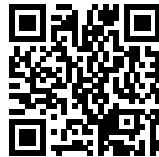# Machine Learning I

Bjoern Andres
bjoern.andres@tu-dresden.de

Machine Learning for Computer Vision
Faculty of Computer Science
TU Dresden

Version $0.4\beta$

# Contents

# Chapter 1

# Introduction

## 1.1 Notation

We shall use the following notation:

- We write "iff" as shorthand for "if and only if"
- For any $m \in \mathbb{N}$, we define $[m] = \{0, \ldots, m-1\}$.
- For any set $A$, we denote by $2^A$ the power set of $A$
- For any set $A$ and any $m \in \mathbb{N}$, we denote by $\binom{A}{m} = \{B \in 2^A \mid |B| = m\}$ the set of all $m$-elementary subsets of $A$
- For any sets $A, B$, we denote by $B^A$ the set of all maps from $A$ to $B$

# Chapter 2

# Supervised learning

## 2.1 Intuition

Informally, supervised learning is the problem of finding in a family $g : \Theta \to Y^X$ of functions, one $g_\theta : X \to Y$ that minimizes a weighted sum of two objectives:

1. $g$ deviates little from a finite set $\{(x_s, y_s)\}_{s \in S}$ of input-output-pairs
2. $g$ has low complexity, as quantified by a function $R : \Theta \to \mathbb{R}_0^+$

We note that the family $g$ can have meaning beyond a mere parameterization of functions from $X$ to $Y$. For instance, $\Theta$ can be a set of forms, $g$ the functions defined by these forms, and $R$ the length of forms. In that case, supervised learning is really an optimization problem over forms of functions, and $R$ penalizes the complexity of these forms. Moreover, $g$ can be chosen so as to constrain the set of functions from $X$ to $Y$ in the first place.

We concentrate exclusively on the special case where $Y$ is finite. In fact, we concentrate on the case where $Y = \{0, 1\}$ in this chapter and reduce more general cases to this case in Chapter 4.

Moreover, we allow ourselves to take a detour by not optimizing over a family $g : \Theta \to \{0, 1\}^X$ directly but instead optimizing over a family $f : \Theta \to \mathbb{R}^X$ and defining $g$ w.r.t. $f$ via a function $L : \mathbb{R} \times \{0, 1\} \to \mathbb{R}_0^+$, called a *loss function*, such that

$$\forall \theta \in \Theta \; \forall x \in X : \quad g_\theta(x) = \operatorname*{argmin}_{\hat{y} \in \{0,1\}} \; L(f_\theta(x), \hat{y}) \; . \tag{2.1}$$

## 2.2 Definition

**Definition 1** For any $S \neq \varnothing$ finite, called a set of *samples*, any $X \neq \varnothing$, called an *attribute space* and any $x : S \to X$, the tuple $(S, X, x)$ is called *unlabeled data*.

For any $y : S \to \{0, 1\}$, given in addition and called a *labeling*, the tuple $(S, X, x, y)$ is called *labeled data*.

**Definition 2** For any labeled data $T = (S, X, x, y)$, any $\Theta \neq \varnothing$ and family of functions $f : \Theta \to \mathbb{R}^X$, any $R : \Theta \to \mathbb{R}_0^+$, called a *regularizer*, any $L : \mathbb{R} \times \{0, 1\} \to \mathbb{R}_0^+$, called a *loss function*, and any $\lambda \in \mathbb{R}_0^+$, called a *regularization parameter*, we define the following optimization problems:

- The instance of the *supervised learning problem* w.r.t. $T, \Theta, f, R, L$ and $\lambda$ is defined as

$$\inf_{\theta \in \Theta} \quad \lambda R(\theta) + \frac{1}{|S|} \sum_{s \in S} L(f_\theta(x_s), y_s) \tag{2.2}$$

- The instance of the *exact supervised learning problem* w.r.t. $T, \Theta, f$ and $R$ is defined as

$$\inf_{\theta \in \Theta} \quad R(\theta) \tag{2.3}$$

$$\text{subject to} \quad \forall s \in S : \quad f_\theta(x_s) = y_s \tag{2.4}$$

- The instance of the *bounded regularity problem* w.r.t. $T, \Theta, f, R$ and $m$ is to decide whether there exists a $\theta \in \Theta$ such that

$$R(\theta) \leq m \tag{2.5}$$

$$\forall s \in S: \quad f_\theta(x_s) = y_s \tag{2.6}$$

**Definition 3** For any unlabeled data $T = (S, X, x)$, any $\hat{f} : X \to \mathbb{R}$ and any $L : \mathbb{R} \times \{0,1\} \to \mathbb{R}_0^+$, the instance of the *inference problem* w.r.t. $T, f$ and $L$ is defined as

$$\min_{y' \in \{0,1\}^S} \sum_{s \in S} L(\hat{f}(x_s), y'_s) \tag{2.7}$$

**Lemma 1** *The solutions to the inference problem are the* $y : S \to \{0,1\}$ *such that*

$$\forall s \in S: \quad y_s \in \underset{\hat{y} \in \{0,1\}}{\mathrm{argmin}} \ L(\hat{f}(x_s), \hat{y}) \ . \tag{2.8}$$

*Moreover, if*

$$\hat{f}(X) \subseteq \{0,1\} \tag{2.9}$$

*and*

$$\forall r \in \mathbb{R} \ \forall \hat{y} \in \{0,1\}: \quad L(r, \hat{y}) = \begin{cases} 0 & \text{if } r = \hat{y} \\ 1 & \text{otherwise} \end{cases} \tag{2.10}$$

*then*

$$\forall s \in S: \quad y'_s = \hat{f}(x_s) \ . \tag{2.11}$$

PROOF  Generally, we have

$$\min_{y \in \{0,1\}^S} \sum_{s \in S} L(\hat{f}(x_s), y_s) = \sum_{s \in S} \min_{y_s \in \{0,1\}} L(\hat{f}(x_s), y_s) \tag{2.12}$$

By (2.9), $L(\hat{f}(x_s), \hat{f}(x_s))$ is well-defined for any $s \in S$. By (2.10) and non-negativity of $L$, we have

$$\forall y_s \in \{0,1\}: \quad L(\hat{f}(x_s), \hat{f}(x_s)) = 0 \leq L(\hat{f}(x_s), y_s) \ . \tag{2.13}$$

Thus, $y_s = \hat{f}(x_s)$ is optimal for any $s \in S$.

We note that the exact supervised learning problem formalizes a philosophical principle known as Ockham's razor.

# Chapter 3

# Deciding

## 3.1 Disjunctive normal forms

### 3.1.1 Data

Throughout Section 3.1, we consider binary attributes. More specifically, we consider some finite set $V \neq \varnothing$ and labeled data $T = (S, X, x, y)$ such that $X = \{0, 1\}^V$. Hence, $x\colon S \to \{0, 1\}^V$ and $y\colon S \to \{0, 1\}$.

### 3.1.2 Familiy of functions

Throughout Section 3.1, we identify $\Theta$ with a set of disjunctive normal forms. More specifically, we consider $\Gamma = \{(V_0, V_1) \in 2^V \times 2^V \mid V_0 \cap V_1 = \varnothing\}$ and $\Theta = 2^\Gamma$ and the following definition.

**Definition 4** For any $\theta \in \Theta$ and the $f_\theta\colon \{0, 1\}^V \to \{0, 1\}$ such that

$$\forall x \in \{0, 1\}^V\colon \quad f_\theta(x) = \bigvee_{(V_0, V_1) \in \theta} \prod_{v \in V_0} (1 - x_v) \prod_{v \in V_1} x_v \ , \tag{3.1}$$

the form on the r.h.s. of (3.1) is called the *disjunctive normal form (DNF)* defined by $V$ and $\theta$. The function $f_\theta$ is said to be defined by the DNF. If there exists a $k \in \mathbb{N}$ such that $\forall (V_0, V_1) \in \theta\colon |V_0 \cup V_1| \leq k$, the DNF is also called a *k-DNF*.

For $\theta \in \Theta$ and the $g_\theta\colon \{0, 1\}^V \to \{0, 1\}$ such that

$$\forall x \in \{0, 1\}^V\colon \quad g_\theta(x) = \prod_{(V_0, V_1) \in \theta} \left( \bigvee_{v \in V_0} (1 - x_v) \vee \bigvee_{v \in V_1} x_v \right) \tag{3.2}$$

the form on the r.h.s. of (3.2) is called the *conjunctive normal form (CNF)* defined by $V$ and $\theta$. The function $g_\theta$ is said to be defined by this CNF. If there exists a $k \in \mathbb{N}$ such that $\forall (V_0, V_1) \in \theta\colon |V_0 \cup V_1| \leq k$, the CNF is also called a *k-CNF*.

### 3.1.3 Regularization

**Definition 5** The functions $R_d, R_l : \Theta \to \mathbb{N}_0$ whose values are defined below for any $\theta \in \Theta$ are called the *depth* and *length*, resp., of the forms defined by $\theta$.

$$R_d(\theta) = \max_{(V_0, V_1) \in \theta} (|V_0| + |V_1|) \tag{3.3}$$

$$R_l(\theta) = \sum_{(V_0, V_1) \in \theta} (|V_0| + |V_1|) \tag{3.4}$$

### 3.1.4   Loss function

We consider the loss function $L$ such that

$$\forall r \in \mathbb{R} \; \forall \hat{y} \in \{0,1\}\colon \quad L(r,\hat{y}) = \begin{cases} 0 & r = \hat{y} \\ 1 & \text{otherwise} \end{cases} . \tag{3.5}$$

### 3.1.5   Learning problem

**Definition 6** For any $R \in \{R_l, R_d\}$ and any $\lambda \in [0,1)$, the instance of the *supervised learning problem of DNFs* with respect to $T, L, R$ and $\lambda$ has the form

$$\min_{\theta \in \Theta} \quad \lambda R(\theta) + \frac{1-\lambda}{|S|} \sum_{s \in S} L(f_\theta(x_s), y_s) \tag{3.6}$$

In order to examine its computational complexity, we consider the related decision problems:

**Definition 7** Let $m \in \mathbb{N}$. The instance of the *bounded depth DNF problem (*DEPTH-$m$-DNF*)* w.r.t. $T$ and $m$ is to decide whether there exists a $\theta \in \Theta$ such that

$$R_d(\theta) \leq m \tag{3.7}$$

$$\forall s \in S\colon \quad f_\theta(x_s) = y_s . \tag{3.8}$$

The instance of the *bounded length DNF problem (*LENGTH-$m$-DNF*)* w.r.t. $T$ and $m$ is to decide whether there exists a $\theta \in \Theta$ such that

$$R_l(\theta) \leq m \tag{3.9}$$

$$\forall s \in S\colon \quad f_\theta(x_s) = y_s . \tag{3.10}$$

We relate these problems to SET-COVER.

**Definition 8 (Haussler (1988))** For any instance $(S', \Sigma, m)$ of SET-COVER, the *Haussler data* $T = (S, X, x, y)$ induced by $(S', \Sigma, m)$ is the labeled data such that

- $S = S' \cup \{1\}$
- $X = \{0,1\}^\Sigma$
- $x_1 = 1^\Sigma$ and

$$\forall s \in S' \; \forall \sigma \in \Sigma\colon \quad x_s(\sigma) = \begin{cases} 0 & \text{if } s \in \sigma \\ 1 & \text{otherwise} \end{cases} \tag{3.11}$$

- $y_1 = 1$ and $\forall s \in S'\colon y_s = 0$

**Lemma 2 (Haussler (1988))** *For any instance $(S', \Sigma, m)$ of* SET-COVER*, the Haussler data* $T = (S, X, x, y)$ *induced by* $(S', \Sigma, m)$*, and any* $\Sigma' \in \binom{\Sigma}{m}$*:*

$$\bigcup_{\sigma \in \Sigma'} \sigma = S' \quad \Leftrightarrow \quad \forall s \in S'\colon \prod_{\sigma \in \Sigma'} x_s(\sigma) = 0$$

PROOF

$$\bigcup_{\sigma \in \Sigma'} \sigma = S'$$

$$\Leftrightarrow \quad \forall s \in S' \; \exists \sigma \in \Sigma'\colon \quad s \in \sigma \tag{3.12}$$

$$\Leftrightarrow \quad \forall s \in S' \; \exists \sigma \in \Sigma'\colon \quad x_s(\sigma) = 0 \tag{3.13}$$

$$\Leftrightarrow \quad \forall s \in S'\colon \quad \prod_{\sigma \in \Sigma'} x_s(\sigma) = 0 \tag{3.14}$$

**Theorem 1** *a)* SET-COVER $\leq$ DEPTH-$m$-DNF

  *b)* SET-COVER $\leq$ LENGTH-$m$-DNF

PROOF The proof is for any $R \in \{R_d, R_l\}$.

 Let $(S', \Sigma, m)$ any instance of SET-COVER.

 Let $T = (S, X, x, y)$ the Haussler data induced by $(S', \Sigma, m)$.

 We show: There exists a cover $\Sigma' \subseteq \Sigma$ of $S'$ with $|\Sigma'| \leq m$ iff there exists a $\theta \in \Theta$ such that $R(\theta) \leq m$ and $\forall s \in S \colon f_\theta(x_s) = y_s$.

 ($\Rightarrow$) Let $\Sigma' \subseteq \Sigma$ a cover of $S$ and $|\Sigma'| \leq m$.

 Let $V_0 = \varnothing$ and $V_1 = \Sigma'$ and $\theta = \{(V_0, V_1)\}$. Thus,

$$\forall x' \in X \colon \quad f_\theta(x') = \prod_{\sigma \in \Sigma'} x'(\sigma) \tag{3.15}$$

 On the one hand, $f(S') = 0$, by Lemma 2, and $f(1^\Sigma) = 1$, by definition of $f_\theta$. Thus, $\forall s \in S \colon f(x_s) = y_s$.

 On the other hand, $R(\theta) = m$.

 ($\Leftarrow$) Let $\theta \in \Theta$ such that $R(\theta) \leq m$ and $\forall s \in S \colon f_\theta(x_s) = y_s$.

 There exists a $(\Sigma_0, \Sigma_1) \in \theta$ such that $\Sigma_0 = \varnothing$, because $1 = y_1 = f_\theta(x_1) = f_\theta(1^\Sigma)$.

 Moreover

$$\forall s \in S' \colon \quad f(x_s) = 0$$

$$\Rightarrow \quad \forall s \in S' \colon \quad \bigvee_{(V_0, V_1) \in \theta} \prod_{v \in V_0} (1 - x_s(v)) \prod_{v \in V_1} x_s(v) = 0 \tag{3.16}$$

$$\Rightarrow \quad \forall s \in S' \; \forall (V_0, V_1) \in \theta \colon \quad \prod_{v \in V_0} (1 - x_s(v)) \prod_{v \in V_1} x_s(v) = 0 \tag{3.17}$$

Thus, for $(\varnothing, \Sigma_1) \in \theta$ in particular:

$$\forall s \in S' \colon \quad \prod_{\sigma \in \Sigma_1} x_s(\sigma) = 0 \tag{3.18}$$

And by virtue of Lemma 2:

$$\bigcup_{\sigma \in \Sigma_1} \sigma = S' \tag{3.19}$$

 Furthermore, $|\Sigma_1| \leq R(\theta) = m$.

### 3.1.6   Inference problem

For any $\theta \in \Theta$, the inference problem w.r.t. $f_\theta, L$ and any suitable $S', X', x'$ is solved by computing $f_\theta(x'_s)$ for any $s \in S'$, by Lemma 1.

### 3.1.7   Inference algorithm

Computing $f_\theta(x'_s)$ requires evaluating the form (3.1). The number of operations is bounded by $R_l(\theta) = O(|\theta| R_d(\theta)) = O(|\theta||V|)$.

## 3.2   Binary decision trees

### 3.2.1   Data

Throughout Section 3.2, we again consider some finite set $V \neq \varnothing$ and labeled data $T = (S, X, x, y)$ such that $X = \{0, 1\}^V$. Hence, $x \colon S \to \{0, 1\}^V$ and $y \colon S \to \{0, 1\}$.

### 3.2.2   Familiy of functions

**Definition 9** A tuple $(V, Y, D, D', d^*, E, \delta, v, y)$ is called a $V$-variate $Y$-valued *binary decision tree (BDT)* iff the following conditions hold:

- $V$ is finite and non-empty, called the set of *variables*

- $Y$ is finite, called the set of *values*

- $(D \cup D', E)$ is a finite, non-empty, directed binary tree

- $d^* \in D \cup D'$ is the unique root of this tree

- $\delta : E \to \{0, 1\}$

- Every $d \in D'$ is a leaf and every $d \in D$ has precisely two out-edges $e = (d, d'), e' = (d, d'')$ such that $\delta(e) = 0$ and $\delta(e) = 1$

- $v : D \to V$ assigning to each interior node a variable

- $y : D' \to Y$ assigning to each leaf a value

For any BDT $(V, Y, D, D', d^*, E, \delta, v, y)$, any $d \in D$ and any $j \in \{0, 1\}$, let $d_{\downarrow j} \in D \cup D'$ the unique node such that $e = (d, d_{\downarrow j}) \in E$ and $\delta(e) = j$.

Throughout Section 3.2, we consider $Y = \{0, 1\}$.

**Definition 10** For any BDT $\theta = (V, Y, D, D', d^*, E, \delta, v, y)$ and any $d \in D \cup D'$, the tuple $\theta[d] = (V, Y, D_2, D_2', d, E', \delta', v', y')$ with $(D_2 \cup D_2', E')$ the subtree of $(D \cup D', E)$ rooted at $d$ and with $\delta'$, $v'$ and $y'$ the restrictions of $\delta$, $v$ and $y$ to the subsets $D_2$, $D_2'$ and $E'$ is called the *binary decision subtree* of $\theta$ rooted at $d$.

**Lemma 3** *For any BDT $\theta = (V, Y, D, D', d^*, E, \delta, v, y)$ and any $d \in D \cup D'$, the binary decision subtree $\theta[d]$ is itself a $V$-variate $Y$-valued BDT.*

**Definition 11** For any BDT $\theta = (V, Y, D, D', d^*, E, \delta, v, y)$, the function defined by $\theta$ is the $f_\theta : \{0, 1\}^V \to Y$ such that $\forall x \in \{0, 1\}^V$:

$$f_\theta(x) = \begin{cases} y(d^*) & \text{if } D = \varnothing \\ (1 - x_{v(d^*)}) f_{\theta[d_{\downarrow 0}^*]}(x) + x_{v(d^*)} f_{\theta[d_{\downarrow 1}^*]}(x) & \text{otherwise} \end{cases} . \tag{3.20}$$

### 3.2.3   Regularization

**Definition 12** For any BDT $\theta = (V, Y, D, D', d^*, E, \delta, v, y)$, the *depth* of $\theta$ is the natural number $R(\theta) \in \mathbb{N}$ such that

$$R(\theta) = \begin{cases} 0 & \text{if } D = \varnothing \\ 1 + \max\{R(\theta[d_{\downarrow 0}^*]), R(\theta[d_{\downarrow 1}^*])\} & \text{otherwise} \end{cases} . \tag{3.21}$$

### 3.2.4   Loss function

We consider the loss function $L$ such that

$$\forall r \in \mathbb{R} \; \forall \hat{y} \in \{0, 1\}: \quad L(r, \hat{y}) = \begin{cases} 0 & r = \hat{y} \\ 1 & \text{otherwise} \end{cases} . \tag{3.22}$$

### 3.2.5 Learning problem

**Definition 13** For any $m \in \mathbb{N}$, the *bounded depth BDT problem (*DEPTH-$m$-BDT*)* w.r.t. $T$ and $m$ is to decide whether there exists a BDT $\theta = (V, Y, D, D', d^*, E, \delta, v, y')$ such that

$$R(\theta) \leq m \tag{3.23}$$

$$\forall s \in S: \quad f_\theta(x_s) = y_s \ . \tag{3.24}$$

**Theorem 2** EC-3 $\leq_p$ DEPTH-$m$-BDT

PROOF For any instance $(S', \Sigma)$ of EC-3 and the $n \in \mathbb{N}$ such that $|S'| = 3n$, we construct the instance of DEPTH-$m$-BDT such that

- $V = \Sigma$

- $S = S' \uplus \{0\}$

- $x : S \to \{0,1\}^\Sigma$ such that $x_0 = 0$ and

$$\forall s \in S' \ \forall \sigma \in \Sigma: \quad x_s(\sigma) = \begin{cases} 1 & \text{if } s \in \sigma \\ 0 & \text{otherwise} \end{cases} \tag{3.25}$$

- $y : S \to \{0,1\}$ such that $y_0 = 0$ and $\forall s \in S' : y_s = 1$.

- $m = n$

Next, we show that the instance EC-3 has a solution iff the instance of DEPTH-$m$-BDT has a solution.

($\Rightarrow$) Let $\Sigma' \subseteq \Sigma$ a solution to the instance of EC-3.

Consider any order on $\Sigma'$ and the bijection $\sigma' : [n] \to \Sigma'$ induced by this order. We show that the BDT $\theta$ depicted below solves the instance of DEPTH-$m$-BDT.



The BDT satisfies (3.23) as $R(\theta) = m$.

The BDT satisfies (3.24) because (i) $f_\theta(x_0) = 0 = y_0$ and (ii) at each of the $m$ interior nodes, three additional elements of $S'$ are mapped to 1. Thus, all $3m$ many elements $s \in S'$ are mapped to 1. That is $\forall s \in S' : f_\theta(x_s) = 1 = y_s$.

($\Leftarrow$) Let $\theta = (V, Y, D, D', d^*, E, \delta, \sigma, y')$ a $V$-variate BDT that solves the instance of DEPTH-$m$-BDT.

W.l.o.g., we assume, for any interior node $d \in D$, that $d_{\downarrow 1}$ is a leaf and $y'(d_{\downarrow 1}) = 1$. Hence, $\theta$ is of the form depicted below.

Therefore,

$$\forall x \in X: \quad f_\theta(x) = \begin{cases} 1 & \text{if } \exists j \in [N]: x(\sigma_j) = 1 \\ 0 & \text{otherwise} \end{cases} , \qquad (3.26)$$

and thus,

$$\forall s \in S: \quad f_\theta(x_s) = \begin{cases} 1 & \text{if } \exists j \in [N]: s \in \sigma_j \\ 0 & \text{otherwise} \end{cases} , \qquad (3.27)$$

by definition of $x$ in (3.25). Consequently,

$$\sigma(D) = \bigcup_{j=0}^{N-1} \sigma_j = S' , \qquad (3.28)$$

by definition of $y$ such that $y(S') = 1$.

Moreover,

$$3m = |S'| \overset{(3.28)}{=} \left| \bigcup_{j=0}^{N-1} \sigma_j \right| \leq \sum_{j=0}^{N-1} |\sigma_j| = \sum_{j=0}^{N-1} 3 = 3N \overset{(3.23)}{\leq} 3m .$$

Therefore,

$$\forall \{j,l\} \in \binom{[N]}{2}: \quad \sigma_k \cap \sigma_l = \varnothing , \qquad (3.29)$$

by Lemma 15. Hence,

$$\sigma(D) = \cup_{j=1}^{N-1} \sigma_j$$

is a solution to the instance of EC-3 defined by $(S', \Sigma)$, by (3.28) and (3.29).

**Lemma 4** DEPTH-$m$-BDT *is* NP-*complete.*

PROOF DEPTH-$m$-BDT $\in$ NP, as solutions can be verified efficiently.
DEPTH-$m$-BDT is NP-hard, by Theorem 2 and Lemma 13.

**Corollary 1** *Discriminative learning of BDTs is* NP-*hard.*

### 3.2.6  Learning algorithm

We introduce a popular heuristic for constructing BDTs.

**Definition 14** For any finite set $S$, any functions $y, y' : S \to Y$ and the set

$$G(y, y') = \left\{ \{s, t\} \in \binom{S}{2} \;\middle|\; \begin{array}{l} (y(s) = y(t) \wedge y'(s) \neq y'(t)) \\ \vee \; (y(s) \neq y(t) \wedge y'(s) = y'(t)) \end{array} \right\} \tag{3.30}$$

the function $\Delta : Y^S \times Y^S \to [0, 1]$ such that

$$\Delta(y, y') = \frac{|G(y, y')|}{\binom{|S|}{2}} \tag{3.31}$$

is called the *impurity pseudometric* of functions $Y^S$.

**Exercise 1** *(i) Show that the impurity pseudometric is indeed a pseudometric.*
*(ii) Show that the impurity pseudometric is not necessarily a metric.*

**Definition 15** For any of the given labeled data $T = (S, X, x, y)$, the *greedy impurity minimizing algorithm* constructs a BDT $\theta = (V, Y = \{0, 1\}, D, D', d^*, E, \delta, v, y')$ by initializing $D = D' = E = \varnothing$ and executing the procedure grow$(d^*, S)$.

| grow$(d, S')$ | |
| :--- | ---: |
| if $\forall s, s' \in S' : y_s = y_{s'}$ | |
| $\quad\quad D' := D' \cup \{d\}$ | add leaf $d$ |
| $\quad\quad$ choose $s \in S'$ | |
| $\quad\quad y'_d := y_s$ | |
| else | |
| $\quad\quad D := D \cup \{d\}$ | add interior node $d$ |
| $\quad\quad$ choose $\hat{v} \in \min\limits_{v \in V} \Delta(y, x_v)$ | |
| $\quad\quad$ for $j \in \{0, 1\}$ | |
| $\quad\quad\quad$ let $d_j \notin D \cup D'$ | |
| $\quad\quad\quad$ grow$(d_j, \{s \in S' \mid x_s(\hat{v}) = j\})$ | recurse |
| $\quad\quad\quad e := (d, d_j)$ | |
| $\quad\quad\quad E := E \cup \{e\}$ | |
| $\quad\quad\quad \delta_e := j$ | |

### 3.2.7 Inference problem

For any BDT $\theta = (V, Y, D, D', d^*, E, \delta, v, y')$, the inference problem w.r.t. $f_\theta, L$ and any suitable $S', X', x'$ is solved by computing $f_\theta(x'_s)$ for any $s \in S'$, by Lemma 1.

### 3.2.8 Inference algorithm

The inference problem is solved by evaluating the form (3.20). The time complexity is $O\left(R(\theta)\right)$.

## 3.3 Linear functions

### 3.3.1 Data

Throughout Section 3.3, we consider real attributes. More specifically, we consider some finite set $V \neq \varnothing$ and labeled data $T = (S, X, x, y)$ with $X = \mathbb{R}^V$. Hence, $x : S \to \mathbb{R}^V$ and $y : S \to \{0, 1\}$.

### 3.3.2 Familiy of functions

Throughout Section 3.3, we consider linear functions. More specifically, we consider $\Theta = \mathbb{R}^V$ and $f : \Theta \to \mathbb{R}^X$ such that

$$\forall \theta \in \Theta \; \forall \hat{x} \in X : \quad f_\theta(\hat{x}) = \langle \theta, \hat{x} \rangle \;. \tag{3.32}$$

### 3.3.3   Probabilistic model

**Random variables**

- For any $s \in S$, let $X_s$ be a random variable whose realization is a vector $x_s \in \mathbb{R}^V$, called the *attribute vector* of $s$

- For any $s \in S$, let $Y_s$ be a random variable whose realization is a binary number $y_s \in \{0, 1\}$, called the *label* of $s$

- For any $v \in V$, let $\Theta_v$ be a random variable whose realization is a real number $\theta_v \in \mathbb{R}$, called a *parameter*

**Conditional independence assumptions**

We assume a probability distribution that factorizes according to the Bayesian net depicted below.



**Factorization**

- Firstly:

$$P(X, Y, \Theta) = \prod_{s \in S} P(Y_s \mid X_s, \Theta) P(X_s) \prod_{v \in V} P(\Theta_v) \tag{3.33}$$

- Secondly:

$$
\begin{aligned}
P(\Theta \mid X, Y) &= \frac{P(X, Y, \Theta)}{P(X, Y)} \\
&= \frac{P(Y \mid X, \Theta)\, P(X)\, P(\Theta)}{P(X, Y)} \\
&\propto P(Y \mid X, \Theta)\, P(\Theta) \\
&= \prod_{s \in S} P(Y_s \mid X_s, \Theta) \prod_{v \in V} P(\Theta_v)
\end{aligned} \tag{3.34}
$$

**Forms**

We consider:

- The *logistic distribution*

$$\forall s \in S : \qquad p_{Y_s \mid X_s, \Theta}(1) = \frac{1}{1 + 2^{-f_\theta(x_s)}} \tag{3.35}$$

- A $\sigma \in \mathbb{R}^+$ and the *normal distribution*:

$$\forall v \in V : \qquad p_{\Theta_v}(\theta_v) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\theta_v^2 / 2\sigma^2} \tag{3.36}$$

### 3.3.4 Learning problem

**Lemma 5 (Logistic regression)** *Estimating maximally probable parameters $\theta$, given attributes $x$ and labels $y$, i.e.,*

$$\operatorname*{argmax}_{\theta \in \mathbb{R}^m} \quad p_{\Theta|X,Y}(\theta, x, y)$$

*is identical to the supervised learning problem w.r.t. $L$, $R$ and $\lambda$ such that*

$$\forall r \in \mathbb{R} \; \forall \hat{y} \in \{0,1\}: \quad L(r, \hat{y}) = -\hat{y}r + \log\left(1 + 2^r\right) \tag{3.37}$$

$$\forall \theta \in \Theta: \quad R(\theta) = \|\theta\|_2^2 \tag{3.38}$$

$$\lambda = \frac{\log e}{2\sigma^2} \tag{3.39}$$

PROOF Firstly,

$$\operatorname*{argmax}_{\theta \in \mathbb{R}^m} \quad p_{\Theta|X,Y}(\theta, x, y)$$

$$\overset{(3.34)}{=} \operatorname*{argmax}_{\theta \in \mathbb{R}^m} \quad \prod_{s \in S} p_{Y_s|X_s,\Theta}(y_s, x_s, \theta) \prod_{v \in V} p_{\Theta_v}(\theta_v)$$

$$= \operatorname*{argmax}_{\theta \in \mathbb{R}^m} \quad \sum_{s \in S} \log p_{Y_s|X_s,\Theta}(y_s, x_s, \theta) + \sum_{v \in V} \log p_{\Theta_v}(\theta_v) \tag{3.40}$$

Substituting in (3.40) the linearization

$$\begin{aligned} &\log p_{Y_s|X_s,\Theta}(y_s, x_s, \theta) \\ =\;& y_s \log p_{Y_s|X_s,\Theta}(1, x_s, \theta) + (1 - y_s) \log p_{Y_s|X_s,\Theta}(0, x_s, \theta) \\ =\;& y_s \log \frac{p_{Y_s|X_s,\Theta}(1, x_s, \theta)}{p_{Y_s|X_s,\Theta}(0, x_s, \theta)} + \log p_{Y_s|X_s,\Theta}(0, x_s, \theta) \end{aligned} \tag{3.41}$$

as well as (3.35) and (3.36) yields the form (3.42) below that is called the instance of the $l_2$-regularized *logistic regression problem* with respect to $x$, $y$ and $\sigma$.

$$\operatorname*{argmin}_{\theta \in \mathbb{R}^m} \quad \sum_{s \in S} \left( -y_s \langle \theta, x_s \rangle + \log\left(1 + 2^{\langle \theta, x_s \rangle}\right) \right) + \frac{\log e}{2\sigma^2} \|\theta\|_2^2 \tag{3.42}$$

**Exercise 2** *a) Derive (3.42) from (3.40) using (3.41), (3.35) and (3.36)*
*b) Is the objective function of (3.42) convex?*

### 3.3.5 Inference problem

**Lemma 6** *Estimating maximally probable labels $y$, given attributes $x'$ and parameters $\theta$, i.e.,*

$$\operatorname*{argmax}_{y \in \{0,1\}^S} \quad p_{Y|X,\Theta}(y, x', \theta) \tag{3.43}$$

*is identical to the inference problem w.r.t. $f$ and $L$. It has the solution*

$$\forall s \in S' : \quad y_s = \begin{cases} 1 & \text{if } f_\theta(x'_s) > 0 \\ 0 & \text{otherwise} \end{cases} \tag{3.44}$$

PROOF  Firstly,

$$\underset{y\in\{0,1\}^{S'}}{\mathrm{argmax}}\quad p_{Y|X,\Theta}(y,x',\theta)$$

$$= \underset{y\in\{0,1\}^{S'}}{\mathrm{argmax}}\quad \prod_{s\in S'} p_{Y_s|X_s,\Theta}(y_s,x'_s,\theta)$$

$$= \underset{y\in\{0,1\}^{S'}}{\mathrm{argmax}}\quad \sum_{s\in S'} \log p_{Y_s|X_s,\Theta}(y_s,x'_s,\theta)$$

$$= \underset{y\in\{0,1\}^{S'}}{\mathrm{argmax}}\quad \sum_{s\in S'} \left( y_s \log \frac{p_{Y_s|X_s,\Theta}(1,x'_s,\theta)}{p_{Y_s|X_s,\Theta}(0,x'_s,\theta)} + \log p_{Y_s|X_s,\Theta}(0,x'_s,\theta) \right)$$

$$= \underset{y\in\{0,1\}^{S'}}{\mathrm{argmin}}\quad \sum_{s\in S'} \left( -y_s f_\theta(x'_s) + \log\left(1 + 2^{f_\theta(x'_s)}\right) \right)$$

$$= \underset{y\in\{0,1\}^{S'}}{\mathrm{argmin}}\quad \sum_{s\in S'} L(f_\theta(x'_s), y_s) \ .$$

Secondly,

$$\min_{y\in\{0,1\}^{S'}} \sum_{s\in S'} \left( -y_s f_\theta(x'_s) + \log\left(1 + 2^{f_\theta(x'_s)}\right) \right) = \sum_{s\in S'} \max_{y_s\in\{0,1\}} y_s f_\theta(x'_s) \ .$$

### 3.3.6   Inference algorithm

The inference problem is solved by computing independently for each $s \in S'$ the label

$$y_s = \begin{cases} 1 & \text{if } \langle \theta, x'_s \rangle > 0 \\ 0 & \text{otherwise} \end{cases} \ . \tag{3.45}$$

The time complexity is $O(|V||S'|)$.

# Chapter 4

# Semi-supervised and unsupervised learning

## 4.1 Intuition

So far, we have considered learning problems w.r.t. labeled data $(S, X, x, y)$ where, for every $s \in S$, a label $y_s \in \{0, 1\}$ is given, and inference problems w.r.t. unlabeled data $(S', X', x)$ where no label is given and every combination of labels $y' : S \to \{0, 1\}$ is a feasible solution.

Next, we consider learning problems where not every label is given and inference problems where not every combination of labels is feasible. Unlike before, the data we look at in both problems coincides, consisting of tuples $(S, X, x, \mathcal{Y})$ where $\mathcal{Y} \subseteq \{0, 1\}^S$ is a set of feasible labelings. In particular, $\mathcal{Y} = \{0, 1\}^S$ is the special case of unlabeled data, and $|\mathcal{Y}| = 1$ is the special case of labeled data. Non-trivial choices of $\mathcal{Y}$ allow us to express problems of learning and inferring finite structures such as maps (Chapter 5), equivalence relations (Chapter 6) and orders (Chapter 7).

## 4.2 Definition

**Definition 16** For any $S \neq \varnothing$ finite, called a set of *samples*, any $X \neq \varnothing$, called an *attribute space*, any $x : S \to X$ and any $\varnothing \neq \mathcal{Y} \subseteq \{0, 1\}^S$, called a set of *feasible labelings*, the tuple $T = (S, X, x, \mathcal{Y})$ is called *constrained data*.

**Definition 17** For any constrained data $T = (S, X, x, \mathcal{Y})$, any $\Theta \neq \varnothing$ and family of functions $f : \Theta \to \mathbb{R}^X$, any $R : \Theta \to \mathbb{R}_0^+$, called a *regularizer*, any $L : \mathbb{R} \times \{0, 1\} \to \mathbb{R}_0^+$, called a *loss function* and any $\lambda \in \mathbb{R}_0^+$, called a *regularization parameter*, the instance of the *learning and inference problem* w.r.t. $T, \Theta, f, R, L$ and $\lambda$ is defined as

$$\min_{y \in \mathcal{Y}} \inf_{\theta \in \Theta} \quad \lambda R(\theta) + \frac{1}{|S|} \sum_{s \in S} L(f_\theta(x_s), y_s) \tag{4.1}$$

The special case of one-elementary $\mathcal{Y} = \{y\}$ is called the *supervised learning problem*.
The special case of one-elementary $\Theta = \{\hat{\theta}\}$ written below is called the *inference problem*.

$$\min_{y \in \mathcal{Y}} \quad \sum_{s \in S} L(f_\theta(x_s), y_s) \tag{4.2}$$

The mixed integer optimization problem (4.1) is an topic of machine learning research beyond the scope of this lecture. Special cases of the binary optimization problem (4.2) have been studied intensively and are discussed in the following chapters.

# Chapter 5

# Classifying

## 5.1  Maps

For any finite set $A \neq \varnothing$ whose elements we seek to classify and any finite set $B \neq \varnothing$ of class labels, we are interested in *maps* $\varphi : A \to B$ that assign to every element $a \in A$ precisely one class label $\varphi(a) \in B$. Maps are precisely those subsets of $\varphi \subseteq A \times B$ that satisfy

$$\forall a \in A \; \exists b \in B : \quad (a, b) \in \varphi \tag{5.1}$$
$$\forall a \in A \; \forall b, b' \in B : \quad (a, b) \in \varphi \wedge (a, b') \in \varphi \Rightarrow b = b' \; . \tag{5.2}$$

They are characterized by those functions $y : A \times B \to \{0, 1\}$ that satisfy

$$\forall a \in A : \quad \sum_{b \in B} y_{ab} = 1 \; . \tag{5.3}$$

We reduce the problem of learning and inferring maps to the problem of learning and inferring decisions, by choosing constrained data with

$$S = A \times B \tag{5.4}$$

$$\mathcal{Y} = \left\{ y \colon A \times B \to \{0, 1\} \;\middle|\; \forall a \in A \colon \sum_{b \in B} y_{ab} = 1 \right\} \; . \tag{5.5}$$

## 5.2  Linear functions

### 5.2.1  Data

Throughout Section 5.2, we consider some finite set $V \neq \varnothing$ and constrained data $(S, X, x, \mathcal{Y})$ with $S = A \times B$ as in (5.4), $X = B \times \mathbb{R}^V$, and $\mathcal{Y}$ as in (5.5). More specifically, we assume that, for any $(a, b) \in A \times B$, the class label $b$ is the first attribute of $(a, b)$, i.e.,

$$\forall a \in A \; \forall b \in B \; \exists \hat{x} \in \mathbb{R}^V : \quad x_{ab} = (b, \hat{x}) \tag{5.6}$$

As a special case, we consider labeled data where we are given just one $\mathcal{Y} = \{y\}$ with $y$ satisfying the constraints (5.3).

### 5.2.2  Familiy of functions

Throughout Section 5.2, we consider linear functions. More specifically, we consider $\Theta = \mathbb{R}^{B \times V}$ and $f : \Theta \to \mathbb{R}^X$ such that

$$\forall \theta \in \Theta \; \forall b \in B \; \forall \hat{x} \in \mathbb{R}^V : \quad f_\theta((b, \hat{x})) = \sum_{v \in V} \theta_{bv} \, \hat{x}_v = \langle \theta_{b\cdot}, \hat{x} \rangle \; . \tag{5.7}$$
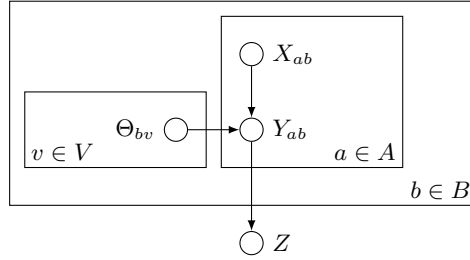
### 5.2.3   Probabilistic model

**Random variables**

- For any $(a, b) \in A \times B$, let $X_{ab}$ be a random variable whose realization is a vector $x_{ab} \in B \times \mathbb{R}^V$, called the *attribute vector* of $(a, b)$.

- For any $(a, b) \in A \times B$, let $Y_{ab}$ be a random variable whose realization is a binary number $y_{ab} \in \{0, 1\}$, called the *decision* of classifying $a$ as $b$

- For any $b \in B$ and any $v \in V$, let $\Theta_{bv}$ be a random variable whose realization is a real number $\theta_{bv} \in \mathbb{R}$, called a *parameter*

- Let $Z$ be a random variable whose realization is a subset $z \subseteq \{0, 1\}^{A \times B}$. For multiple label classification, we are interested in $z = \mathcal{Y}$, the set of the characteristic functions of all maps from $A$ to $B$.

**Conditional independence assumptions**

We assume a probability distribution that factorizes according to Bayesian net depicted below.



**Factorization**

These conditional independence assumptions imply the following factorizations:

- Firstly:

$$P(X, Y, Z, \Theta) = P(Z \mid Y) \prod_{(a,b) \in A \times B} P(Y_{ab} \mid X_{ab}, \Theta) \prod_{(b,v) \in B \times V} P(\Theta_{bv}) \prod_{(a,b) \in A \times B} P(X_{ab}) \tag{5.8}$$

- Secondly:

$$
\begin{aligned}
P(\Theta \mid X, Y, Z) &= \frac{P(X, Y, Z, \Theta)}{P(X, Y, Z)} \\
&= \frac{P(Z \mid Y)\, P(Y \mid X, \Theta)\, P(X)\, P(\Theta)}{P(Z \mid X, Y)\, P(X, Y)} \\
&= \frac{P(Z \mid Y)\, P(Y \mid X, \Theta)\, P(X)\, P(\Theta)}{P(Z \mid Y)\, P(X, Y)} \\
&= \frac{P(Y \mid X, \Theta)\, P(X)\, P(\Theta)}{P(X, Y)} \\
&\propto P(Y \mid X, \Theta)\, P(\Theta) \\
&= \prod_{(a,b) \in A \times B} P(Y_{ab} \mid X_{ab}, \Theta) \prod_{(b,v) \in B \times V} P(\Theta_{bv}) \tag{5.9}
\end{aligned}
$$

- Thirdly,

$$
\begin{aligned}
P(Y \mid X, Z, \theta) &= \frac{P(X, Y, Z, \Theta)}{P(X, Z, \Theta)} \\
&= \frac{P(Z \mid Y)\,P(Y \mid X, \Theta)\,P(X)\,P(\Theta)}{P(X, Z, \Theta)} \\
&\propto P(Z \mid Y)\,P(Y \mid X, \Theta) \\
&= P(Z \mid Y) \prod_{(a,b) \in A \times B} P(Y_{ab} \mid X_{ab}, \Theta)
\end{aligned}
\tag{5.10}
$$

**Forms**

Here, we consider:

- The *logistic distribution*

$$
\forall (a,b) \in A \times B: \qquad p_{Y_{ab} \mid X_{ab}, \Theta}(1) = \frac{1}{1 + 2^{-f_\theta(x_{ab})}}
\tag{5.11}
$$

- A $\sigma \in \mathbb{R}^+$ and the *normal distribution*:

$$
\forall (b,v) \in B \times V: \qquad p_{\Theta_{bv}}(\theta_{bv}) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\theta_{bv}^2 / 2\sigma^2}
\tag{5.12}
$$

- A uniform distribution on a subset:

$$
\forall z \subseteq \{0,1\}^{A \times B}: \quad p_{Z \mid Y}(z) \propto \begin{cases} 1 & \text{if } y \in z \\ 0 & \text{otherwise} \end{cases}
\tag{5.13}
$$

Note that $p_{Z \mid Y}(\mathcal{Y})$ is non-zero iff the relation $y^{-1}(1) \subseteq A \times B$ is a map.

### 5.2.4 Learning problem

**Lemma 7** *Estimating maximally probable parameters $\theta$, given attributes $x$ and decisions $y$, i.e.,*

$$
\underset{\theta \in \mathbb{R}^{B \times V}}{\mathrm{argmax}} \quad p_{\Theta \mid X, Y}(\theta, x, y)
$$

*is identical to the supervised learning problem w.r.t. $L$, $R$ and $\lambda$ such that*

$$
\forall r \in \mathbb{R} \ \forall \hat{y} \in \{0,1\}: \quad L(r, \hat{y}) = -\hat{y} r + \log\left(1 + 2^r\right)
\tag{5.14}
$$

$$
\forall \theta \in \Theta: \qquad R(\theta) = \|\theta\|_2^2
\tag{5.15}
$$

$$
\lambda = \frac{\log e}{2\sigma^2}
\tag{5.16}
$$

*Moreover, this problem separates into $|B|$ independent supervised learning problems, each w.r.t. parameters in $\mathbb{R}^V$, with $L$ and $\lambda$ as above, and with*

$$
\forall \theta' \in \mathbb{R}^V: \qquad R'(\theta') = \|\theta'\|_2^2
\tag{5.17}
$$

PROOF Analogous to the case of binary classification from Section 3.3, we now obtain:

$$
\begin{aligned}
&\underset{\theta \in \mathbb{R}^{B \times V}}{\mathrm{argmax}} \quad p_{\Theta \mid X, Y}(\theta, x, y) \\
&= \underset{\theta \in \mathbb{R}^{B \times V}}{\mathrm{argmin}} \sum_{(a,b) \in A \times B} \left( -y_{ab} f_\theta(x_{ab}) + \log\left(1 + 2^{f_\theta(x_{ab})}\right) \right) + \frac{\log e}{2\sigma^2} \|\theta\|_2^2 \ .
\end{aligned}
\tag{5.18}
$$

Consider the unique $x' : A \times B \to \mathbb{R}^V$ such that, for any $(a, b) \in A \times B$, we have $x_{ab} = (b, x'_{ab})$.

Problem (5.18) separates into $|B|$ many $l_2$-regularized logistic regression problems, one for each $b \in B$, because

$$\min_{\theta \in \mathbb{R}^{B \times V}} \sum_{(a,b) \in A \times B} \left( -y_{ab} \langle \theta_{b\cdot}, x'_{ab} \rangle + \log \left( 1 + 2^{\langle \theta_{b\cdot}, x'_{ab} \rangle} \right) \right) + \frac{\log e}{2\sigma^2} \|\theta\|_2^2$$

$$= \min_{\theta \in \mathbb{R}^{B \times V}} \sum_{b \in B} \left( \sum_{a \in A} \left( -y_{ab} \langle \theta_{b\cdot}, x'_{ab} \rangle + \log \left( 1 + 2^{\langle \theta_{b\cdot}, x'_{ab} \rangle} \right) \right) + \frac{\log e}{2\sigma^2} \|\theta_{b\cdot}\|_2^2 \right)$$

$$= \sum_{b \in B} \min_{\theta_{b\cdot} \in \mathbb{R}^V} \left( \sum_{a \in A} \left( -y_{ab} \langle \theta_{b\cdot}, x'_{ab} \rangle + \log \left( 1 + 2^{\langle \theta_{b\cdot}, x'_{ab} \rangle} \right) \right) + \frac{\log e}{2\sigma^2} \|\theta_{b\cdot}\|_2^2 \right) \quad .$$

### 5.2.5   Inference problem

**Lemma 8** *For any constrained data as defined above, any $\theta \in \mathbb{R}^{B \times V}$ and any $\hat{y} : A \times B \to \{0, 1\}$, $\hat{y}$ is a solution to the inference problem*

$$\min_{y \in \mathcal{Y}} \sum_{(a,b) \in A \times B} L(f_\theta(x_{ab}), y_{ab}) \tag{5.19}$$

*iff there exists an $\varphi : A \to B$ such that*

$$\forall a \in A : \quad \varphi(a) \in \max_{b \in B} \langle \theta_{b\cdot}, x'_{ab} \rangle \tag{5.20}$$

*and*

$$\forall (a, b) \in A \times B : \quad \hat{y}_{ab} = 1 \Leftrightarrow \varphi(a) = b \quad . \tag{5.21}$$

PROOF

$$\sum_{(a,b) \in A \times B} L(f_\theta(x_{ab}), y_{ab})$$

$$= \sum_{(a,b) \in A \times B} \left( L(f_\theta(x_{ab}), 1) \, y_{ab} + L(f_\theta(x_{ab}), 0) \, (1 - y_{ab}) \right)$$

$$= \sum_{(a,b) \in A \times B} \left( L(f_\theta(x_{ab}), 1) - L(f_\theta(x_{ab}), 0) \right) y_{ab} + \text{const.}$$

$$= \sum_{(a,b) \in A \times B} \left( -f_\theta(x_{ab}) \right) y_{ab} \qquad\qquad \text{by (5.14)}$$

$$= \sum_{(a,b) \in A \times B} \left( -\langle \theta_{b\cdot}, x'_{ab} \rangle \right) y_{ab} \qquad\qquad x_{ab} = (b, x'_{ab})$$

$$= \sum_{a \in A} \sum_{b \in B} \left( -\langle \theta_{b\cdot}, x'_{ab} \rangle \right) y_{ab}$$

### 5.2.6   Inference algorithm

The inference problem is solved by solving (5.20) independently for each $a \in A$. The time complexity is $O(|A||B||V|)$.

# Chapter 6

# Clustering

## 6.1 Partitions and equivalence relations

Throughout this chapter, we consider some finite set $A \neq \varnothing$ that we seek to partition into subsets. Hence, our feasible solutions are the *partitions* of $A$. A partition $\Pi$ of $A$ is a collection $\Pi \subseteq 2^A$ of non-empty and pairwise disjoint subsets of $A$ whose union is $A$.

The partitions of $A$ are characterized by the equivalence relations on $A$. An equivalence relation $\equiv$ on $A$ is a binary relation $\equiv \subseteq A \times A$ that is reflexive, symmetric and transitive. For any partition $\Pi$ of $A$, the equivalence relation $\equiv_\Pi$ induced by $\Pi$ is such that

$$\forall a, a' \in A: \quad a \equiv_\Pi a' \Leftrightarrow \exists U \in \Pi: a \in U \land a' \in U \ . \tag{6.1}$$

In turn, the equivalence relations on $A$ are characterized by those $y : \binom{A}{2} \to \{0, 1\}$ that satisfy

$$\forall \{a, b, c\} \in \binom{A}{3}: \quad y_{\{a,b\}} + y_{\{b,c\}} - 1 \leq y_{\{a,c\}} \ . \tag{6.2}$$

For any partition $\Pi$ of $A$, consider the $y^\Pi : \binom{A}{2} \to \{0, 1\}$ such that

$$\forall \{a, a'\} \in \binom{A}{2}: \quad y^\Pi_{\{a,a'\}} = 1 \Leftrightarrow \exists U \in \Pi: a \in U \land a' \in U \ . \tag{6.3}$$

Herein, for any $\{a, b\} \in \binom{A}{2}$, the decision $y_{\{a,b\}} = 1$ means that $a$ and $b$ are in the same cluster, whereas the decision $y_{\{a,b\}} = 0$ means that $a$ and $b$ in distinct clusters.

We reduce the problem of learning and inferring partitions to the problem of learning and inferring decisions by choosing constrained data with

$$S = \binom{A}{2} \tag{6.4}$$

$$\mathcal{Y} = \left\{ y : \binom{A}{2} \to \{0, 1\} \ \middle| \ (6.2) \right\} \ . \tag{6.5}$$

One can think of the set $S$ as the set of edges in the complete graph with nodes $A$ and without self-edges.

## 6.2 Correlation clustering

### 6.2.1 Data

Throughout Section 6.2, we consider some finite set $V \neq \varnothing$ and constrained data $(S, X, x, \mathcal{Y})$ with $S = \binom{A}{2}$ as in (6.4), $X = \mathbb{R}^V$ and $\mathcal{Y}$ as in (6.5). As a special case, we consider labeled data, i.e., just one $\mathcal{Y} = \{y\}$ with $y$ satisfying the constraints (6.2).

## 6.2.2   Familiy of functions

Throughout Section 5.2, we consider linear functions. More specifically, we consider $\Theta = \mathbb{R}^V$ and $f : \Theta \to \mathbb{R}^X$ such that

$$\forall \theta \in \Theta \ \forall \hat{x} \in \mathbb{R}^V : \quad f_\theta(\hat{x}) = \langle \theta, \hat{x} \rangle \ . \tag{6.6}$$
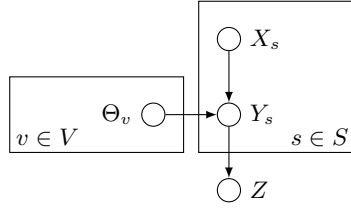
## 6.2.3   Probabilistic model

### Random variables

- For any $\{a, a'\} \in S$, let $X_{\{a,a'\}}$ be a random variable whose realization is a vector $x_{\{a,a'\}} \in \mathbb{R}^V$, called the *attribute vector* of the pair $\{a, a'\}$.

- For any $\{a, a'\} \in S$, let $Y_{\{a,a'\}}$ be a random variable whose realization is a binary number $y_{\{a,a'\}} \in \{0, 1\}$, called the *decision* of joining $a$ and $a'$ in the same cluster.

- For any $v \in V$, let $\Theta_v$ be a random variable whose realization is a real number $\theta_v \in \mathbb{R}$, called a *parameter*

- Let $Z$ be a random variable whose realization is a subset $z \subseteq \{0, 1\}^S$. For clustering, we are interested in $z = \mathcal{Y}$, a characterization of all partitions of $A$.

### Conditional independence assumptions

We assume a probability distribution that factorizes according to the Bayesian net depicted below.



### Factorization

These conditional independence assumptions imply the following factorizations:

- Firstly:

$$P(X, Y, Z, \Theta) = P(Z \mid Y) \prod_{s \in S} P(Y_s \mid X_s, \Theta) \prod_{s \in S} P(X_s) \prod_{v \in V} P(\Theta_v) \tag{6.7}$$

- Secondly:

$$\begin{aligned}
P(\Theta \mid X, Y, Z) &= \frac{P(X, Y, Z, \Theta)}{P(X, Y, Z)} \\
&= \frac{P(Z \mid Y) \, P(Y \mid X, \Theta) \, P(X) \, P(\Theta)}{P(Z \mid X, Y) \, P(X, Y)} \\
&= \frac{P(Z \mid Y) \, P(Y \mid X, \Theta) \, P(X) \, P(\Theta)}{P(Z \mid Y) \, P(X, Y)} \\
&= \frac{P(Y \mid X, \Theta) \, P(X) \, P(\Theta)}{P(X, Y)} \\
&\propto P(Y \mid X, \Theta) \, P(\Theta) \\
&= \prod_{s \in S} P(Y_s \mid X_s, \Theta) \prod_{v \in V} P(\Theta_v) \tag{6.8}
\end{aligned}$$

- Thirdly,

$$
\begin{aligned}
P(Y \mid X, Z, \theta) &= \frac{P(X, Y, Z, \Theta)}{P(X, Z, \Theta)} \\
&= \frac{P(Z \mid Y)\, P(Y \mid X, \Theta)\, P(X)\, P(\Theta)}{P(X, Z, \Theta)} \\
&\propto P(Z \mid Y)\, P(Y \mid X, \Theta) \\
&= P(Z \mid Y) \prod_{s \in S} P(Y_s \mid X_s, \Theta) \qquad (6.9)
\end{aligned}
$$

**Forms**

Here, we consider:

- The *logistic distribution*

$$
\forall s \in S : \qquad p_{Y_s \mid X_s, \Theta}(1) = \frac{1}{1 + 2^{-f_\theta(x_s)}} \qquad (6.10)
$$

- A $\sigma \in \mathbb{R}^+$ and the *normal distribution*:

$$
\forall v \in V : \qquad p_{\Theta_v}(\theta_v) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\theta_v^2/2\sigma^2} \qquad (6.11)
$$

- A uniform distribution on a subset:

$$
\forall z \subseteq \{0,1\}^S : \quad p_{Z \mid Y}(z) \propto \begin{cases} 1 & \text{if } y \in z \\ 0 & \text{otherwise} \end{cases} \qquad (6.12)
$$

Note that $p_{Z\mid Y}(\mathcal{Y})$ is non-zero iff $y$ induces an equivalence relations and hence a partition of $A$.

### 6.2.4 Learning problem

**Corollary 2** *Estimating maximally probable parameters $\theta$, given attributes $x$ and labels $y$, i.e.,*

$$
\underset{\theta \in \mathbb{R}^m}{\operatorname{argmax}} \quad p_{\Theta \mid X, Y}(\theta, x, y)
$$

*is identical to the supervised learning problem w.r.t. $L$, $R$ and $\lambda$ such that*

$$
\forall r \in \mathbb{R} \; \forall \hat{y} \in \{0,1\} : \quad L(r, \hat{y}) = -\hat{y}r + \log(1 + 2^r) \qquad (6.13)
$$

$$
\forall \theta \in \Theta : \qquad R(\theta) = \|\theta\|_2^2 \qquad (6.14)
$$

$$
\lambda = \frac{\log e}{2\sigma^2} \qquad (6.15)
$$

### 6.2.5 Inference problem

**Corollary 3** *For any constrained data as defined above and any $\theta \in \mathbb{R}^V$, the inference problem takes the form of* CORRELATION-CLUSTERING

$$
\min_{y: \binom{A}{2} \to \{0,1\}} \sum_{\{a,a'\} \in S} \left( -\langle \theta, x_{\{a,a'\}} \rangle \right) y_{\{a,a'\}} \qquad (6.16)
$$

$$
\text{subject to} \quad \forall \{a,b,c\} \in \binom{A}{3} : \quad y_{\{a,b\}} + y_{\{b,c\}} - 1 \le y_{\{a,c\}} \; . \qquad (6.17)
$$

CORRELATION-CLUSTERING has been studied intensively, notably by Bansal et al. (2004) and Demaine et al. (2006).

**Lemma 9 (Bansal et al. (2004))** CORRELATION-CLUSTERING *is* NP-*hard.*

Bansal et al. (2004) establish NP-hardness of CORRELATION-CLUSTERING by a reduction of $k$-TERMINAL-CUT whose NP-hardness is an important result of Dahlhaus et al. (1994).

### 6.2.6   Inference algorithm

Below, we discuss three local search algorithms for the correlation clustering problem, i.e., heuristics whose outputs need not be optimal.

For simplicity, we define $c : S \to \mathbb{R}$ such that

$$\forall \{a, a'\} \in S: \quad c_{\{aa'\}} = -\langle \theta, x_{\{a,a'\}} \rangle \ , \tag{6.18}$$

and write the objective function $\varphi : \{0,1\}^S \to \mathbb{R}$ such that

$$\forall y \in \{0,1\}^S: \quad \varphi(y) = \sum_{\{a,a'\} \in S} c_{\{a,a'\}} \, y_{\{a,a'\}} \tag{6.19}$$

**Greedy joining**

The greedy joining algorithm starts from any initial partition and searches for partitions with lower objective value by joining pairs of clusters recursively. As clusters can only grow and the number of clusters decreases by precisely one in every step, one typically starts from the partition $\Pi_0$ of $A$ into one-elementary subsets.

**Definition 18** For any partition $\Pi$ of $A$, and any $B, C \in \Pi$, let $\mathrm{join}_{BC}[\Pi]$ the partition of $A$ obtained by joining the sets $B$ and $C$ in $\Pi$, i.e.

$$\mathrm{join}_{BC}[\Pi] = (\Pi \setminus \{B, C\}) \cup \{B \cup C\} \ . \tag{6.20}$$

**Algorithm 1** The greedy joining algorithm is defined by the recursion below.

---

$\Pi' = \text{greedy-joining}(\Pi)$

---

choose $\{B, C\} \in \underset{\{B', C'\} \in \binom{\Pi}{2}}{\mathrm{argmin}} \ \varphi(y^{\mathrm{join}_{B'C'}[\Pi]}) - \varphi(y^{\Pi})$

if $\varphi(y^{\mathrm{join}_{BC}[\Pi]}) - \varphi(y^{\Pi}) < 0$

     $\Pi' := \text{greedy-joining}(\mathrm{join}_{BC}[\Pi])$

else

     $\Pi' := \Pi$

---

**Exercise 3**  *a) Write the difference $\varphi(y^{\mathrm{join}_{B'C'}[\Pi]}) - \varphi(y^{\Pi})$ in terms of the c defined in (6.18).*
   *b) Implement greedy joining efficiently.*
   *c) Establish a bound on the time complexity of your implementation.*

**Greedy moving**

The greedy moving algorithm starts from any initial partition, e.g., the fixed point of greedy joining, and seeks to lower the objective function by recursively moving individual elements from one cluster to another, or to new clusters. When an element is moved to a new cluster, the number of clusters can increase. When the last element is moved from a cluster, the number of clusters can decrease.

**Definition 19** For any partition $\Pi$ of $A$, any $a \in A$ and any $U \in \Pi \cup \{\varnothing\}$, let $\mathrm{move}_{aU}[\Pi]$ the partition of $A$ obtained by moving the element $a$ to a subset $U \cup \{a\}$ in $\Pi$.

$$\mathrm{move}_{aU}[\Pi] = \Pi \setminus \{U\} \setminus \{W \in \Pi \mid a \in W\} \cup \{U \cup \{a\}\} \cup \bigcup_{\{W \in \Pi \mid a \in W \wedge \{a\} \neq W\}} \{W \setminus \{a\}\} \ . \tag{6.21}$$

**Algorithm 2** The greedy moving algorithm is defined by the recursion below.

---

$\Pi' = $ greedy-moving$(\Pi)$

---

choose $(a, U) \in \underset{(a', U') \in A \times (\Pi \cup \{\varnothing\})}{\operatorname{argmin}} \varphi(y^{\operatorname{move}_{a'U'}[\Pi]}) - \varphi(y^{\Pi})$

if $\varphi(y^{\operatorname{move}_{aU}[\Pi]}) - \varphi(y^{\Pi}) < 0$

    $\Pi' := $ greedy-moving$(\operatorname{move}_{aU}[\Pi])$

else

    $\Pi' := \Pi$

---

**Exercise 4** *a) Write the difference $\varphi(y^{\operatorname{move}_{aU}[\Pi_t]}) - \varphi(y^{\Pi_t})$ in terms of the c defined in (6.18).*
  *b) Implement greedy moving.*

### Greedy moving using the technique of Kernighan and Lin (1970)

Both algorithms discussed above terminate as soon as no transformation (join and move, resp.) leads to a partition with strictly lower objective value. This can be suboptimal in case transformations that increase the objective value at one point in the recursion can lead to transformations that decrease the objective value at later points in the recursion and the decrease overcompensates the increase. One technique introduced by Kernighan and Lin (1970) for generalizing greedy algorithms with the goal of escaping such sub-optimal fixed points is applied below to greedy moving.

**Algorithm 3** The greedy moving algorithm generalized by the technique of Kernighan and Lin (1970) is defined by the recursion below.

---

$\Pi' = $ greedy-moving-kl$(\Pi)$

---

$\Pi_0 := \Pi$

$\delta_0 := 0$

$A_0 := A$

$t := 0$

repeat                                                               (build sequence of moves)

    choose $(a_t, U_t) \in \underset{(a, U) \in A_t \times (\Pi \cup \{\varnothing\})}{\operatorname{argmin}} \varphi(y^{\operatorname{move}_{aU}[\Pi_t]}) - \varphi(y^{\Pi_t})$

    $\Pi_{t+1} := \operatorname{move}_{a_t U_t}[\Pi_t]$

    $\delta_{t+1} := \varphi(y^{\Pi_{t+1}}) - \varphi(y^{\Pi_t}) < 0$

    $A_{t+1} := A_t \setminus \{a_t\}$                                       (move $a_t$ only once)

    $t := t + 1$

until $A_t = \varnothing$

$\hat{t} := \min \underset{t' \in \{0, \dots, |A|\}}{\operatorname{argmin}} \sum_{\tau = 0}^{t'} \delta_\tau$                           (choose sub-sequence)

if $\sum_{\tau = 0}^{\hat{t}} \delta_\tau < 0$

    $\Pi' := $ greedy-moving-kl$(\Pi_{\hat{t}})$                                  (recurse)

else

    $\Pi' := \Pi$                                                          (terminate)

---

**Exercise 5** *a) Implement greedy moving using the technique of Kernighan and Lin (1970).*
  *b) Generalize the greedy joining algorithm using the technique of Kernighan and Lin (1970).*
  *c) Implement greedy joining using the technique of Kernighan and Lin (1970).*

# Chapter 7

# Ordering

## 7.1 Orders

Throughout this chapter, we consider some finite set $A \neq \varnothing$ that we seek to order. Hence, our feasible solutions are *strict (total) orders* on $A$. A strict order on $A$ is a binary relation $< \subseteq A \times A$ that is trichotomous and transitive, i.e., it satisfies the following conditions:

$$\forall a \in A : \quad \neg a < a \tag{7.1}$$

$$\forall \{a, b\} \in \binom{A}{2} : \quad a < b \ \text{ xor } \ b < a \tag{7.2}$$

$$\forall \{a, b, c\} \in \binom{A}{3} : \quad a < b \ \wedge \ b < c \ \Rightarrow \ a < c \tag{7.3}$$

On the one hand, the strict orders on $A$ are characterized by the bijections $\alpha : \{0, \dots, |A|-1\} \to A$. For any such bijection, consider the order $<_\alpha$ such that

$$\forall a, b \in A : \quad a < b \ \Leftrightarrow \ \alpha^{-1}(a) < \alpha^{-1}(b) \ . \tag{7.4}$$

On the other hand, the strict orders on $A$ are characterized by those

$$y : \{(a, b) \in A \times A \mid a \neq b\} \to \{0, 1\} \tag{7.5}$$

that satisfy the following conditions:

$$\forall \{a, b\} \in \binom{A}{2} : \quad y_{ab} + y_{ba} = 1 \tag{7.6}$$

$$\forall \{a, b, c\} \in \binom{A}{3} : \quad y_{ab} + y_{bc} - 1 \leq y_{ac} \tag{7.7}$$

We reduce the problem of learning and inferring orders to the problem of learning and inferring decisions, by choosing constrained data with

$$S = \{(a, b) \in A \times A \mid a \neq b\} \tag{7.8}$$

$$\mathcal{Y} = \{y \in \{0, 1\}^S \mid (7.6) \wedge (7.7)\} \ . \tag{7.9}$$

One can think of the set $S$ as the set of edges in the complete digraph with nodes $A$ and without self-edges.

## 7.2 Linear ordering

### 7.2.1 Data

Throughout Section 7.2, we consider some finite set $A \neq \varnothing$ and constrained data $(S, X, x, \mathcal{Y})$ with $S = \{(a, b) \in A \times A \mid a \neq b\}$ as in (7.8), $X = \mathbb{R}^V$ and $\mathcal{Y}$ as in (7.9). As a special case, we consider labeled data, i.e., just one $\mathcal{Y} = \{y\}$ with $y$ satisfying the constraints (7.6) and (7.7).

## 7.2.2   Learning

The learning of linear orders is analogous to the learning of equivalence relations (Chapter 6).

## 7.2.3   Inference problem

**Corollary 4** *For any constrained data as defined above and any $\theta \in \mathbb{R}^V$, the inference problem takes the form of* LINEAR-ORDERING, *i.e.*

$$\min_{y \in \{0,1\}^S} \sum_{(a,b) \in S} \left(-\langle \theta, x_{ab} \rangle\right) y_{ab} \tag{7.10}$$

$$\text{subject to} \quad \forall \{a,b\} \in \binom{A}{2}: \quad y_{ab} + y_{ba} = 1 \tag{7.11}$$

$$\forall \{a,b,c\} \in \binom{A}{3}: \quad y_{ab} + y_{bc} - 1 \leq y_{ac} \ . \tag{7.12}$$

LINEAR-ORDERING has been studied intensively. A comprehensive survey is by Martí and Reinelt (2011). Important early research is by Grötschel et al. (1984).

**Lemma 10** LINEAR-ORDERING *is* NP-*hard.*

## 7.2.4   Inference algorithms

Below, we discuss two local search algorithms for the linear ordering problem, i.e., heuristics whose outputs need not be optimal.

For simplicity, we define $c : S \to \mathbb{R}$ such that

$$\forall (a,b) \in S: \quad c_{ab} = -\langle \theta, x_{ab} \rangle \ , \tag{7.13}$$

and write the objective function $\varphi : \{0,1\}^S \to \mathbb{R}$ such that

$$\forall y \in \{0,1\}^S: \quad \varphi(y) = \sum_{(a,b) \in S} c_{ab} \, y_{ab} \tag{7.14}$$

**Greedy transposition**

The greedy transposition algorithm starts from an initial strict order and searches for strict orders with lower objective value by swapping pairs of elements.

**Definition 20** For any bijection $\alpha : [|A|] \to A$ and any $j, k \in [|A|]$, let transpose$_{jk}[\alpha]$ the bijection obtained from $\alpha$ by swapping $\alpha_j$ and $\alpha_k$, i.e.

$$\forall l \in [|A|]: \quad \text{transpose}_{jk}[\alpha](l) = \begin{cases} \alpha_k & \text{if } l = j \\ \alpha_j & \text{if } l = k \\ \alpha_l & \text{otherwise} \end{cases} . \tag{7.15}$$

**Algorithm 4** The greedy transposition algorithm is defined by the recursion below.

---
$\alpha' = \text{greedy-transposition}(\alpha)$

---
choose $(j,k) \in \underset{0 \leq j' < k' \leq |A|}{\text{argmin}} \ \varphi(y^{\text{transpose}_{j'k'}[\alpha]}) - \varphi(y^\alpha)$

if $\varphi(y^{\text{transpose}_{jk}[\alpha]}) - \varphi(y^\alpha) < 0$

$\quad \alpha' := \text{greedy-transposition}(\text{transpose}_{jk}[\alpha])$

else

$\quad \alpha' := \alpha$

---

**Exercise 6** *a) Write the difference $\varphi(y^{\text{transpose}_{j'k'}[\alpha]}) - \varphi(y^\alpha)$ in terms of the c defined in (7.13).*
*b) Implement greedy transposition.*

**Greedy transposition using the technique of Kernighan and Lin (1970)**

As in the case of greedy heuristics for the correlation clustering problem, the greedy transposition algorithm for the linear ordering problem has been generalized using the idea of Kernighan and Lin (1970).

**Algorithm 5 (Martí and Reinelt (2011))** The greedy transposition algorithm generalized by the technique of Kernighan and Lin (1970) is defined by the recursion below.

---

$\alpha' = \text{greedy-transposition-kl}(\alpha)$

---

$\alpha^0 := \alpha$
$\delta_0 := 0$
$J_0 := [|A|]$
$t := 0$
repeat $\hspace{6cm}$ (build sequence of swaps)
$\quad$ choose $(j,k) \in \underset{\{(j',k')\in J^2 | j'<k'\}}{\text{argmin}} \varphi(y^{\text{transpose}_{j'k'}[\alpha^t]}) - \varphi(y^{\alpha^t})$
$\quad \alpha^{t+1} := \text{transpose}_{jk}[\alpha_t]$
$\quad \delta_{t+1} := \varphi(y^{\alpha^{t+1}}) - \varphi(y^{\alpha^t}) < 0$
$\quad J_{t+1} := J_t \setminus \{j,k\} \hspace{4cm}$ (move $\alpha_j$ and $\alpha_k$ only once)
$\quad t := t+1$
until $|J_t| < 2$

$\hat{t} := \min \underset{t'\in\{0,...,|A|\}}{\text{argmin}} \sum_{\tau=0}^{t'} \delta_\tau \hspace{4cm}$ (choose sub-sequence)

if $\sum_{\tau=0}^{\hat{t}} \delta_\tau < 0$
$\quad \alpha' := \text{greedy-transposition-kl}(\alpha^{\hat{t}}) \hspace{4cm}$ (recurse)
else
$\quad \alpha' := \alpha \hspace{6cm}$ (terminate)

---

**Exercise 7** *Implement greedy transposition using the technique of Kernighan and Lin (1970).*

# Appendix A

# Combinatorial problems

## A.1 Satisfiability

**Definition 21** For any CNF defined by $V$ and $\theta$ as in Definition 4, and for the Boolean function $f_\theta$ defined by this form as in Definition 4, deciding whether there exists an $x \in \{0,1\}^V$ such that $f_\theta(x) = 1$ is called the instance of the *satisfiability problem (*SAT*)* with respect to $V$ and $\theta$.

Any instance of SAT with respect to a 3-CNF defined by $V$ and $\theta$ as in Definition 4 is additionally called an instance of the *3-satisfiability problem (*3-SAT*)* with respect to $V$ and $\theta$.

**Theorem 3 (Cook (1971))** SAT *is* NP-*complete.*

**Theorem 4 (Karp (1972))** SAT $\leq_p$ 3-SAT.

**Lemma 11 (Karp (1972))** 3-SAT *is* NP-*complete.*

PROOF  3-SAT $\in$ NP, as solutions can be verified efficiently.
3-SAT is NP-hard by Theorems 3 and 4.

## A.2 Matching

**Definition 22** For any $m \in \mathbb{N}$, any finite, non-empty sets $S_0, \ldots, S_{m-1}$ and any $T \subseteq S_0 \times \cdots \times S_{m-1}$ the set $T$ is called a *matching* of $S_0, \ldots, S_{m-1}$ iff, for all distinct $(s_0, \ldots, s_{m-1}), (s'_0, \ldots, s'_{m-1}) \in T$ and all $j \in [m]$, we have $s_j \neq s'_j$.

The matching is called *perfect* iff, in addition, for every $j \in [m]$ and every $s_j \in S_j$, there exists precisely one $s' \in T$ such that $s'_j = s_j$.

**Lemma 12** *If a perfect matching as in Definition 22 exists, we have* $|S_0| = \cdots = |S_{m-1}|$.

**Definition 23** For any $m \in \mathbb{N}$, any finite, non-empty sets $S_0, \ldots, S_{m-1}$ such that $|S_0| = \cdots = |S_{m-1}|$ and any $T \subseteq S_0 \times \cdots \times S_{m-1}$, deciding whether there exists a perfect matching $T' \subseteq T$ of $S_0, \ldots, S_{m-1}$ is called the instance of the *$m$-dimensional perfect matching problem ($m$-*PM*)* with respect to $m, S_0, \ldots, S_{m-1}$ and $T$.

**Theorem 5** 3-SAT $\leq_p$ 3-PM

PROOF  Consider any instance of 3-SAT defined by a 3-CNF defined by $V$ and $\theta$ as in Definition 4.
Let $n = |\theta|$ the number of clauses of the 3-CNF.
Choose any order on $\theta$ to obtain a bijection $\theta' : [n] \to \theta$.
Define an instance of 3-PM by three sets $A, B, C$ of equal cardinality and one set $T \subseteq A \times B \times C$ of triples constructed as follows:

- For any clause index $c \in [n]$ and any variable $v \in V$, let

$$(a_c^v, \ b_c^v, \ \bar{x}_c^v) \in T \tag{A.1}$$

$$(a_{c+1 \bmod n}^v, \ b_c^v, \ x_c^v) \in T \tag{A.2}$$

- For any clause index $c \in [n]$ and the $(V_0, V_1) = \theta_c'$, distinguish between the variables in the clause, that is, the set $V_0 \cup V_1$, and the variables not in the clause, that is, the set $V \setminus (V_0 \cup V_1)$:

    - For each variable $v \in V \setminus (V_0 \cup V_1)$, let

$$(\alpha_c^v, \ \beta_c^v, \ \bar{x}_c^v) \in T \tag{A.3}$$

$$(\alpha_c^v, \ \beta_c^v, \ x_c^v) \in T \tag{A.4}$$

    - For each variable $v \in V_0$, let

$$(a_c, \ b_c, \ \bar{x}_c^v) \in T \tag{A.5}$$

    - For each variable $v \in V_1$, let

$$(a_c, \ b_c, \ x_c^v) \in T \tag{A.6}$$

    - If $|V_0 \cup V_1| \geq 2$, then, for each variable $v \in V_0 \cup V_1$, let

$$(\alpha_c, \ \beta_c, \ x_c^v) \in T \tag{A.7}$$

$$(\alpha_c, \ \beta_c, \ \bar{x}_c^v) \in T \tag{A.8}$$

    - If $|V_0 \cup V_1| = 3$, then, for each variable $v \in V_0 \cup V_1$, let

$$(\alpha_c', \ \beta_c', \ x_c^v) \in T \tag{A.9}$$

$$(\alpha_c', \ \beta_c', \ \bar{x}_c^v) \in T \tag{A.10}$$

In order to verify that $|A| = |B| = |C|$, we count the elements defined in the triples above:

| Triples | Elements | $|A|$ | $|B|$ | $|C|$ |
|---|---|---|---|---|
| (A.1), (A.2) | $a_c^v$ | $n|V|$ | | |
| (A.1), (A.2) | $b_c^v$ | | $n|V|$ | |
| (A.1), (A.2) | $\bar{x}_c^v$ | | | $n|V|$ |
| (A.1), (A.2) | $x_c^v$ | | | $n|V|$ |
| (A.3), (A.4) | $\alpha_c^v$ | $\sum\limits_{(V_0,V_1)\in\theta} (|V| - |V_0| - |V_1|)$ | | |
| (A.3), (A.4) | $\beta_c^v$ | | $\sum\limits_{(V_0,V_1)\in\theta} (|V| - |V_0| - |V_1|)$ | |
| (A.5)–(A.10) | $a_c, \alpha_c, \alpha_c'$ | $\sum\limits_{(V_0,V_1)\in\theta} (|V_0| + |V_1|)$ | | |
| (A.5)–(A.10) | $b_c, \beta_c, \beta_c'$ | | $\sum\limits_{(V_0,V_1)\in\theta} (|V_0| + |V_1|)$ | |
| | Total | $2n|V|$ | $2n|V|$ | $2n|V|$ |

Next, we show that the instance of 3-SAT defined by $V$ and $\theta$ has a solution iff the instance of 3-PM defined by $A, B, C$ and $T$ has a solution.

($\Rightarrow$) Let $\hat{x} \in \{0,1\}^V$ a solution to the instance of 3-SAT defined by $V$ and $\theta$.

Hence, for any clause $(V_0, V_1) \in \theta$, we can choose one variable $v_c \in V_0 \cup V_1$ such that

$$(v_c \in V_0 \ \wedge \ \hat{x}(v_c) = 0) \ \vee \ (v_c \in V_1 \ \wedge \ \hat{x}(v_c) = 1) \ . \tag{A.11}$$

Consider any $T' \subseteq T$ such that:

- For any clause index $c \in [n]$ and any variable $v \in V$:

$$(a_c^v, \ b_c^v, \ \bar{x}_c^v) \in T' \ \Leftrightarrow \ \hat{x}(v) = 1 \tag{A.12}$$

$$(a_{c+1 \bmod n}^v, \ b_c^v, \ x_c^v) \in T' \ \Leftrightarrow \ \hat{x}(v) = 0 \tag{A.13}$$

- For any clause index $c \in [n]$, the $(V_0, V_1) = \theta_c'$ and any variable $v \in V \setminus (V_0 \cup V_1)$:

$$(\alpha_c^v, \ \beta_c^v, \ \bar{x}_c^v) \in T' \ \Leftrightarrow \ \hat{x}(v) = 0 \tag{A.14}$$

$$(\alpha_c^v, \ \beta_c^v, \ x_c^v) \in T' \ \Leftrightarrow \ \hat{x}(v) = 1 \tag{A.15}$$

- For any clause index $c \in [n]$:

  - For the $v_c \in V_0 \cup V_1$ chosen in (A.11):

$$(a_c, \ b_c, \ \bar{x}_c^{v_c}) \in T' \ \Leftrightarrow \ \hat{x}(v_c) = 0 \tag{A.16}$$

$$(a_c, \ b_c, \ x_c^{v_c}) \in T' \ \Leftrightarrow \ \hat{x}(v_c) = 1 \tag{A.17}$$

  - If $|V_0 \cup V_1| \geq 2$, for precisely one variable $v_c' \in (V_0 \cup V_1) \setminus \{v_c\}$:

$$(\alpha_c, \ \beta_c, \ \bar{x}_c^{v_c'}) \in T' \ \Leftrightarrow \ \hat{x}(v_c') = 0 \tag{A.18}$$

$$(\alpha_c, \ \beta_c, \ x_c^{v_c'}) \in T' \ \Leftrightarrow \ \hat{x}(v_c') = 1 \tag{A.19}$$

  - If $|V_0 \cup V_1| = 3$, for the remaining variable $v_c'' \in (V_0 \cup V_1) \setminus \{v_c, v_c'\}$:

$$(\alpha_c', \ \beta_c', \ \bar{x}_c^{v_c''}) \in T' \ \Leftrightarrow \ \hat{x}(v_c'') = 0 \tag{A.20}$$

$$(\alpha_c', \ \beta_c', \ x_c^{v_c''}) \in T' \ \Leftrightarrow \ \hat{x}(v_c'') = 1 \tag{A.21}$$

We show that $T'$ is a solution to the instance of 3-PM defined by $A, B, C$ and $T$, by verifying that each element of $A, B$ and $C$ occurs in precisely one triple in $T'$:

- For any $c \in [n]$, the $(V_0, V_1) = \theta_c'$ and any $v \in V$:

| Element | Triple | |
|---|---|---|
| $a_c^v$ | $(a_c^v, \ b_c^v, \ \bar{x}_c^v)$ | if $\hat{x}(v) = 1$ |
| | $(a_c^v, \ b_{c-1 \bmod n}^v, \ x_{c-1 \bmod n}^v)$ | if $\hat{x}(v) = 0$ |
| $b_c^v$ | $(a_c^v, \ b_c^v, \ \bar{x}_c^v)$ | if $\hat{x}(v) = 1$ |
| | $(a_{c+1 \bmod n}^v, \ b_c^v, \ x_c^v)$ | if $\hat{x}(v) = 0$ |
| $\bar{x}_v^c$ | $(a_c^v, \ b_c^v, \ \bar{x}_c^v)$ | if $\hat{x}(v) = 1$ |
| | $(\alpha_c^v, \ \beta_c^v, \ \bar{x}_c^v)$ | if $\hat{x}(v) = 0$ and $v \in V \setminus (V_0 \cup V_1)$ |
| | $(a_c, \ b_c, \ \bar{x}_c^{v_c})$ | if $\hat{x}(v) = 0$ and $v = v_c$ |
| | $(a_c, \ b_c, \ \bar{x}_c^{v_c'})$ | if $\hat{x}(v) = 0$ and $|V_0 \cup V_1| \geq 2$ and $v = v_c'$ |
| | $(a_c, \ b_c, \ \bar{x}_c^{v_c''})$ | if $\hat{x}(v) = 0$ and $|V_0 \cup V_1| = 3$ and $v = v_c''$ |
| $x_v^c$ | $(a_{c+1 \bmod n}^v, \ b_c^v, \ x_c^v)$ | if $\hat{x}(v) = 0$ |
| | $(\alpha_c^v, \ \beta_c^v, \ x_c^v)$ | if $\hat{x}(v) = 1$ and $v \in V \setminus (V_0 \cup V_1)$ |
| | $(a_c, \ b_c, \ x_c^{v_c})$ | if $\hat{x}(v) = 1$ and $v = v_c$ |
| | $(a_c, \ b_c, \ x_c^{v_c'})$ | if $\hat{x}(v) = 1$ and $|V_0 \cup V_1| \geq 2$ and $v = v_c'$ |
| | $(a_c, \ b_c, \ x_c^{v_c''})$ | if $\hat{x}(v) = 1$ and $|V_0 \cup V_1| = 3$ and $v = v_c''$ |

- For any $c \in [n]$, the $(V_0, V_1) = \theta_c'$ and any $v \in V \setminus (V_0 \cup V_1)$:

| Element | Triple | |
| --- | --- | --- |
| $\alpha_c^v$ | $(\alpha_c^v,\ \beta_c^v,\ \bar{x}_c^v)$ | if $\hat{x}(v) = 0$ |
| | $(\alpha_c^v,\ \beta_c^v,\ x_c^v)$ | if $\hat{x}(v) = 1$ |
| $\beta_c^v$ | $(\alpha_c^v,\ \beta_c^v,\ \bar{x}_c^v)$ | if $\hat{x}(v) = 0$ |
| | $(\alpha_c^v,\ \beta_c^v,\ x_c^v)$ | if $\hat{x}(v) = 1$ |

- For any $c \in [n]$:

| Element | Triple | |
| --- | --- | --- |
| $a_c$ | $(a_c,\ b_c,\ \bar{x}_c^{v_c})$ | if $\hat{x}(v_c) = 0$ |
| | $(a_c,\ b_c,\ x_c^{v_c})$ | if $\hat{x}(v_c) = 1$ |
| $b_c$ | $(a_c,\ b_c,\ \bar{x}_c^{v_c})$ | if $\hat{x}(v_c) = 0$ |
| | $(a_c,\ b_c,\ x_c^{v_c})$ | if $\hat{x}(v_c) = 1$ |

- For any $c \in [n]$ with $(V_0, V_1) = \theta_c'$ such that $|V_0 \cup V_1| \geq 2$:

| Element | Triple | |
| --- | --- | --- |
| $\alpha_c$ | $(\alpha_c,\ \beta_c,\ \bar{x}_c^{v_c'})$ | if $\hat{x}(v_c') = 0$ |
| | $(\alpha_c,\ \beta_c,\ x_c^{v_c'})$ | if $\hat{x}(v_c') = 1$ |
| $\beta_c$ | $(\alpha_c,\ \beta_c,\ \bar{x}_c^{v_c'})$ | if $\hat{x}(v_c') = 0$ |
| | $(\alpha_c,\ \beta_c,\ x_c^{v_c'})$ | if $\hat{x}(v_c') = 1$ |

- For any $c \in [n]$ with $(V_0, V_1) = \theta_c'$ such that $|V_0 \cup V_1| = 3$:

| Element | Triple | |
| --- | --- | --- |
| $\alpha_c'$ | $(\alpha_c',\ \beta_c',\ \bar{x}_c^{v_c''})$ | if $\hat{x}(v_c'') = 0$ |
| | $(\alpha_c',\ \beta_c',\ x_c^{v_c''})$ | if $\hat{x}(v_c'') = 1$ |
| $\beta_c'$ | $(\alpha_c',\ \beta_c',\ \bar{x}_c^{v_c''})$ | if $\hat{x}(v_c'') = 0$ |
| | $(\alpha_c',\ \beta_c',\ x_c^{v_c''})$ | if $\hat{x}(v_c'') = 1$ |

($\Leftarrow$) Let $T' \subseteq T$ be any solution to the instance of 3-PM defined by $A, B, C$ and $T$. Moreover, let $\hat{x} \in \{0,1\}^V$ such that

$$\forall v \in V: \quad \hat{x}_v = \begin{cases} 1 & \text{if } (a_0^v,\ b_0^v,\ \bar{x}_0^v) \in T' \\ 0 & \text{otherwise} \end{cases}. \tag{A.22}$$

We show that $\hat{x}$ is a solution to the instance of 3-SAT defined by $V$ and $\theta$:

To begin with, we observe the clause index 0 in (A.22) being an arbitrary choice, because, for any clause index $c \in [n]$, we have by construction of $T$:

$$(a_c^v,\ b_c^v,\ \bar{x}_c^v) \in T'$$
$$\Rightarrow\ (a_{c+1\ \mathrm{mod}\ n}^v,\ b_c^v,\ x_c^v) \notin T' \tag{A.23}$$
$$\Rightarrow\ (a_{c+1\ \mathrm{mod}\ n}^v,\ b_{c+1\ \mathrm{mod}\ n}^v,\ \bar{x}_{c+1\ \mathrm{mod}\ n}^v) \in T' \tag{A.24}$$

By induction, we conclude for any clause index $c \in [n]$:

$$(a_c^v,\ b_c^v,\ \bar{x}_c^v) \in T'\ \Leftrightarrow\ (a_0^v,\ b_0^v,\ \bar{x}_0^v) \in T' \tag{A.25}$$

For any clause index $c \in [n]$ and the $(V_0, V_1) = \theta_c'$, there exists a triple in $T'$ that contains $a_c$ (because $T'$ is a solution to the instance of 3-PM).

Thus, and by construction of $T$, there exists a $v \in V$ such that one of the following statements holds:

$$v \in V_0 \ \wedge \ (a_c, \ b_c, \ \bar{x}_c^v) \in T' \tag{A.26}$$

$$v \in V_1 \ \wedge \ (a_c, \ b_c, \ x_c^v) \in T' \tag{A.27}$$

We consider the two cases separately:

- If $v \in V_0$, we have

$$(a_c, \ b_c, \ \bar{x}_c^v) \in T'$$
$$\Rightarrow \ (a_c^v, \ b_c^v, \ \bar{x}_c^v) \notin T' \tag{A.28}$$
$$\Rightarrow \ (a_0^v, \ b_0^v, \ \bar{x}_0^v) \notin T' \qquad \text{by (A.25)} \tag{A.29}$$
$$\Rightarrow \ \hat{x}(v) = 0 \ . \tag{A.30}$$

The clause indexed by $c$ and defined by $(V_0, V_1)$ is satisfied because $v \in V_0$ and $\hat{x}(v) = 0$.

- If $v \in V_0$, we have

$$(a_c, \ b_c, \ x_c^v) \in T'$$
$$\Rightarrow \ (a_{c+1 \bmod n}^v, \ b_c^v, \ x_c^v) \notin T' \tag{A.31}$$
$$\Rightarrow \ (a_c^v, \ b_c^v, \ \bar{x}_c^v) \in T' \tag{A.32}$$
$$\Rightarrow \ (a_0^v, \ b_0^v, \ \bar{x}_0^v) \in T' \qquad \text{by (A.25)} \tag{A.33}$$
$$\Rightarrow \ \hat{x}(v) = 1 \ . \tag{A.34}$$

The clause indexed by $c$ and defined by $(V_0, V_1)$ is satisfied because $v \in V_1$ and $\hat{x}(v) = 1$.

**Exercise 8** *Follow the proof of Theorem 5 in order to:*

*(a) construct the instance of* 3-PM *for the instance of* 3-SAT *given by the 3-CNF $(x_1 \vee (1 - x_2) \vee x_3) \cdot ((1 - x_1) \vee x_2 \vee x_4)$.*

*(b) construct, for any solution to this instance of* 3-SAT, *the solution to the instance of* 3-PM.

**Lemma 13** 3-PM *is* NP-*complete.*

PROOF 3-PM $\in$ NP, as solutions can be verified efficiently.

3-PM is NP-hard, by Theorem 5 and Lemma 11.

## A.3 Packing

**Lemma 14** *For any finite set $S$ and any $\varnothing \notin \Sigma \subseteq 2^S$, we have*

$$\left| \bigcup_{U \in \Sigma} U \right| \leq \sum_{U \in \Sigma} |U| \ . \tag{A.35}$$

PROOF For $|\Sigma| = 0$, the statement is obviously correct.

For $|\Sigma| > 0$, there exists a $V \in \Sigma$ and

$$\left| \bigcup_{U \in \Sigma} U \right| = \left| V \cup \bigcup_{U \in \Sigma \setminus \{V\}} U \right| = |V| + \left| \bigcup_{U \in \Sigma \setminus \{V\}} U \right| - \left| V \cap \bigcup_{U \in \Sigma \setminus \{V\}} U \right|$$

$$\leq |V| + \left| \bigcup_{U \in \Sigma \setminus \{V\}} U \right|$$

$$\leq |V| + \sum_{U \in \Sigma \setminus \{V\}} |U| \qquad \text{by induction}$$

$$= \sum_{U \in \Sigma} |U| \ .$$

**Definition 24** For any finite set $S$ and any $\varnothing \notin \Sigma \subseteq 2^S$, the set $\Sigma$ is called a *packing* of $S$ iff the elements of $\Sigma$ are pairwise disjoint, i.e.:

$$\forall \{U, U'\} \in \binom{S}{2}: \quad U \cap U' = \varnothing \ . \tag{A.36}$$

**Lemma 15** *For any finite set $S$ and any $\varnothing \notin \Sigma \subseteq 2^S$, the set $\Sigma$ is a packing of $S$ iff*

$$\left| \bigcup_{U \in \Sigma} U \right| = \sum_{U \in \Sigma} |U| \ . \tag{A.37}$$

PROOF ($\Rightarrow$) For $|\Sigma| = 0$, the statement (A.37) is obviously correct.

For $|\Sigma| > 0$, there exists a $V \in \Sigma$ and

$$\left| \bigcup_{U \in \Sigma} U \right| = \left| V \cup \bigcup_{U \in \Sigma \setminus \{V\}} U \right| = |V| + \left| \bigcup_{U \in \Sigma \setminus \{V\}} U \right| - \left| V \cap \bigcup_{U \in \Sigma \setminus \{V\}} U \right|$$

$$= |V| + \left| \bigcup_{U \in \Sigma \setminus \{V\}} U \right| - \left| \bigcup_{U \in \Sigma \setminus \{V\}} \underbrace{(V \cap U)}_{\varnothing} \right|$$

$$= |V| + \left| \bigcup_{U \in \Sigma \setminus \{V\}} U \right| \qquad\qquad \text{by (A.36)}$$

$$= |V| + \sum_{U \in \Sigma \setminus \{V\}} |U| \qquad\qquad \text{by induction}$$

$$= \sum_{U \in \Sigma} |U|$$

($\Leftarrow$) For $|\Sigma| = 0$, the statement (A.37) is obviously correct.

For $|\Sigma| > 0$ and any $V \in \Sigma$:

$$\left| \bigcup_{U \in \Sigma} U \right| = \left| V \cup \bigcup_{U \in \Sigma \setminus \{V\}} U \right|$$

$$= |V| + \left| \bigcup_{U \in \Sigma \setminus \{V\}} U \right| - \left| V \cap \bigcup_{U \in \Sigma \setminus \{V\}} U \right|$$

$$= |V| + \left| \bigcup_{U \in \Sigma \setminus \{V\}} U \right| - \left| \bigcup_{U \in \Sigma \setminus \{V\}} (V \cap U) \right|$$

Therefore:

$$\left| \bigcup_{U \in \Sigma \setminus \{V\}} (V \cap U) \right| = |V| + \left| \bigcup_{U \in \Sigma \setminus \{V\}} U \right| - \left| \bigcup_{U \in \Sigma} U \right|$$

$$\leq |V| + \sum_{U \in \Sigma \setminus \{V\}} |U| - \left| \bigcup_{U \in \Sigma} U \right| \qquad\qquad \text{by Lemma 14}$$

$$= \sum_{U \in \Sigma} |U| - \left| \bigcup_{U \in \Sigma} U \right|$$

$$= \sum_{U \in \Sigma} |U| - \sum_{U \in \Sigma} |U| \qquad\qquad \text{by (A.37)}$$

$$= 0$$

Thus:

$$\forall U, V \in \binom{\Sigma}{2}: \quad U \cap V = \varnothing$$

## A.4 Covering

**Definition 25** For any set $S$ and any $\varnothing \notin \Sigma \subseteq 2^S$, the set $\Sigma$ is called a *cover* of $S$ iff

$$\bigcup_{U \in \Sigma} U = S \ . \tag{A.38}$$

A cover $\Sigma$ of $S$ is called *exact* (and a *partition* of $S$) iff it is also a packing of $S$, i.e., if the elements of $\Sigma$ are pairwise disjoint, i.e.:

$$\forall U, U' \in \binom{\Sigma}{2}: \quad U \cap U' = \varnothing \tag{A.39}$$

**Definition 26** Let $S$ be any set, and let $\varnothing \notin \Sigma \subseteq 2^S$.

Deciding whether there exists a $\Sigma' \subseteq \Sigma$ such that $\Sigma'$ is an exact cover of $S$ is called the instance of the *exact cover problem (*EC*)* with respect to $S$ and $\Sigma$.

Additionally, if $|S|$ is an integer multiple of three and any $U \in \Sigma$ is such that $|U| = 3$, the instance of EC with respect to $S$ and $\Sigma$ is also called the instance of the *exact cover by 3-sets problem (*EC-3*)* with respect to $S$ and $\Sigma$.

**Lemma 16** 3-PM $\leq_p$ EC-3

PROOF Consider any instance of 3-PM defined by $A, B, C$ and $T$. Without loss of generality, assume that $A, B, C$ are pairwise disjoint.

Now, consider the instance of EC-3 defined by

$$S = A \cup B \cup C \tag{A.40}$$
$$\Sigma = \{\{a, b, c\} \subseteq S \mid (a, b, c) \in T\} \tag{A.41}$$

Firstly, $|S|$ is an integer multiple of three, as $|A| = |B| = |C|$.

Secondly, every $U \in \Sigma$ is such that $|U| = 3$, by construction.

Thirdly, $\Sigma' \subseteq \Sigma$ is an exact cover of $S$ iff the triples defined uniquely by the elements of $\Sigma'$ are a perfect 3-dimensional matching of $A \times B \times C$, by construction.

**Lemma 17** EC-3 *is* NP*-complete.*

PROOF EC-3 $\in$ NP, as solutions can be verified efficiently.

EC-3 is NP-hard, by Lemmata 13 and 16.

**Lemma 18** EC *is* NP*-complete.*

PROOF EC $\in$ NP, as solutions can be verified efficiently.

EC is NP-hard, as EC-3 is NP-hard and EC-3 $\subseteq$ EC.

**Definition 27** For any $k \in \mathbb{N}$, deciding whether there exists a $\Sigma' \subseteq \Sigma$ such that (i) $\Sigma'$ is a cover of $S$ and (ii) $|\Sigma'| \leq k$ is called the instance of the *set cover problem (*SET-COVER*)* with respect to $S, \Sigma$ and $k$.

**Lemma 19** EC-3 $\leq$ SET-COVER

PROOF For any instance $(S, \Sigma)$ of EC-3, let $m \in \mathbb{N}$ such that $|S| = 3m$.

Consider the instance of SET-COVER defined by $(S, \Sigma, m)$.

We show for any $\Sigma' \subseteq \Sigma$: $\Sigma'$ solves the instance of EC-3 iff $\Sigma'$ solves the instance of SET-COVER.

($\Rightarrow$) If $\Sigma'$ is a solution to the instance of EC-3, we have $|\Sigma'| \leq m$ because

$$3m = |S| \stackrel{\text{(A.38)}}{=} \left| \bigcup_{U \in \Sigma'} U \right| \stackrel{\text{(A.39) and Lemma 15}}{=} \sum_{U \in \Sigma'} |U| = \sum_{U \in \Sigma'} 3 = 3|\Sigma'| \ .$$

($\Leftarrow$) If $\Sigma'$ is a solution to the instance of SET-COVER, we have

$$\left| \bigcup_{U \in \Sigma'} U \right| = \sum_{U \in \Sigma'} |U|$$

because

$$3m = |S| \stackrel{\text{(A.38)}}{=} \left| \sum_{U \in \Sigma'} U \right| \stackrel{\text{Lemma 14}}{\leq} \sum_{U \in \Sigma'} |U| = \sum_{U \in \Sigma'} 3 = 3|\Sigma'| \stackrel{|\Sigma'| \leq m}{\leq} 3m \ .$$

Thus, $\Sigma'$ is a packing of $S$, by Lemma 15.

Hence, $\Sigma'$ is a solution to the instance of EC-3.

## A.5   Coloring

**Definition 28** Let $G = (V, E)$ be any graph.

For any $m \in \mathbb{N}$ and any $c : V \to E$, the map $c$ is called an *m-coloring* of $G$ iff $\forall \{v, w\} \in E : c(v) \neq c(w)$. $G$ is called *m-colorable* iff there exists an $m$-coloring of $G$.
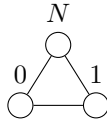
$G$ is called *m-chromatic*, and $m$ is called the *chromatic number* of $G$, iff $m$ is the minimal natural number such that $G$ is $m$-colorable.

**Definition 29** For any graph $G = (V, E)$ and any $m \in \mathbb{N}$, deciding whether $G$ is *m-colorable* is called the instance of the *m coloring problem (m-*COLORING*)* with respect to $G$ and $m$.
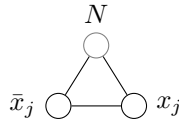
**Theorem 6 (Karp (1972))** 3-SAT $\leq_p$ 3-COLORING

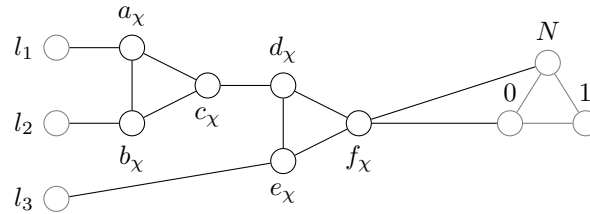PROOF Given any instance of 3-SAT, we define a graph $G = (V, E)$ as follows:

- Let $V$ contain a complete subgraph of three special nodes labeled 0, 1 and $N$. This subgraph is commonly called a *palette*.



- For every variable $x_j$ in the given instance of 3-SAT, let $V$ contain two nodes labeled $x_j$ and $\bar{x}_j$, respectively. Let both these nodes be connected to the palette node labeled $N$. Moreover, let the nodes labeled $x_j$ and $\bar{x}_j$ be connected by an edge:

- For any clause $(l_1 \lor l_2 \lor l_3) = \chi$ in the given instance of 3-SAT, let $V$ contain the following subgraph in which $l_j$ denotes the node labeled $x_k$ iff $l_j = x_k$, and $l_j$ denotes the node labeled $\bar{x}_k$ iff $l_j = 1 - x_k$:



Observe that the size of $G$ is polynomially bounded by the size of the given instance of 3-SAT.

**Exercise 9** *Complete the proof sketched above by showing that the instance of 3-SAT has a solution iff (!) $G$ is 3-colorable.*

**Lemma 20 (Karp (1972))** 3-COLORING *is* NP-*complete.*

PROOF 3-COLORING $\in$ NP, as solutions can be verified efficiently.
3-COLORING is NP-hard by Lemma 11 and Theorem 6.

**Lemma 21** $m$-COLORING *is* NP-*complete.*

PROOF $m$-COLORING $\in$ NP, as solutions can be verified efficiently.
$m$-COLORING is NP-hard, as 3-COLORING is NP-hard and 3-COLORING $\subseteq m$-COLORING.

# Bibliography

Bansal, N., A. Blum, and S. Chawla
  2004. Correlation clustering. *Machine Learning*, 56(1–3):89–113.

Cook, S. A.
  1971. The complexity of theorem-proving procedures. In *Proceedings of the Third Annual ACM Symposium on Theory of Computing*.

Dahlhaus, E., D. S. Johnson, C. H. Papadimitriou, P. D. Seymour, and M. Yannakakis
  1994. The complexity of multiterminal cuts. *SIAM Journal on Computing*, 23(4):864–894.

Demaine, E. D., D. Emanuel, A. Fiat, and N. Immorlica
  2006. Correlation clustering in general weighted graphs. *Theoretical Computer Science*, 361(2):172–187.

Garey, M. R. and D. S. Johnson
  1990. *Computers and Intractability; A Guide to the Theory of NP-Completeness*. USA: W. H. Freeman & Co.

Grötschel, M., M. Jünger, and G. Reinelt
  1984. A cutting plane algorithm for the linear ordering problem. *Operations Research*, 32(6):1195–1220.

Haussler, D.
  1988. Quantifying inductive bias: Ai learning algorithms and valiant's learning framework. *Artificial Intelligence*, 36(2):177–221.

Karp, R. M.
  1972. Reducibility among combinatorial problems. In *Complexity of Computer Computations*, R. E. Miller and J. W. Thatcher, eds., The IBM Research Symposia Series, Pp. 85–103. Plenum Press, New York.

Kernighan, B. W. and S. Lin
  1970. An efficient heuristic procedure for partitioning graphs. *Bell Systems Technical Journal*, 49(2):291–307.

Martí, R. and G. Reinelt
  2011. *The Linear Ordering Problem*. Springer.