Task
Eg. Sorting

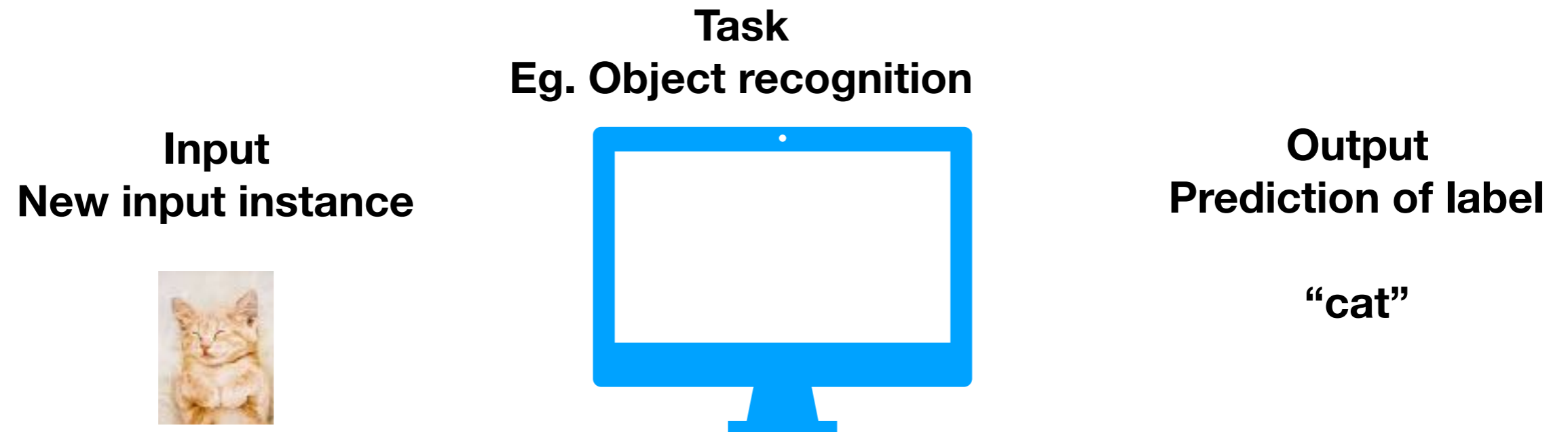
Input
Eg: 2, 4, 3, 8, 7



Output
Eg: 2, 3, 4, 7, 8



MACHINE LEARNING



MACHINE LEARNING

Input
New input instance



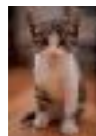
Task
Eg. Object recognition



Output
Prediction of label

“cat”

Input
A set of input/output pairs



, “cat”



, “dog”



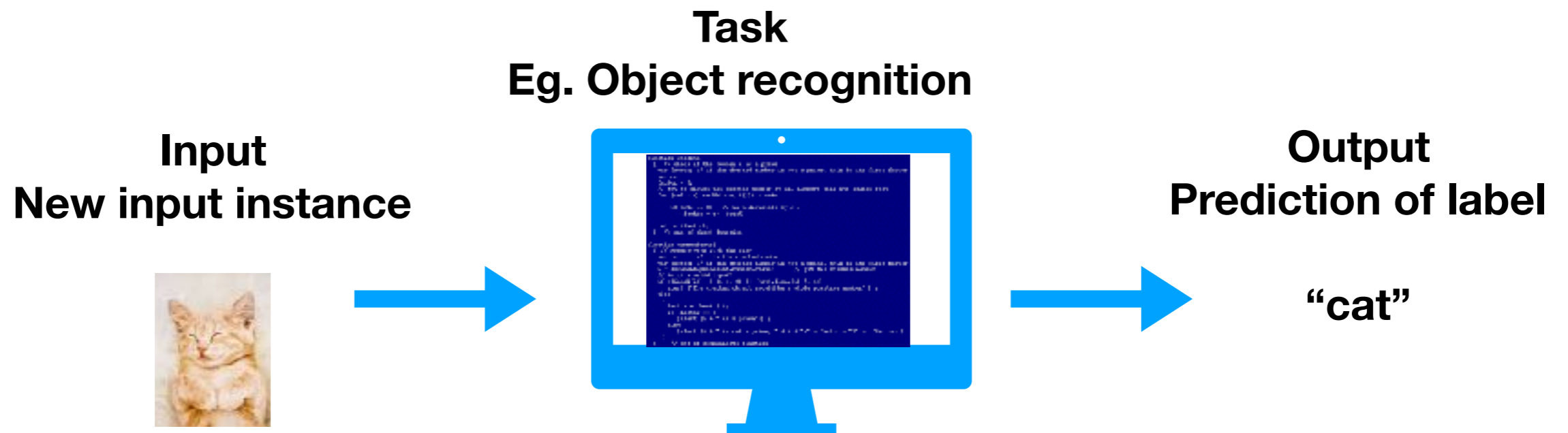
, “dog”



, “cat”

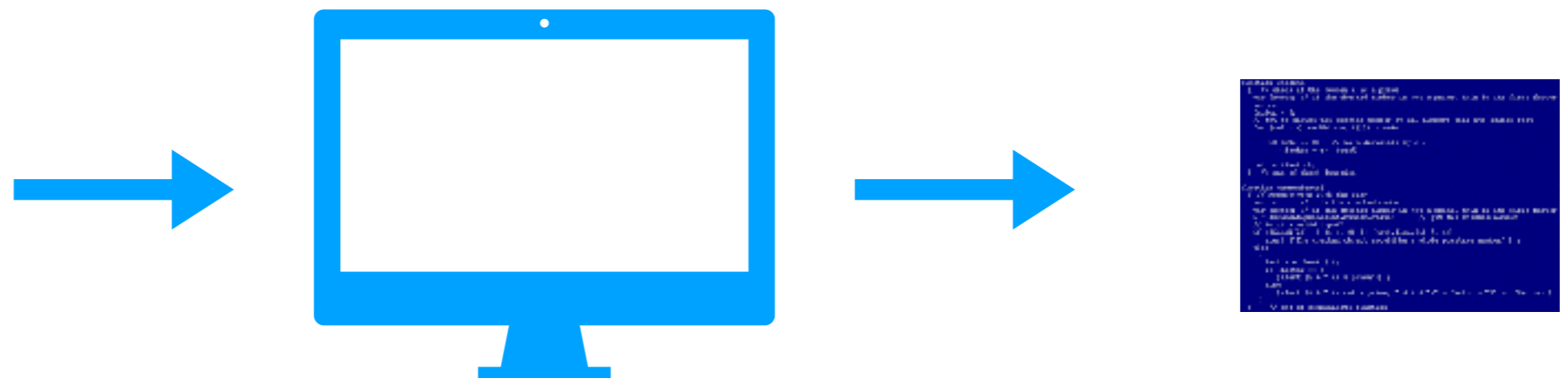


MACHINE LEARNING



Input
A set of input/output pairs

-  , “cat”
-  , “dog”
-  , “dog”
-  , “cat”



WHAT IS MACHINE LEARNING

“The field of study that gives computers the ability to learn without being explicitly programmed” - Arthur Lee Samuel

“A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E ” - Tom Mitchell

LEARNING PROBLEM : BASIC NOTATION

- Input space/ feature space : \mathcal{X}
(Eg. bag-of-words, n-grams, vector of grey-scale values, user-movie pair to rate)
- Output space/ label space \mathcal{Y}
(Eg. $\{\pm 1\}$, $[K]$, \mathbb{R} -valued output, structured output)
- Loss function : $\ell : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}$
(Eg. 0-1 loss $\ell(y', y) = \mathbf{1}\{y' \neq y\}$, sq-loss $\ell(y', y) = (y - y')^2$), absolute loss $\ell(y', y) = |y - y'|$
Measures performance/cost per instance (inaccuracy of prediction/ cost of decision).

TWO SCENARIOS

Universe of instances



U

$$f_{i^*} \left(\begin{array}{c} \text{[Dog Image]} \end{array} \right) = \text{"not cat"}$$
$$f_{i^*} \left(\begin{array}{c} \text{[Cat Image]} \end{array} \right) = \text{"cat"}$$

i^* in $[N]$ is unknown

Amongst set of models $\{f_1, \dots, f_N\}$

There is the perfect model f_{i^*}

Two Scenarios

SCENARIO I

Universe of instances



U

Draw n instances from the universe at random and label them

$$S = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

S is called training set!

x_i 's are images taken from the universe

$$y_i = f_{i^*}(x_i)$$

Learning algorithm has access to the models $\{f_1, \dots, f_N\}$

Goal: return a model with small classification error

SCENARIO I

What should the learning algorithm be?

What kind of guarantee can we provide on its error?

How does our guarantee (bound) on error depend on N the number of models, on n the number of samples we drew?

SCENARIO I

Algorithm: return any classifier that is consistent with S

Return: $\hat{f}_S \in \{f_i : \forall t \in [n], f_i(x_t) = y_t\}$

Error bound:

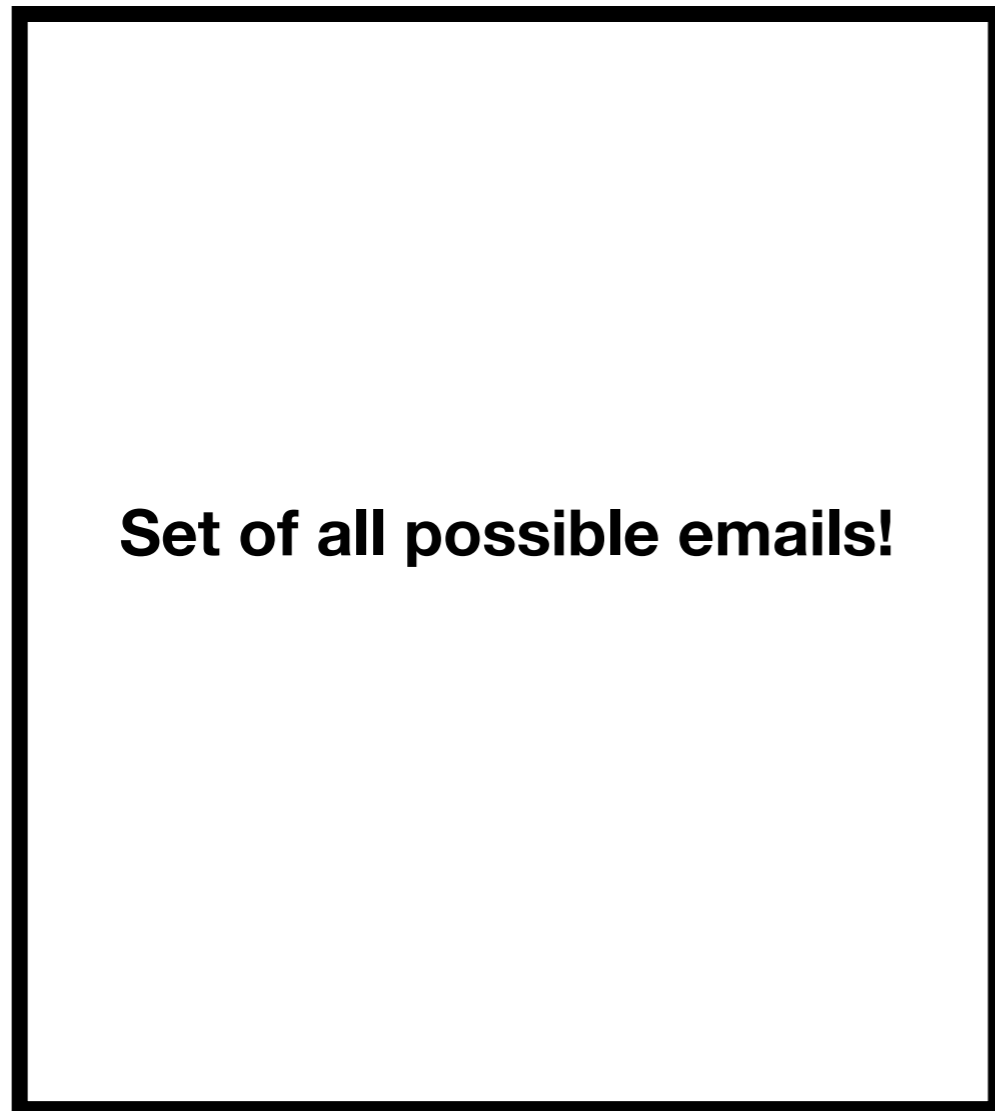
For any $\delta > 0$, with probability at least $1 - \delta$ over draws of S ,

$$P(\hat{f}_S(x) \neq y) \leq \frac{\log(N/\delta)}{n}$$

PAC: Probably Approximately Correct

SCENARIO II

Universe of instances



U

On each round t :

Email x_t is composed, possibly by spammer!

System classifies email as \hat{y}_t

True label $y_t = f_{i^*}(x_t)$ revealed

We get feedback every round. But spammer can pick next email.

Goal: Make as few mistakes as possible.

SCENARIO II

What should the learning algorithm do?

What is the bound on total number of mistakes made?

SCENARIO II

How about using the same algorithm from scenario 1 for each t (re-run)?

How many mistakes would it make?

SCENARIO II

Algorithm:

Pick $\mathcal{F}_t = \{f_i : i \in [N], \forall s < t, f_i(x_s) = y_s\}$

Set $\hat{y}_t = \text{Majority}(\{f(x_t) : f \in \mathcal{F}_t\})$

Mistake Bound:

$$\sum_t \mathbf{1}\{\hat{y}_t \neq y_t\} \leq \log_2 N$$

Why?

WHAT IS IN THIS COURSE?

1. Statistical Learning theory

A. Generalization Error, Training Vs Test loss, Model Complexity

B. PAC model and VC theory

C. Rademacher Complexity and Uniform convergence

D. Role of Regularization in learning, model selection and validation

E. Algorithmic Stability

2. Online learning

1. Online Bit prediction and Cover's result

2. Perceptron, winnow

3. Online experts problem

4. Online Gradient descent and Mirror descent

3. Boosting

4. Stochastic Optimization and Learning: including understanding Stochastic Gradient Descent

5. Bandit problems: both stochastic and adversarial settings

6. A primer to theory of deep learning: new challenges

7. Computational Learning theory: Computational hardness or learning, proper vs improper learning

8. Societal aspects of ML: Differential Privacy, Right to be Forgotten, Fairness and ML

GRADING

- 3% Class participation
- 4 Assignments worth 40% of your grades
- One prelims worth 30% of your grades
- One Term project worth 27% of your grades
- For CS 5783 additional 2 reading assignment + quizzes on them this will be 10% of grade (prelims 25% and Proj 22%). CS 4783 students can also optionally take this.

ROUGH TIMELINE

- Assignments: there are tentative and subject to changes
 - HW1: Jan 31, HW2: Feb 16, HW3: Mar 9, HW4: Apr 18
 - Each assignment has roughly a week
- Prelims: will be held end of march
- Project: Initial proposal due mid semester (march), there will be a project brainstorming lecture in April. Final report due at the end