

Mathematical Foundations of ML (CS 4785/5783)

Lecture 1

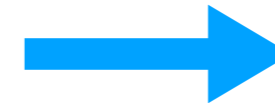
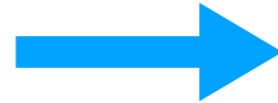
Setting up the Learning Problem

<http://www.cs.cornell.edu/Courses/cs4783/2023fa/notes01.pdf>

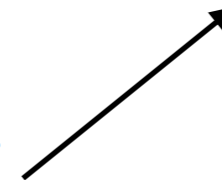
TRADITIONAL COMPUTER SCIENCE

Task
Eg. Sorting

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



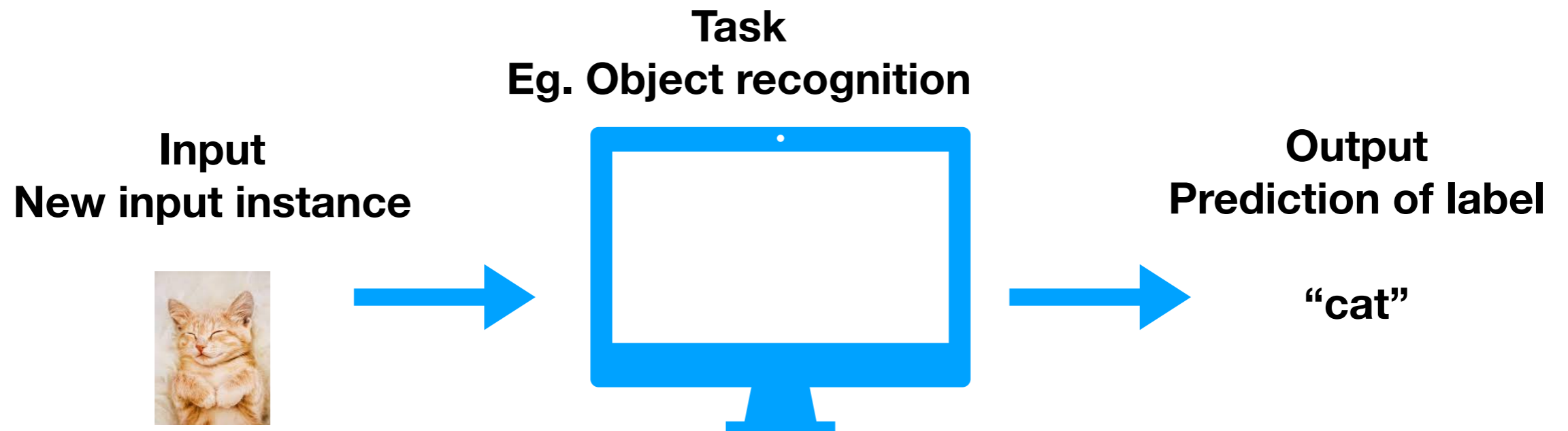
Output
Eg: 2, 3, 4, 7, 8



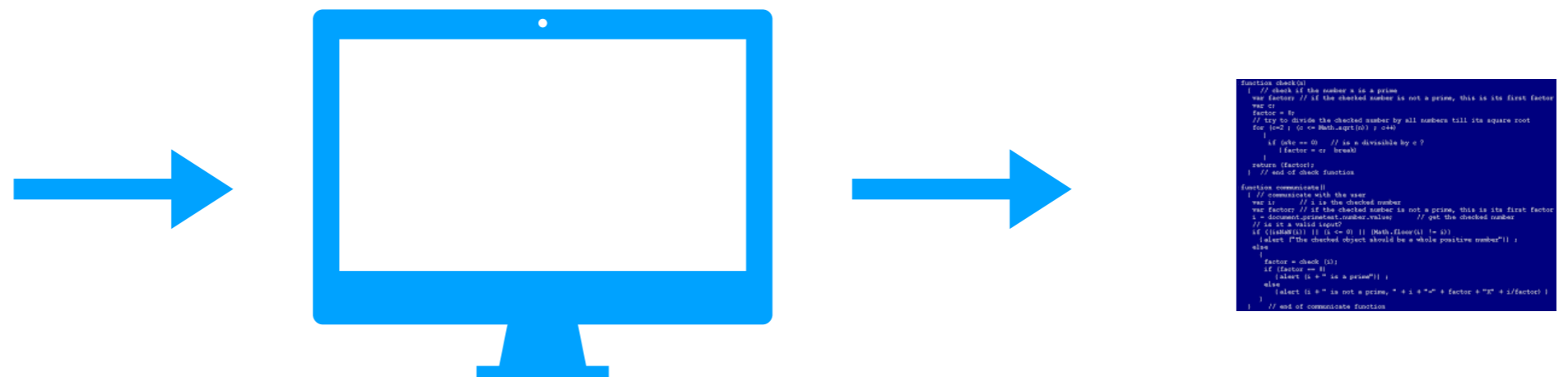
```
function check()
| // check if the number n is a prime
var factor // if the checked number is not a prime, this is its first factor
var n;
factor = n;
// try to divide the checked number by all numbers till its square root
for (var i = 2; i <= Math.sqrt(n); i++)
| if (n % i == 0) // is n divisible by i?
| factor = i; break;
| return (factor);
| // end of check function

function communicate()
| // communicate with the user
var n; // n is the checked number
var factor; // if the checked number is not a prime, this is its first factor
n = document.getElementById("number").value // get the checked number
// is it a valid input?
if (isNaN(n) || n <= 0 || Math.floor(n) != n)
| alert ("The checked object should be a whole positive number!");
| else
| factor = check (n);
| if (factor == n)
| | alert (n + " is a prime!");
| else
| | alert (n + " is not a prime, * " + n + " = " + factor + "*" + n / factor);
| // end of communicate function
|
```

MACHINE LEARNING



Input
A set of input/output pairs



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:

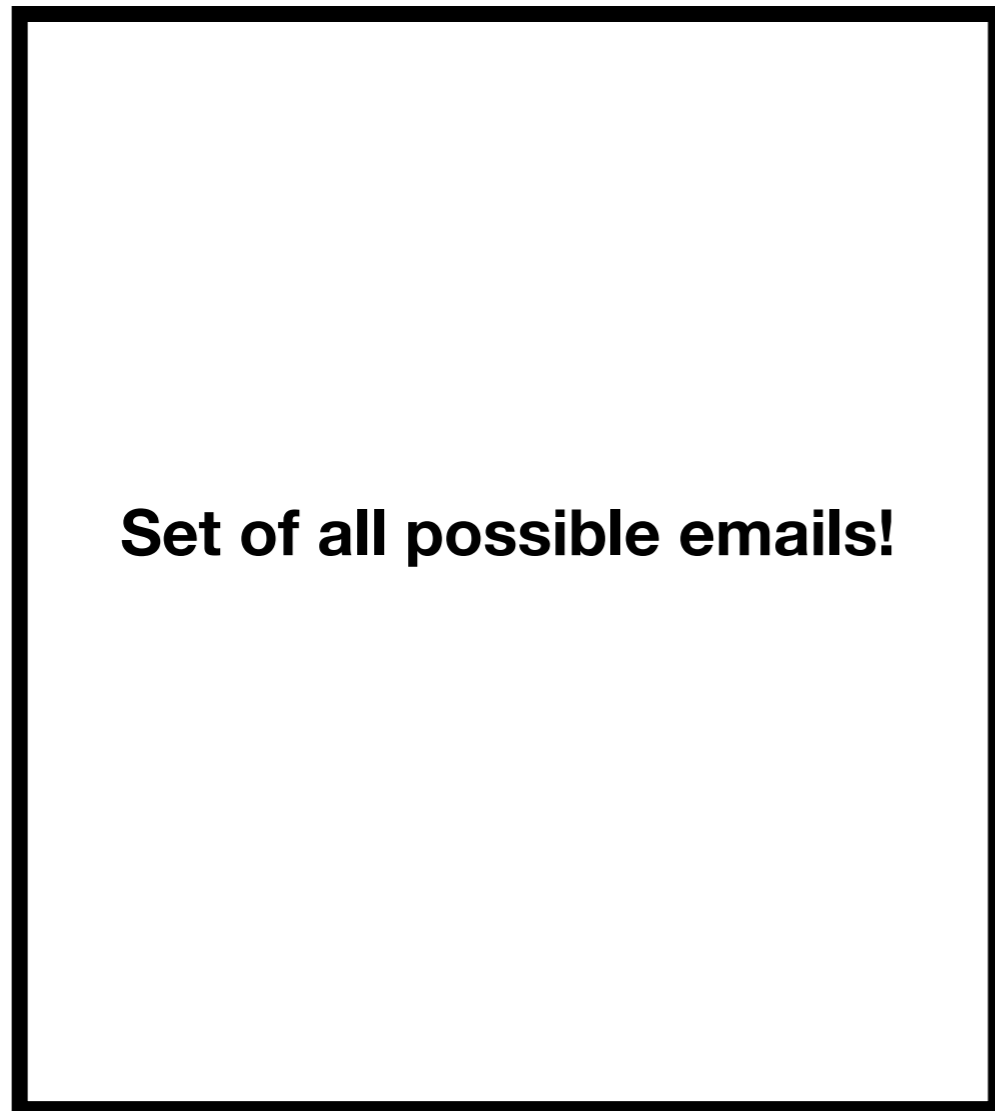
For any $\delta > 0$, with probability at least $1 - \delta$ over draws of S ,

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

PAC: Probably Approximately Correct

SCENARIO II

Universe of instances



Set of all possible emails!

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. Perceptron

2. Online experts problem

3. 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: Aug 23, HW2: Sep 13, HW3: Oct 11, HW4: Nov 13
 - Each assignment has roughly a week
- Prelims: Oct 25th in class
- Project: Initial proposal due mid semester (early oct), there will be a project brainstorming lecture in November (Nov 6th tentative). Final report due exam week