CS2109S Tutorial 4 Decision Trees and Linear Models

(AY 24/25 Semester 2)

February 21, 2025

(Prepared by Benson)

Contents

Decision Trees

Recap Q1. Decisions That Matter

Linear Regression and Gradient Descent

- Q2. Linear Regression Model Fitting
- Q3. Examining Cost Functions
- Q4. Choosing Learning Rates

Bonus. Learning Rate Schedulers (Practical)

Recap: Decision Tree Learning

Information content / "surprisal": The informational value of communicating that an event happened.

$$I(e) = \log \left(\frac{1}{p}\right) = -\log p \text{ bits}$$
frequency of occurrence

Entropy: The expected amount of information conveyed by identifying the outcome of a random trial.

$$H(X) = \sum_{e \in E} P(e)I(e) = -\sum_{e \in E} P(e) \log P(e)$$

Information gain: Difference between the entropy before the split and the expected entropy (of a sample) after the split.

$$IG(D, A) = H(D) - \left[\sum_{\nu \in A} \frac{|D_{\nu}|}{|D|} H(D_{\nu})\right]$$

Remainder (minimize this)

(a) Construct the best decision tree to classify the final outcome (Decision) from the three features Income, Credit History, and Debt.



Income	Credit History	Debt	Decision
Over 10k	Bad	Low	Reject
Over 10k	Good	High	Approve
0 - 10k	Good	Low	Approve
Over 10k	Good	Low	Approve
Over 10k	Good	Low	Approve
Over 10k	Good	Low	Approve
0 - 10k	Good	Low	Approve
Over 10k	Bad	Low	Reject
Over 10k	Good	High	Approve
0 - 10k	Bad	High	Reject

(b) Construct the best decision tree. Calculate the information gain values and remainders at each stage.

remainder(Income) =
$$\frac{3}{10}I(2,1) + \frac{7}{10}I(6,1)$$

= $\frac{3}{10}(0.9183) + \frac{7}{10}(0.5917)$
= 0.690

remainder(Credit History) =
$$\frac{7}{10}I(7,0) + \frac{3}{10}I(1,2)$$

= $\frac{7}{10}(0) + \frac{3}{10}(0.9183)$
= 0.275

remainder(Debt) =
$$\frac{7}{10}I(6,1) + \frac{3}{10}I(1,2)$$

= $\frac{7}{10}(0.5917) + \frac{3}{10}(0.9183)$
= 0.690

)

	Inc	com	е	Cr	edit l	His	tory	0	Debt	D	ecision
	Ove	er 1	0k		Ba	nd			Low	A	pprove
ĺ	Over 10k 0 - 10k			Good			ł	High		pprove	
	0 - 10k Over 10k			Go	od			Low		pprove	
	Ove	er 1	0k		Go	od			Low	A	pprove
ĺ	Ονε	er 1	0k		Go	od			Low	Α	pprove
	Over 10k		Good				Low		pprove		
	0 -	10	k		Go	od			Low	A	pprove
ĺ	Ονε	er 1	0k		Ba	hd			Low		Reject
ĺ	Ove	Over 10k		Good			ł	High	Α	pprove	
ĺ	0 - 10k		k	Bad			ł	High		Reject	
	+ -	0	1	L	2		3		4		5
	1	0	1	L	0.91	83	0.811	3	0.721) (0.6500
	2	0	0.9	183	1		0.971	0	0.918	3 (0.8631
	3	0	0.8	113	0.97	10	1		0.985	2 (0.9544
	4	0	0.73	219	0.91	83	0.985	2	1	(0.9911
	5	0	0.6	500	0.863	31	0.954	4	0.991	L	1
	6	0	0.5	917	0.81	13	0.918	3	0.971) (0.9940
	7	0	0.54	436	0.76	42	0.881	3	0.945	7 (0.9799

(b) Construct the best decision tree. Calculate the information gain values and remainders at each stage.



remainder(Income) =
$$\frac{2}{3}I(1,1) + \frac{1}{3}I(0,1)$$

= $\frac{2}{3}(1) + \frac{1}{3}(0) = 0.667$
remainder(Debt) = $\frac{2}{3}I(1,1) + \frac{1}{3}I(0,1)$
= $\frac{2}{3}(1) + \frac{1}{3}(0) = 0.667$

	Inc	com	e (Credit His	story	Debt	Decision
Over 10k Over 10k			0k	Bad		Low	Approve
	Ove	er 1	0k	Good		High	Approve
	0 -	10	k	Good	Good		Approve
	Over 10k			Good		Low	Approve
	Ove	er 1	0k	Good		Low	Approve
	Ove	er 1	0k	Good		Low	Approve
	0 -	10	k	Good		Low	Approve
	Over 10k			Bad		Low	Reject
	Ove	er 1	0k	Good		High	Approve
	0 -	10	k	Bad		High	Reject
	+ -	0	1	2	3	4	5
	1	0	1	0.9183	0.811	3 0.721	9 0.6500
	2	0	0.9183	31	0.971	0 0.918	3 0.8631
	3	0	0.8113	3 0.9710	1	0.985	2 0.9544
	4	0	0.7219	9 0.9183	0.985	2 1	0.9911
	5	0	0.650	0.8631	0.954	4 0.991	1 1
	6	0	0.591	7 0.8113	0.918	3 0.971	0.9940
	7	0	0.543	0.7642	0.881	3 0.945	7 0.9799

(b) Construct the best decision tree. Calculate the information gain values and remainders at each stage.



Inc	com	e	Cr	edit	His	tory	[Debt	0	Decision	1
Ove	er 1	0k		Ba	hd			Low	ļ	Approve	
Over 10k				Go	od		I	High		Approve	
0 -	10	k		Go	od			Low	1	Approve	
Ove	er 1	0k		Go	od			Low	1	Approve	
Ove	er 1	0k		Go	od			Low	1	Approve	
Ove	er 1	0k		Go	od			Low	1	Approve	
0 -	10	k		Go	od			Low	1	Approve	
Ove	er 1	0k		Ba	hd			Low		Reject	
Ove	Over 10k		Good				High		Approve		
0 -	0 - 10k		Bad			I	High		Reject		
+	0	1	L	2		3		4		5	
1	0	1	L	0.91	83	0.811	.3	0.721	9	0.6500	
2	0	0.9	183	1		0.971	0	0.918	3	0.8631	
3	0	0.8	113	0.97	10	1		0.985	2	0.9544	
4	0	0.7	219	0.91	83	0.985	52	1		0.9911	
5	0	0.6	500	0.86	31	0.954	4	0.991	1	1	
6	0	0.5	917	0.81	13	0.918	33	0.971	C	0.9940	
7	0	0.5	436	0.76	42	0.881	.3	0.945	7	0.9799	

(c) What is the decision made by the decision tree in part (b) for a person with an income over 10k, a bad credit history, and low debt?



The decision might either reject or accept the person.

Inc	om	e C	redit His	tory	Debt	Decision
Ονε	er 1	0k	Bad		Low	Approve
Over 10k			Good		High	Approve
0 -	10	k	Good		Low	Approve
Over 10k			Good		Low	Approve
Ονε	er 1	0k	Good		Low	Approve
Ονε	er 1	0k	Good		Low	Approve
0 -	10	k	Good		Low	Approve
Ονε	er 1	0k	Bad		Low	Reject
Ονε	er 1	0k	Good		High	Approve
0 -	10	k	Bad		High	Reject
+ -	0	1	2	3	4	5
1	0	1	0.9183	0.811	3 0.721	9 0.6500
2	0	0.9183	1	0.971	0 0.918	3 0.8631
3	0	0.8113	0.9710	1	0.985	2 0.9544
4	0	0.7219	0.9183	0.985	2 1	0.9911
5	0	0.6500	0.8631	0.954	4 0.991	1 1
6	0	0.5917	0.8113	0.918	3 0.971	0.9940
7	0	0 5436	0 7642	0.881	3 0 945	7 0 0700

(d) The situation in part (c) demonstrates inconsistent data in the decision tree i.e. the attribute does not provide information to differentiate the two classes. In practice, it is usually left undecided. What are some ways to mitigate inconsistent data?

Solutions:

- Pruning (Min-sample / Max-depth).
- Pre-processing of data to remove outliers that create noise.
- Select only relevant features.
- Collect more data on new features to clearly differentiate the inconsistent classes.

	Inc	com	e	Cr	edit Hi	story	D	ebt	Decision	
	Income Over 10k O - 10k Ver 10k 0 - 10k 1 2 3 0 3 0		0k	Bad			L	ow	Approve	
	Ove	er 1	0k		Good		Н	igh	Approve	
Income Over 10k Over 10k			k	Good				ow	Approve	
			0k	Good				ow	Approve	
	Ove	er 1	0k		Good		L	ow	Approve	
	Ove	er 1	0k		Good		L	ow	Approve	
	0 -	10	k		Good		L	ow	Approve	
	Ove	er 1	0k		Bad		L	ow	Reject	
	Ove	Over 10k			Good			igh	Approve	
	0 -	10	k		Bad		Н	igh	Reject	
	+ -	0	1		2	3		4	5	
	1	0	1		0.9183	0.811	.3 (0.7219	0.6500	
	2	0	0.91	83	1	0.971	0 (0.9183	8 0.8631	
	3	0	0.81	13	0.9710	1	(0.9852	0.9544	
	4	0	0.72	19	0.9183	0.985	2	1	0.9911	
	5	0	0.65	00	0.8631	0.954	4 (0.9911	1	
	6	0	0.59	17	0.8113	0.918	3 (0.9710	0.9940	
	7		0 54	26	0 7642	0 001	2 1	0.457	7 0 0 7 0 0	

(e) Let's consider a scenario where you desire a Decision Tree with each leaf node representing a minimum of 3 training data points. Derive the tree by pruning the tree you previously obtained in part (b). Which data(s) do you think are likely outlier(s)?



The first person is probably the outlier.

Inc	com	e C	redit His	tory	Debt	Decision
Ove	er 1	0k	Bad		Low	Approve
Ove	er 1	0k	Good		High	Approve
0 -	10	k	Good		Low	Approve
Ove	er 1	0k	Good		Low	Approve
Ove	er 1	0k	Good		Low	Approve
Ove	er 1	0k	Good		Low	Approve
0 -	10	k	Good		Low	Approve
Ove	er 1	0k	Bad		Low	Reject
Ove	er 1	0k	Good		High	Approve
0 -	10	k	Bad		High	Reject
+	0	1	2	3	4	5
1	0	1	0.9183	0.811	3 0.721	9 0.6500
2	0	0.9183	1	0.971	0 0.918	3 0.8631
3	0	0.8113	0.9710	1	0.985	2 0.9544
4	0	0.7219	0.9183	0.985	2 1	0.9911
5	0	0.6500	0.8631	0.954	4 0.991	1 1
6	0	0.5917	0.8113	0.918	3 0.971	0.9940
7	0	0.5436	0.7642	0.881	3 0.945	7 0.9799

Q2. Linear Regression Model Fitting

(a) Apply the Normal Equation formula to obtain a linear regression model that minimizes MSE of the data points.

$$w = (\boldsymbol{X}^{ op} \boldsymbol{X})^{-1} \boldsymbol{X}^{ op} \boldsymbol{y}$$

Solution.

· .

$$\boldsymbol{X} = \begin{bmatrix} 1 & 6 & 4 & 11 \\ 1 & 8 & 5 & 15 \\ 1 & 12 & 9 & 25 \\ 1 & 2 & 1 & 3 \end{bmatrix}, \ \boldsymbol{y} = \begin{bmatrix} 20 \\ 30 \\ 50 \\ 7 \end{bmatrix}$$
$$\boldsymbol{w} = (\boldsymbol{X}^{\top} \boldsymbol{X})^{-1} \boldsymbol{X}^{\top} \boldsymbol{y} = \begin{bmatrix} 4 & -5.5 & -7 & 7 \end{bmatrix}^{\top}$$
$$\hat{\boldsymbol{y}} = 4 - 5.5 \boldsymbol{x}_1 - 7 \boldsymbol{x}_2 + 7 \boldsymbol{x}_3.$$

<i>x</i> ₁	<i>x</i> ₂	<i>x</i> 3	у
6	4	11	20
8	5	15	30
12	9	25	50
2	1	3	7

Q2. Linear Regression Model Fitting

- (b) Normal Equation needs the calculation of (X^TX)⁻¹. But sometimes this matrix is not invertible. When will that happen, and what should we do in that situation?
- When is $\mathbf{X}^{\top}\mathbf{X}$ invertible?
 - When all columns are linearly independent.
 - If some columns are linearly dependent, there are infinitely many solutions (thus it's impossible to find an "unique optimal solution").
- Almost linearly dependent columns create problems too...
 - A slight change in values drastically affects the result.
- **Solution**: Use gradient descent instead.

(It is also a good practice to remove linearly dependent features.)

Lemma. $\mathbf{X}^{\top}\mathbf{X}$ is invertible if and only if all columns of \mathbf{X} are linearly independent. *Proof.*

$$\boldsymbol{X}^{\top}\boldsymbol{X}$$
 is invertible
 $\Leftrightarrow \operatorname{rank}(\boldsymbol{X}^{\top}\boldsymbol{X}) = m$
 $\Leftrightarrow \operatorname{rank}(\boldsymbol{X}) = m$

since
$$\mathbf{X}^{\top}\mathbf{X}$$
 is an $m \times m$ matrix
rank $(\mathbf{X}^{\top}\mathbf{X}) = \operatorname{rank}(\mathbf{X})$

MA1522/2001 Refresher: rank($X^{\top}X$) = rank(X). (Prove this by showing the nullspace of $X^{\top}X$ is equal to the nullspace of X. See tutorial 7 for either course.)

Mean Squared Error: $L(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$ Mean Absolute Error: $L(y, \hat{y}) = \frac{1}{2}|y - \hat{y}|$

(a) Justify your choice of cost function for this problem.



Mean Squared Error:
$$L(y, \hat{y}) = rac{1}{2}(y-\hat{y})^2$$

Mean Absolute Error: $L(y, \hat{y}) = rac{1}{2}|y-\hat{y}|$

MSE penalizes large losses (caused by outliers) heavier than MAE.



Which line corresponds to MSE? Which line corresponds to MAE?

- MSE peanlizes large losses heavily \Rightarrow it will be more sensitive to outliers.
- These outliers could be a result of human error, they should have a smaller impact and MAE is preferred. If we consider outliers as important, MSE is preferred.



- True or False:
 - 1. Consider a dataset where all *y* values are larger than 1. MSE penalizes the outliers more heavily than MAE.
 - 2. Consider a dataset where all *y* values are between 0 and 1. MSE penalizes the outliers more heavily than MAE.



(b) Can you provide examples of cost functions that are better suited to handle outliers more effectively?

Huber loss: MSE \rightarrow MAE.

$$\mathcal{L}(y, \hat{y}) = egin{cases} rac{1}{2}(y - \hat{y})^2 & ext{for } |y - \hat{y}| \leq \delta \ \delta \cdot |y - \hat{y}| - rac{1}{2}\delta^2 & ext{otherwise} \end{cases}$$

Log-cosh loss: $\approx \frac{1}{2}x^2$ for small $x \rightarrow \approx |x| - \log 2$ for large x.

 $L(y, \hat{y}) = \log(\cosh(y - \hat{y}))$



Recap: Gradient Descent

We wish to minimize the loss y by varying x.

▶ Parameters: Initial value x, Learning rate α_{k}^{\dagger} Decides this How far do we change x?



Theorem. If α is small enough, gradient descent would reach a local minima. This would also be the global minima if f(x) is convex. For convex functions, local

minima \Rightarrow global minima.

 ${f o}$ Assume lpha is negative. Which of the following is the equation to perform the updates?

A.
$$x \leftarrow x + \alpha \frac{dy}{dx}$$

B. $x \leftarrow x - \alpha \frac{dy}{dx}$
C. $y \leftarrow y + \alpha \frac{dy}{dx}$
D. $y \leftarrow y - \alpha \frac{dy}{dx}$

Solution. When $\frac{dy}{dx} > 0$, we need to decrease x. Only option A successfully decreases x.





(a) Trace the gradient descent algorithm at the function $y = x^2$, starting from the point a = (5, 25) and using $\alpha = 10, 1, 0.1, 0.01$. Record the value of x over 5 iterations. $\alpha = 0.1$: 1. $\left. \frac{dy}{dx} \right|_{x=5} = 10$ $x = 5 - 0.1 \cdot 10 = 4.$ 2. $\left. \frac{dy}{dx} \right|_{x=4} = 8$ $x = 4 - 0.1 \cdot (8) = 3.2.$ 3. $\frac{dy}{dx}\Big|_{x=3,2} = 6.4$ $x = 3.2 - 0.1 \cdot 6.4 = 2.56$ 3.2 4. $\left. \frac{dy}{dx} \right|_{x=2.56} = 5.12$ 2.56 $x = 2.56 - 0.1 \cdot 5.12 = 2.048$. 2.0485. $\left. \frac{dy}{dx} \right|_{x=2.048} = 4.096$.6384 Х $x = 2.048 - 0.1 \cdot 4.096 = 1.6384$





- (b) During the course of training for a large number of epochs/iterations, what can be done to the value of the learning rate α to enable better convergence?
- Large α helps the model to converge faster, but might cause it to overshoot or even diverge.
- Idea: Vary α to help the model "stabilize". But how?
 - **Solution**: Decrease the learning rate α through the course of training.
 - > This is the logic behind a **learning rate scheduler**.

Bonus. Learning Rate Schedulers (Practical)

The given template code implements the gradient descent algorithm corresponding to Tutorial Q4. Copy the template code from the website.

Experiment with the different Pytorch learning rate schedulers and record your findings.



Bonus. Learning Rate Schedulers (Practical)

Sample Solution.

```
def optimize(lr_scheduler_class, scheduler_params={}, ...):
2
       . . .
3
      scheduler = lr_scheduler_class(optimizer, **scheduler_params)
4
5
      for i in range(num_iterations):
6
           . . .
8
           # Update the learning rate (i.e. take a step)
9
           if isinstance(scheduler, torch.optim.lr_scheduler.ReduceLROnPlateau):
               scheduler.step(metrics=loss.item())
           else:
               scheduler.step()
14
           x_history.append(x.item())
           lr_history.append(scheduler.get_last_lr()[0]) # Get current lr
16
           # optimizer.param_groups[0]['lr'] also works
17
18
19
       . . .
```

Bonus. Learning Rate Schedulers (Practical) Sample Solution.

