

1. **Decision Tree Example** You're stuck in a forest with nothing to eat. Suddenly, you spot a mushroom but you don't know if its poisonous. Luckily, you've studied some mushrooms as part of a class to fulfill your undergraduate requirements. Your previous knowledge is summarized by the following chart:

Sample #	IsColorful	IsSmelly	IsSmooth	IsSmall	IsPoisonous
1	0	0	0	1	1
2	0	0	0	0	0
3	1	0	1	1	1
4	1	0	0	0	1
5	0	0	1	0	0
6	0	0	1	0	0
7	1	1	0	0	1
8	1	1	1	0	1
9	0	1	1	0	?

- (a) What is the entropy of IsPoisonous, i.e.,  $H(IsPoisonous)$ ? **Solution:** Define the binary entropy function as follows:

$$H_b(p) = -p \log(p) - (1 - p) \log(1 - p).$$

$$H(IsPoisonous) = H_b\left(\frac{3}{8}\right) = -\left(\frac{5}{8} \log\left(\frac{5}{8}\right) + \frac{3}{8} \log\left(\frac{3}{8}\right)\right) \approx 0.9544$$

- (b) Calculate the conditional entropy of IsPoisonous conditioning on IsColorful. To do this, first compute  $H(IsPoisonous|IsColorful = 0)$  and  $H(IsPoisonous|IsColorful = 1)$ , then weight each term by the probabilities  $P(IsColorful = 0)$  and  $P(IsColorful = 1)$ , respectively. Namely, calculate the following:

$$\begin{aligned} & H(IsPoisonous|IsColorful) \\ &= P(IsColorful = 0)H(IsPoisonous|IsColorful = 0) \\ &+ P(IsColorful = 1)H(IsPoisonous|IsColorful = 1). \end{aligned}$$

**Solution:** Use the given equation, we get:

$$\begin{aligned} & H(IsPoisonous|IsColorful) \\ &= \frac{1}{2}H_b\left(\frac{3}{4}\right) + \frac{1}{2}H_b(1) \approx \frac{0.81127}{2} + 0 = 0.4056. \end{aligned}$$

(c) Similarly, calculate

$$H(IsPoisonous|X), \text{ for } X \in \{IsSmelly, IsSmooth, IsSmall\},$$

i.e., the conditional entropy of IsPoisonous conditioning on the other three features.

**Solution:**

$$H(IsPoisonous|IsSmelly) = \frac{2}{8}H_b(1) + \frac{6}{8}H_b\left(\frac{1}{2}\right) = 0.75.$$

$$H(IsPoisonous|IsSmooth) = \frac{1}{2}H_b\left(\frac{3}{4}\right) + \frac{1}{2}H_b\left(\frac{1}{2}\right) \approx \frac{0.8113}{2} + \frac{1}{2} = 0.9056.$$

$$H(IsPoisonous|IsSmall) = \frac{2}{8}H_b(1) + \frac{6}{8}H_b\left(\frac{1}{2}\right) = 0.75.$$

(d) Calculate the information gain:

$$I(IsPoisonous; X) = H(IsPoisonous) - H(IsPoisonous|X),$$

for

$$X \in \{IsColorful, IsSmelly, IsSmooth, IsSmall\}.$$

**Solution:**

$$I(IsPoisonous; IsColorful) = 0.9544 - 0.4056 = 0.5488;$$

$$I(IsPoisonous; IsSmelly) = 0.9544 - 0.75 = 0.2044;$$

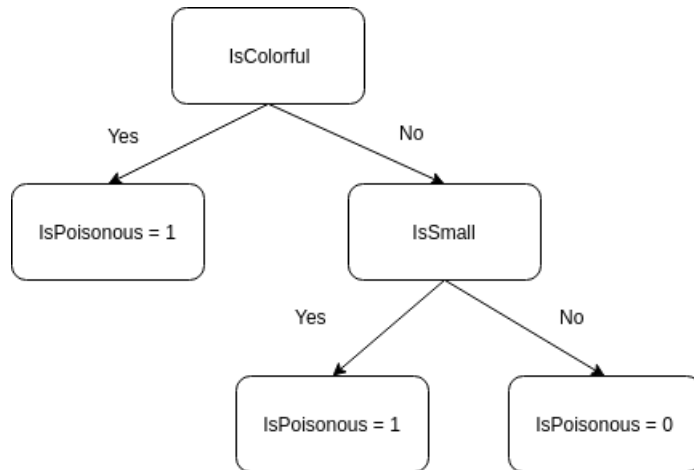
$$I(IsPoisonous; IsSmooth) = 0.9544 - 0.9056 = 0.1407;$$

$$I(IsPoisonous; IsSmall) = 0.9544 - 0.75 = 0.2044.$$

(e) Based on the information gain, determine the first attribute to split on. **Solution:** We choose IsColorful which has the largest information gain.

(f) Make the full decision tree. After each split, treat the sets of samples with  $X = 0$  and  $X = 1$  as two separate sets and redo (b), (c), (d) and (e) on each of them.  $X$  is the feature for previous split and is thus excluded from the available features which can be split on next. Terminate splitting if after the previous split, the entropy of IsPoisonous in the current set is 0. For example, if we choose IsSmall as our first feature to split, we get  $H(IsGoodRestaurant|IsSmall = 1) = 0$ . We thus stop splitting the tree in this branch. Draw the tree and indicate the split at each node.

**Solution:**



After the first split,  $H(IsPoisonous|IsColorful = 1) = 0$  so the tree stops growing on that branch. We are left with the samples that have  $IsColorful = 0$  which is summarized in the following table. We notice that  $H(IsPoisonous|IsSmall) =$

Sample #	IsSmelly	IsSmooth	IsSmall	IsPoisonous
1	0	0	1	1
2	0	0	0	0
5	0	1	0	0
6	0	1	0	0

0 for this reduced set. Therefore splitting using the feature IsSmall maximize the information gain.

(g) Is this mushroom poisonous? Not Poisonous!

2. **Regression Tree** So far, we have only focused on using tree structures for classification. We can also apply them to regression problems. In decision trees, we define the spread of a discrete dataset by using entropy. For real valued sets, we use variance.

For each set  $V$ , we associate a regression value  $u$  that minimizes the variance

$$Var(V) = \sum_{x_i \in V} (x_i - u)^2.$$

- (a) What is the value of  $u$  that minimizes  $Var(V)$ ?

**Solution:** Take the derivative with respect to  $u$ , set it to zero, and you get

$$u = \sum_{x_i \in V} (x_i) / |V|$$

- (b) Assume that a decision tree is trying to split  $V$  into two sets such that  $V_1 \cup V_2 = V$  and  $V_1 \cap V_2 = \emptyset$ . Write the formula for the reduction in variance.

**Solution:**

$$Var(V) - \left( \frac{|V_1|}{|V|} Var(V_1) + \frac{|V_2|}{|V|} Var(V_2) \right)$$

- (c) Example: You've always been told that drinking milk, getting plenty of sleep, eating your vegetables, and regularly exercising makes you grow up big and strong. Given your habits, you want to find out on average, how tall would you get? You ask your older friends whether they they did these things growing up and compile their answers in the following chart:

Sample #	DrinksMilk	SleepsWell	EatsVeggies	Height(cm)
1	0	1	1	200
2	0	1	0	210
3	0	1	0	200
4	1	1	0	180
5	1	0	1	130
6	1	0	0	150
9	1	1	1	?

Using this data, you will construct a regression tree to tell how tall you will get.

- i. What is variance of Height?

**Solution:** First, we get the average which is 178.33 which results in variance of around 5083.33

- ii. Determine the first attribute to split on by determining which attribute gives you the most reduction in variance.

**Solution:**

- DrinksMilk: 666.67
- SleepsWell: 383.33

- EatsVeggies: 2216.67

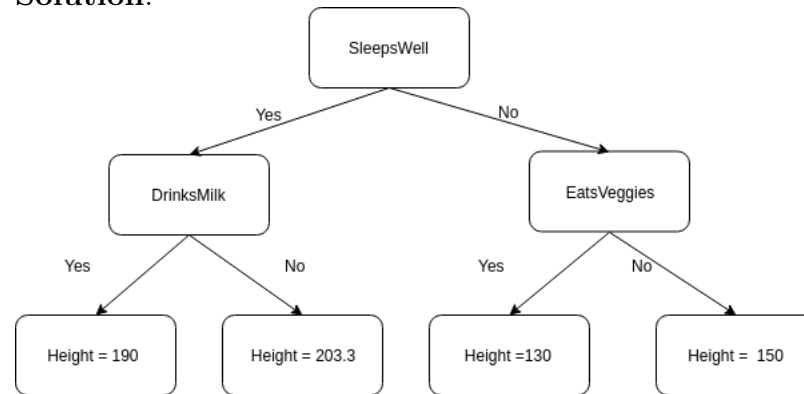
SleepsWell is the best attribute.

- iii. What is the reduction of variance for the previous attribute?

**Solution:** 4700

- iv. Make the full decision tree with max depth 2. Draw the tree and indicate the split at each node and the average at each leaf.

**Solution:**



- v. Now, determine how tall you would get.

**Solution:**

190

3. **Multi-class Classification Least Squares** In this section, you will determine the parameter matrix  $\mathbf{W} \in \mathbb{R}^{m \times p}$  for the Multi-class Least Squares classification.

Given a data matrix  $\mathbf{X} \in \mathbb{R}^{n \times m}$  and target matrix  $\mathbf{T} \in \mathbb{R}^{n \times p}$ , the sum-of-squares error function can be written as

$$Er(\mathbf{W}) = \text{Tr}\{(\mathbf{XW} - \mathbf{T})^T(\mathbf{XW} - \mathbf{T})\}$$

where  $\text{Tr}$  is the trace of a matrix. You can assume that  $\mathbf{X}$  has full rank.

We will solve this problem by setting the derivative with respect to  $\mathbf{W}$  to be zero and solve for  $\mathbf{W}$ . To do this we must first know some matrix derivative properties.

(a) Let  $\mathbf{A}, \mathbf{Z}$  be two matrices. Prove

$$\frac{d\text{Tr}(\mathbf{AZ})}{d\mathbf{Z}} = \mathbf{A}^T$$

**Solution:**

We can do this one variable at a time. Note that

$$\text{Tr}(\mathbf{AZ}) = \sum_i \sum_j z_{ij} a_{ji}.$$

Hence, we get

$$\frac{d\text{Tr}(\mathbf{AZ})}{dz_{ij}} = a_{ji}$$

which proves the result.

(b) Let  $\mathbf{A}, \mathbf{Z}$  be two matrices. Prove

$$\frac{d\text{Tr}(\mathbf{ZAZ}^T)}{d\mathbf{Z}} = \mathbf{ZA}^T + \mathbf{ZA}$$

**Solution:**

We can do this one variable at a time. Note that

$$\begin{aligned} \text{Tr}(\mathbf{ZAZ}^T) &= \sum_i \sum_j \sum_k z_{ij} a_{jk} z_{ik} \\ &= \sum_i \sum_j z_{ij}^2 a_{jk} + \sum_i \sum_j \sum_{k, k \neq j} z_{ij} a_{jk} z_{ik} \end{aligned}$$

Now, we can take the derivative with respect to each element in  $\mathbf{Z}$ .

First, let us differentiate with respect to  $z_{ii}$  and get

$$\begin{aligned} \frac{d\text{Tr}(\mathbf{ZAZ}^T)}{dz_{ii}} &= 2a_{ii}z_{ii} + \sum_{k, k \neq i} a_{ik}z_{ik} + \sum_{k, k \neq i} a_{ki}z_{ik} \\ &= \sum_k z_{ik}(a_{ik} + a_{ki}) \end{aligned}$$

which matches the result for the diagonal elements.

Now, let us differentiate with respect to  $z_{ij}$  where  $i \neq j$  and get

$$\frac{d\text{Tr}(\mathbf{ZAZ}^T)}{dz_{ij}} = \sum_k z_{ik}(a_{jk} + a_{kj})$$

which matches the result for the off-diagonal elements.

- (c) Now, we can take the derivative of  $Er(\mathbf{W})$  and set it to zero. Show that this results in

$$\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}$$

**Solution:**

$$\begin{aligned} Er(\mathbf{W}) &= \text{Tr}\{(\mathbf{XW} - \mathbf{T})^T(\mathbf{XW} - \mathbf{T})\} \\ &= \text{Tr}\{\mathbf{W}^T \mathbf{X}^T \mathbf{XW}\} - \text{Tr}\{\mathbf{W}^T \mathbf{X}^T \mathbf{T}\} - \text{Tr}\{\mathbf{T}^T \mathbf{XW}\} + \text{Tr}\{\mathbf{T}^T \mathbf{T}\} \\ \nabla_{\mathbf{W}} Er(\mathbf{W}) &= 2\mathbf{X}^T \mathbf{XW} - 2\mathbf{X}^T \mathbf{T} = 0 \\ \implies \mathbf{W} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T} \end{aligned}$$