# Advanced Machine Learning

# Imbalanced Learning:

# Cost-Sensitive Learning, Part 1

- Cost-sensitive learning:
  - Classical learning: data sets are balanced, and all errors have equal costs
  - We now assume given, unequal cost
  - And try to minimize them in expectation
- Applications:
  - Medicine — Misdiagnosing as healthy vs. having a disease
  - (Extreme) Weather prediction — Incorrectly predicting that no hurricane occurs
  - Credit granting — Lending to a risky client vs. not lending to a trustworthy client.

Confusion matrix

| | True class | |
|---|---|---|
| | $y = 1$ | $y = -1$ |
| Pred. $\hat{y} = 1$ | TP | FP |
| class $\hat{y} = -1$ | FN | TN |

| Pred. | Truth | |
|---|---|---|
| | Default | Pays Back |
| Default | 0 | 10 |
| Pays Back | 1000 | 0 |

Cost matrix

| | $y = 1$ | $y = -1$ |
|---|---|---|
| Pred. $\hat{y} = 1$ | $C(1,1)$ | $C(1,-1)$ |
| class $\hat{y} = -1$ | $C(-1,1)$ | $C(-1,-1)$ |

- In these examples, **the costs of a false negative is much higher than the costs of a false positive**.

**Learning goals**
- Cost matrix
- Minimum expected cost principle
- Optimal theoretical threshold

- In some applications, the costs are **unknown** → need to be specified by experts, or be learnt.

# COST-SENSITIVE LEARNING: IN A NUTSHELL

- **Input: cost matrix C**
  - Classical learning: data sets are balanced, and all errors have equal costs
  - We now assume given, unequal cost
  - Classification: try to minimize $C(j, k)$ in expectation

| | True Class $y$ | | | |
|---|---|---|---|---|
| | 1 | 2 | $\ldots$ | $g$ |
| 1 | $C(1, 1)$ | $C(1, 2)$ | $\ldots$ | $C(1, g)$ |
| 2 | $C(2, 1)$ | $C(2, 2)$ | $\ldots$ | $C(2, g)$ |
| $\vdots$ | | | | |
| $g$ | $C(g, 1)$ | $C(g, 2)$ | | $C(g, g)$ |

- Applications:
  - Medicine — Misdiagnosing as healthy vs. having a disease
  - (Extreme) Weather prediction — Incorrectly predicting that no hurricane occurs
  - Credit granting — Lending to a risky client vs. not lending to a trustworthy client.
  - In these examples, **the costs of a false negative is much higher than the costs of a false positive.**

- $C(j, k)$ is the cost of classifying class $k$ as $j$,
- 0-1-loss would simply be: $C(j, k) = \mathbb{1}_{[j \neq k]}$

| | Truth | | |
|---|---|---|---|
| Pred. | | Default | Pay Back |
| | Default | 0 | 10 |
| | Pay Back | 100 | 0 |

- **C** designed by experts with domain knowledge
  1. Too low costs: not enough change in model, still costly errors
  2. Too high costs: might never predict costly classes
  - In some applications, the costs are ... specified by experts, or be learnt.

- Common heuristic for imbalanced data sets:
  - $C(j, k) = \frac{n_j}{n_k}$ with $n_k \ll n_j$, misclassifying a minority class $k$ as a majority class $j$
  - $C(j, k) = 1$ with $n_j \ll n_k$, misclassifying a majority class $k$ as a minority class $j$
  - 0 for a correct classification

|  |  | True Class $y$ | | |
|---|---|---|---|---|
| Classification | | 1 | 2 | ... | g |
| | 1 | $C(1,1)$ | $C(1,2)$ | ... | $C(1,g)$ |
| | 2 | | $C(2,2)$ | ... | $C(2,g)$ |
| | ... | | | | ... |
| | g | | $C(g,2)$ | ... | $C(g,g)$ |

- $C(j, k)$ is the cost of classifying class $k$ as $j$,
- 0-1-loss would simply be: $C(j, k) = \mathbb{1}_{[j \neq k]}$
- Imbalanced binary classification:
- C designed by experts with domain knowledge
  1. Too low costs: not enough change in model, still costly errors
  2. Too high costs: might never predict costly classes

|  | True class | |
|---|---|---|
|  | $y = 1$ | $y = -1$ |
| Pred. $\hat{y} = 1$ | 0 | 1 |
| class $\hat{y} = -1$ | $\frac{n_-}{n_+}$ | 0 |

- So: much higher costs for FNs

- Suppose we have: for imbalanced data sets:
  - a cost matrix **C** with $n_k \ll n_j$,
  - knowledge of the true posterior $p(\cdot \mid \mathbf{x})$ majority class $j$
- Predict class j with smallest expected costs when marginalizing over true classes: misclassifying a majority class $k$ as a minority class $j$
  - 0 for a correct classification

$$\mathbb{E}_{K \sim p(\cdot \mid \mathbf{x})}(C(j,K)) = \sum_{k=1}^{g} p(k \mid \mathbf{x})\, C(j,k)$$

- Imbalanced binary classification:
- If we trust we trust a probabilistic classifier, we can convert its scores to labels:

| | True class | |
|---|---|---|
| | $y = 1$ | $y = -1$ |
| Pred. $\hat{y} = 1$ | | 1 |

$$h(\mathbf{x}) := \arg\min_{j=1,\dots,g} \sum_{k=1}^{g} \pi_k(\mathbf{x})\, C(j,k).$$

- So: much higher costs for FNs
- Can be better to take a less probable class ( ▶ Elkan et. al. 2001 )

# OPTIMAL THRESHOLD FOR BINARY CASE

- Optimal decisions do not change if
  - $C$ is multiplied by positive constant
  - $C$ is added with constant shift
- Scale and shift $C$ to get simpler $C'$:

|  | True class | |
|---|---|---|
|  | $y = 1$ | $y = -1$ |
| Pred. $\hat{y} = 1$ | $C'(1, 1)$ | $1$ |
| class $\hat{y} = -1$ | $C'(-1, 1)$ | $0$ |

where

- $C'(-1, 1) = \frac{C(-1,1) - C(-1,-1)}{C(1,-1) - C(-1,-1)}$
- $C'(1, 1) = \frac{C(1,1) - C(-1,-1)}{C(1,-1) - C(-1,-1)}$

- We predict $\mathbf{x}$ as class 1 if

$$\mathrm{E}_{K \sim p(\cdot \mid \mathbf{x})}(C'(1, K)) \leq \mathrm{E}_{K \sim p(\cdot \mid \mathbf{x})}(C'(-1, K))$$

- Can be better to take a less probable class ( ► Elkan et. al. 2001 )

- Let's unroll the expected value and use $C'$:
- Optimal decisions do not change if $C$ is multiplied by a positive const. and a shift is added

$$p(-1 \mid x)C'(1,-1) + p(1 \mid x)C'(1,1) \le p(-1 \mid x)C'(-1,-1) + p(1 \mid x)C'(-1,1)$$

- Scale and shift $C$ to get simpler $C'$:

$$\Rightarrow p(1 \mid x) \ge \frac{1}{C'(-1,1) - C'(1,1) + 1}$$

$$\Rightarrow p(1 \mid x) \ge \frac{C(1,-1) - C(-1,-1)}{C(-1,1) - C(1,1) + C(1,-1) - C(-1,-1)} = c^*$$

| | True class | |
|---|---|---|
| | $K = 1$ | $K = -1$ |
| Pred. $\hat{y} = 1$ | $C'(1,1)$ | $1$ |
| class $y = -1$ | $C'(-1,1)$ | $0$ |

- If even $C(1,1) = C(-1,-1) = 0$, we get:

$$p(1 \mid x) \ge \frac{C(1,-1)}{C(-1,1) + C(1,-1)} = c^*$$

where

- $C'(-1,1) = \frac{C(-1,1) - C(-1,-1)}{C(1,-1) - C(-1,-1)}$

- $C'(1,1) = \frac{C(1,1) - C(-1,-1)}{C(1,-1) - C(-1,-1)}$

- Optimal threshold $c^*$ for probabilistic classifier

- We predict $x$ as class $h(x) := 2 \cdot 1_{[\pi(x) \ge c^*]} - 1$

$$E_{K \sim p(\cdot \mid x)}(C'(1,K)) \le E_{K \sim p(\cdot \mid x)}(C'(-1,K))$$

# OPTIMAL THRESHOLD FOR BINARY CASE

- Let's unroll the expected value and use $\mathbf{C}'$:

$$p(-1 \mid \mathbf{x})C'(1, -1) + p(1 \mid \mathbf{x})C'(1, 1) \leq p(-1 \mid \mathbf{x})C'(-1, -1) + p(1 \mid \mathbf{x})C'(-1, 1)$$

$$\Rightarrow [1 - p(1 \mid \mathbf{x})] \cdot 1 + p(1 \mid \mathbf{x})C'(1, 1) \leq p(1 \mid \mathbf{x})C'(-1, 1)$$

$$\Rightarrow p(1 \mid \mathbf{x}) \geq \frac{1}{C'(-1, 1) - C'(1, 1) + 1}$$

$$\Rightarrow p(1 \mid \mathbf{x}) \geq \frac{C(1, -1) - C(-1, -1)}{C(-1, 1) - C(1, 1) + C(1, -1) - C(-1, -1)} = c^*$$

- If even $C(1, 1) = C(-1, -1) = 0$, we get:

$$p(1 \mid \mathbf{x}) \geq \frac{C(1, -1)}{C(-1, 1) + C(1, -1)} = c^*$$

- Optimal threshold $c^*$ for probabilistic classifier

$$h(\mathbf{x}) := 2 \cdot \mathbb{1}_{[\pi(\mathbf{x}) \geq c^*]} - 1$$