## Part 1: Theoretical Exercises (16 points)

### 1. Gini Impurity

In class, we defined the Gini impurity as

$$\varphi_{Gini}(p) = 1 - \sum_{j=1}^{k} p_j^2, \qquad p \in [0,1]^k,$$

where $p = (p_1, \ldots, p_k)$ represents class proportions in a set of instances. This means that $\sum_{j=1}^{k} p_j = 1$.

1. Prove that

$$\varphi_{Gini}(p) \leq 1 - 1/k.$$

Hint:

- Express the function $f : \mathbb{R}^{k-1} \to \mathbb{R}$:

$$f(p_1, \ldots, p_{k-1}) = \varphi_{Gini}(p_1, \ldots, 1 - \sum_{j=1}^{k}).$$

- Argue that $f$ is bounded from above, hence it has a maximal value in $\mathbb{R}^{k-1}$.
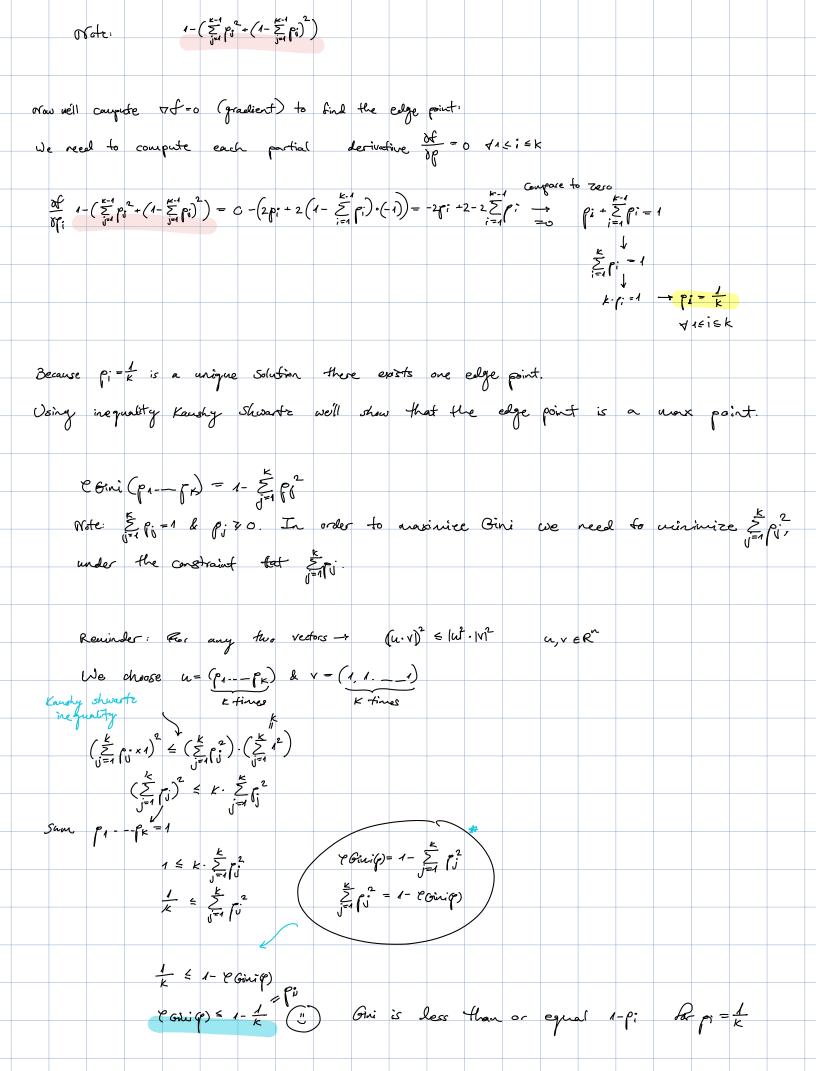- Solve the equation $\nabla f = 0$ and argue that the solution is unique.

(you do not have to follow the hint; all correct and clearly written solutions are acceptable)

Let $Y_1$ and $Y_2$ be two independent random variables, each represnting the class label of a randomly sampled instance from the set. Namely:

$$\Pr[Y_i = j] = p_j, \qquad i \in \{1, 2\}, \qquad j \in \{1 \ldots k\}.$$

2. Prove that Gini impurity is the probability that two randomly sampled instances (with replacement) from the set of instances have different class labels. Namley, that

$$\varphi_{Gini}(p) = \Pr[Y_1 \neq Y_2].$$

---

### Question #1

1. $\ell \, Gini \, (\varphi) = 1 - \sum_{j=1}^{K} p_j^2 \qquad p \in [0,1]^K \qquad \left( \sum_{j=1}^{K} p_j = 1 \quad \text{sum of proportions} \right)$

Prove: $\ell \, Gini \, (\varphi) \leq 1 - \frac{1}{K}$

$\ell \, Gini \, (p_1 \cdots p_k) = 1 - \sum_{j=1}^{K} p_j^2 \quad \& \quad \sum_{j=1}^{K} p_j = 1 \quad \rightarrow \text{therefore} \quad \sum_{j=1}^{K-1} p_j + p_K = 1$

$$p_K = 1 - \sum_{j=1}^{K-1} p_j$$

We'll define a function: $f : \mathbb{R}^{K-1} \to \mathbb{R}$ s.t. $f(p_1 \cdots p_k) = \ell \, Gini \, (p_1 \cdots 1 - \sum_{j=1}^{K-1} p_j)$ (Hint)

We'll also define $S = \sum_{j=1}^{K-1} p_j$. Note that $p_1 + \cdots + p_{K-1} = \sum_{j=1}^{K-1} p_j \rightarrow 1 - \sum_{j=1}^{K-1} p_j - \sum_{j=1}^{K-1} p_j = 1$ or sum of proportions is 1 therefore $p$ is well defined

$1 \geq$ sum of prop$^2 \geq 0$

· We will now prove $f$ is bounded from above:

$$f(p_1 \cdots p_{k-1}) = \ell \, Gini \, \left( p_1 \cdots 1 - \sum_{j=1}^{K-1} p_j \right) = \underbrace{\ell \, Gini \, (p_1 \cdots 1 - S)}_{K-1 \text{ elements}} = \underbrace{1 - \left( \sum_{j=1}^{K-1} p_j^2 + (1 - \sum_{j=1}^{K-1} p_j)^2 \right)}_{\text{By the definition of Gini}} \leq 1 \rightarrow f \text{ is bounded from above by 1.}$$

Note:  $1 - \left( \sum_{j=1}^{k-1} p_j^2 + \left( 1 - \sum_{j=1}^{k-1} p_i \right)^2 \right)$

Now we'll compute $\nabla f = 0$ (gradient) to find the edge point.

We need to compute each partial derivative $\frac{\partial f}{\partial p} = 0 \quad \forall 1 \le i \le k$

$\frac{\partial f}{\partial p_i} \quad 1 - \left( \sum_{j=1}^{k-1} p_j^2 + \left( 1 - \sum_{j=1}^{k-1} p_i \right)^2 \right) = 0 \quad - \left( 2 p_i + 2 \left( 1 - \sum_{i=1}^{k-1} p_i \right) \cdot (-1) \right) = -2 p_i + 2 - 2 \sum_{i=1}^{k-1} p_i \xrightarrow[=0]{\text{Compare to zero}} p_i + \sum_{i=1}^{k-1} p_i = 1$

$$\downarrow$$
$$\sum_{i=1}^{k} p_i = 1$$
$$\downarrow$$
$$k \cdot p_i = 1 \longrightarrow \boxed{p_i = \frac{1}{k}}$$
$$\forall 1 \le i \le k$$

Because $p_i = \frac{1}{k}$ is a unique solution there exists one edge point.

Using inequality Kaushy Shwartz we'll show that the edge point is a max point.

$\ell Gini(p_1 - \cdots p_k) = 1 - \sum_{j=1}^{k} p_j^2$

Note: $\sum_{j=1}^{k} p_j = 1$ & $p_j \ge 0$. In order to maximize Gini we need to minimize $\sum_{j=1}^{k} p_j^2$, under the constraint that $\sum_{j=1}^{k} p_j$.

Reminder: For any two vectors $\rightarrow \quad (u \cdot v)^2 \le |u|^2 \cdot |v|^2 \qquad u, v \in \mathbb{R}^n$

We choose $u = \underbrace{(p_1 \cdots p_k)}_{k \text{ times}}$ & $v = \underbrace{(1, 1, \ldots, 1)}_{k \text{ times}}$

Kaushy shwartz inequality

$\left( \sum_{j=1}^{k} p_j \times 1 \right)^2 \le \left( \sum_{j=1}^{k} p_j^2 \right) \cdot \left( \sum_{j=1}^{k} 1^2 \right)$  $\overset{k}{=}$

$\left( \sum_{j=1}^{k} p_j \right)^2 \le k \cdot \sum_{j=1}^{k} p_j^2$

Sum $p_1 \cdots p_k = 1$

$1 \le k \cdot \sum_{j=1}^{k} p_j^2$

$\frac{1}{k} \le \sum_{j=1}^{k} p_j^2$

$\ell Gini(p) = 1 - \sum_{j=1}^{k} p_j^2$

$\sum_{j=1}^{k} p_j^2 = 1 - \ell Gini(p)$

$\frac{1}{k} \le 1 - \ell Gini(p)$

$\ell Gini(p) \le 1 - \frac{1}{k} \quad \overset{= p_i}{(")} \quad$ Gini is less than or equal $1 - p_i \quad$ for $p_i = \frac{1}{k}$

2. Let $Y_1$ & $Y_2$ be two independent random variables. Each represent the class label of a randomly sampled instance from the set: $\Pr(Y_i = j) = p_j$ $i \in \{1, 2\}$ $j \in \{1 \ldots k\}$

Prove: gini impurity

$$\ell_{Gini}(p) = \Pr[Y_1 \neq Y_2]$$

↑ probability the two samples have different class labels.

$\Pr[Y_1 \neq Y_2] = 1 - \Pr[Y_1 = Y_2]$ (following hypoth)
(independent)

$\Pr[Y_1 = j \wedge Y_2 = j] \overset{\downarrow}{=} \Pr[Y_1 = j] \cdot \Pr[Y_2 = j] = p_j \cdot p_j = p_j^2$

↑ both samples have same label

$j \in \{1 \ldots k\}$ so we'll compute probability for each of the possibilities:

$$\Pr_{\substack{\cup \\ j=1}}^{k}(Y_1 = j \wedge Y_2 = j) = \sum_{j=1}^{k} p_j^2$$

definition ↓

$$\Pr[Y_1 \neq Y_2] = 1 - \Pr[Y_1 = Y_2] = 1 - \sum_{j=1}^{k} p_j^2 \overset{\downarrow}{=} \ell_{Gini}(p)$$

Question #2

1. Using $h(\lambda_1 x_1 + \lambda_2 x_2) \geq \lambda_1 h(x_1) + \lambda_2 \cdot h(x_2)$ $\quad \forall \; x_1, x_2 \in [0,1], \; \lambda_1, \lambda_2 \in [0,1]$ s.t. $\lambda_1 + \lambda_2 = 1$ *

we'll prove the claim by induction:

Base: For $t=2 \rightarrow$ $h\left(\sum_{j=1}^{2} \lambda_j x_j\right) = h\left(\lambda_1 x_1 + \lambda_2 x_2\right) \geq \lambda_1 h(x_1) + \lambda_2 \cdot h(x_2) = \sum_{j=1}^{2} \lambda_j h(x_j)$ *

Hypothesis: We assume that $h\left(\sum_{j=1}^{t-1} \lambda_j x_j\right) \geq \sum_{j=1}^{t-1} \lambda_j h(x_j)$

step: We'll prove $h\left(\sum_{j=1}^{t} \lambda_j x_j\right) \geq \sum_{j=1}^{t} \lambda_j x_j$ $\quad \forall \, t \geq 2$

Base case (2 points) and the induction hypothesis

$h\left(\sum_{j=1}^{t} \lambda_j x_j\right) = h\left(\lambda_t x_t + \sum_{j=1}^{t-1} \lambda_j x_j\right) \geq \lambda_t h(x_t) + \sum_{j=1}^{t-1} \lambda_j h(x_j) = \sum_{j=1}^{t} \lambda_j h(x_j)$ 🙂

2. Using the inequality we'll prove the claim: $IG(s,A) \geq 0$ for 2 labels

data   finite set of attributes          need to prove

$IG(s, A) = H(s) - \sum_{v \in values(A)} \frac{|S_v|}{|S|} H(s_v) \geq 0$

Given: $H(s) = h(p_1) = -\sum_{i=1}^{2} p_i \log(p_i) = -p_1 \log(p_1) - (1-p_1) \log(1-p_1)$

$H(s) = h\left(\sum_{j=1}^{t} \lambda_j s_j\right) \geq \sum_{j=1}^{t} \lambda_j h(s_j) = \sum_{j=1}^{t} \frac{|S_j|}{|S|} H(s_j)$

$1 \leq j \leq t$ values

$\lambda_j = \frac{|S_j|}{|S|}$

based on the jensen inequality

$H(s) \geq \sum_{j=1}^{t} \frac{|S_j|}{|S|} H(s_j)$

$IG(s,A) = H(s) - \sum_{j=1}^{t} \frac{|S_j|}{|S|} H(s_j) \geq 0$ 🙂