**Oleksandr Romanko**, Ph.D.
Senior Research Analyst, Risk Analytics, Watson Financial Services, IBM Canada
Adjunct Professor, University of Toronto

**Xinrun Shan**
Teaching Assistant, University of Toronto

# MIE1624H – Introduction to Data Science and Analytics
## Assignment 1 – Salary Classification Problem

University of Toronto
October 7, 2019

# Kaggle ML & DS Survey 2018 – https://www.kaggle.com/kaggle/kaggle-survey-2018/

# Kaggle ML & DS Survey 2018 – https://www.kaggle.com/kaggle/kaggle-survey-2018/

# Assignment goal

Based on "2018 Kaggle ML & DS Survey" **identify skills** that are **important for data scientists** to posses based on salary level as a target variable. The goal is to **identify skills and qualifications from the survey that are important to learn in data science field in order to have a high salary**. Supervised machine learning algorithm (logistic regression) would be used for the assignment. In your f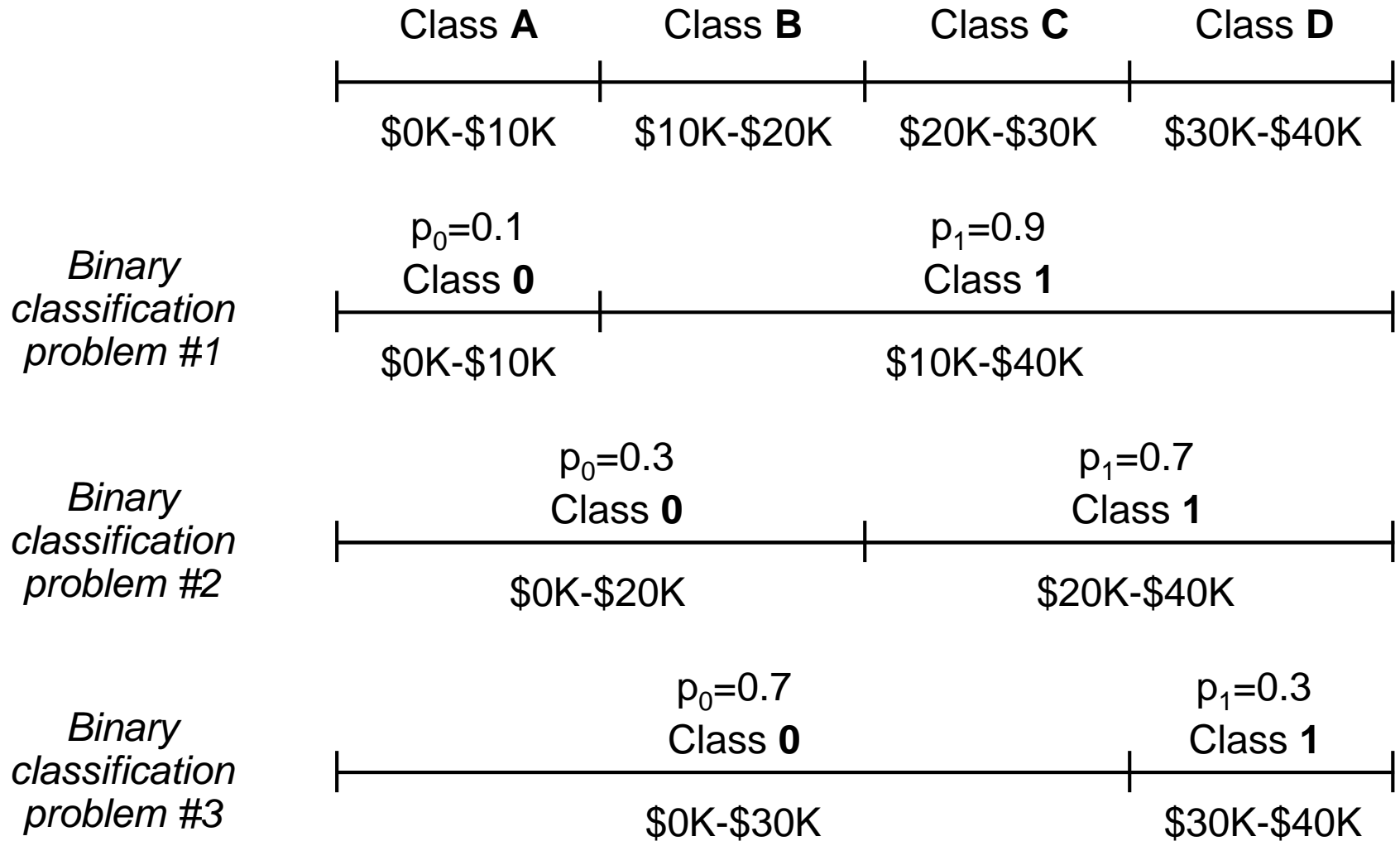ree time, you can also experiment with other supervised learning algorithms, e.g., decision trees, and unsupervised machine learning algorithms, e.g., clustering of skills.

**Possible decisions to make based on modeling results** (not part of the assignment):
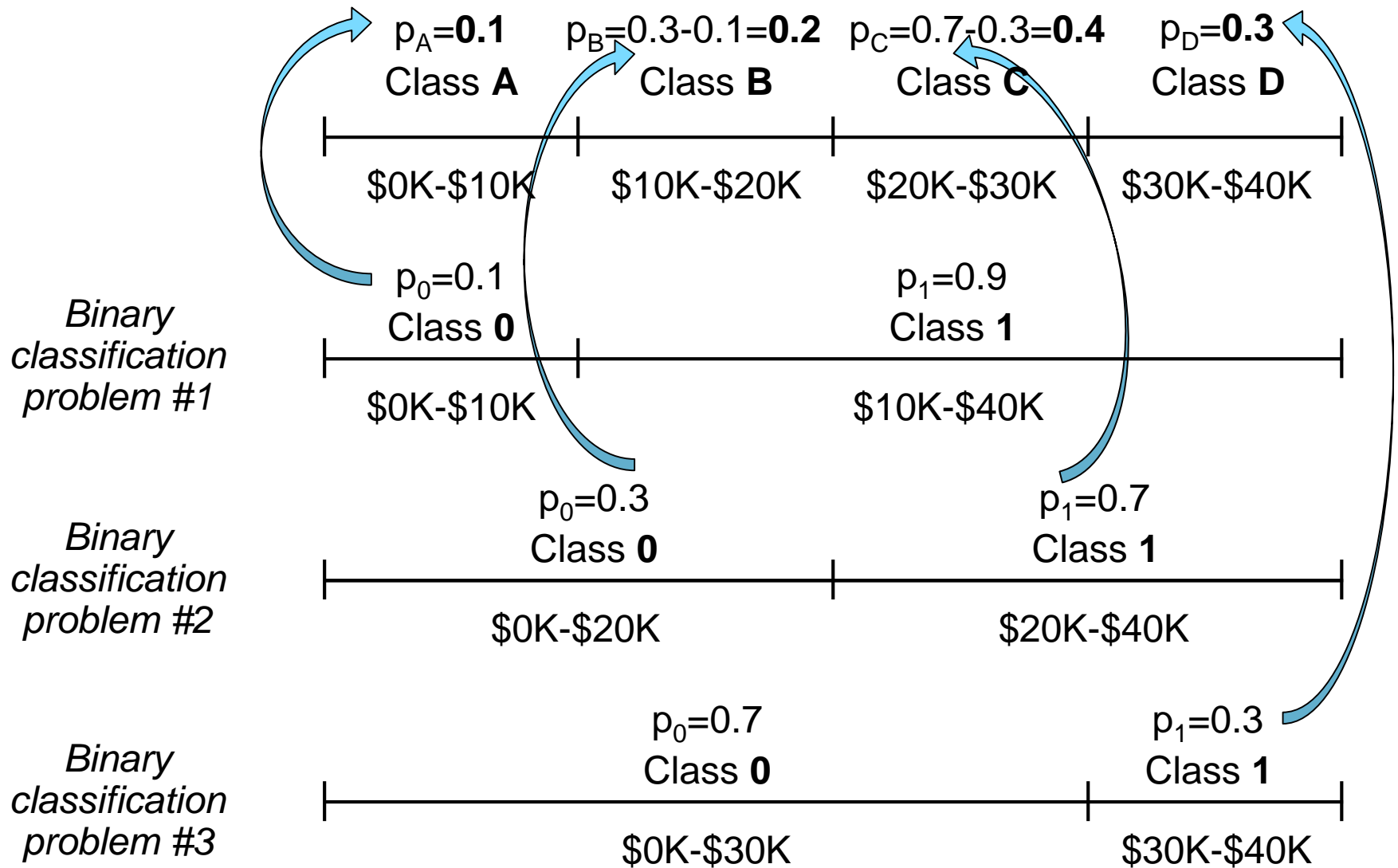- Design a study program at your university, e.g., Master of DS & AI
- Build a recommender system for taking online courses
- Design personalized study program that maximizes your salary
- Any other creative use of insights from the survey

**If you are interested to continue with this type of project as your research project or MEng project, please e-mail the course instructor.**

# Ordinary multi-class classification algorithm

| | Class **A** | Class **B** | Class **C** | Class **D** |
|---|---|---|---|---|
| | \$0K-\$10K | \$10K-\$20K | \$20K-\$30K | \$30K-\$40K |

*Binary classification problem #1*

| $p_0$=0.1 Class **0** | $p_1$=0.9 Class **1** |
|---|---|
| \$0K-\$10K | \$10K-\$40K |

*Binary classification problem #2*

| $p_0$=0.3 Class **0** | $p_1$=0.7 Class **1** |
|---|---|
| \$0K-\$20K | \$20K-\$40K |

*Binary classification problem #3*

| $p_0$=0.7 Class **0** | $p_1$=0.3 Class **1** |
|---|---|
| \$0K-\$30K | \$30K-\$40K |

# Ordinary multi-class classification algorithm

# Ordinary multi-class classification algorithm

| $p_A=\mathbf{0.1}$ | $p_B=\mathbf{0.2}$ | $p_C=\mathbf{0.4}$ | $p_D=\mathbf{0.3}$ |
|---|---|---|---|
| Class **A** | Class **B** | Class **C** | Class **D** |
| $0K-$10K | $10K-$20K | $20K-$30K | $30K-$40K |

*Binary classification problem #1*

| $p_0=0.1$ | $p_1=0.9$ |
|---|---|
| Class **0** | Class **1** |
| $0K-$10K | $10K-$40K |

*Binary classification problem #2*

| $p_0=0.3$ | $p_1=0.7$ |
|---|---|
| Class **0** | Class **1** |
| $0K-$20K | $20K-$40K |

*Binary classification problem #3*

| $p_0=0.7$ | $p_1=0.3$ |
|---|---|
| Class **0** | Class **1** |
| $0K-$30K | $30K-$40K |

# Algorithm accuracy



$100-$125K ⌐⌐ 30% accuracy
$90-$150K ⌐———⌐ 70% accuracy
$80-$200K ⌐————————⌐ 94% accuracy