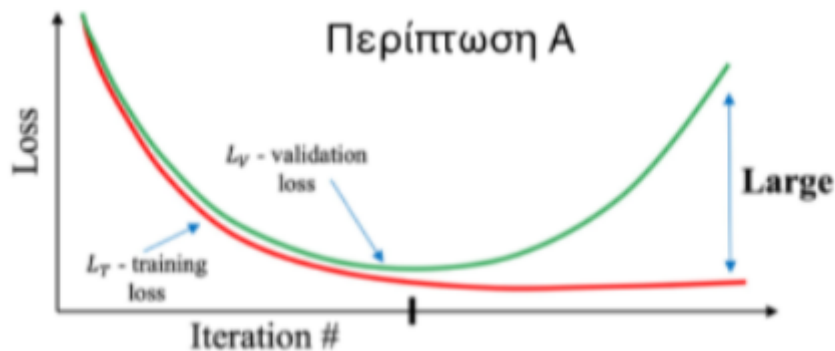


**ΝΕΥΡΩΝΙΚΑ ΔΙΚΤΥΑ ΚΑΙ ΒΑΘΙΑ ΜΑΘΗΣΗ**  
**Ακαδημαϊκό έτος 2023-2024 Σειρά Αναλυτικών Ασκήσεων**

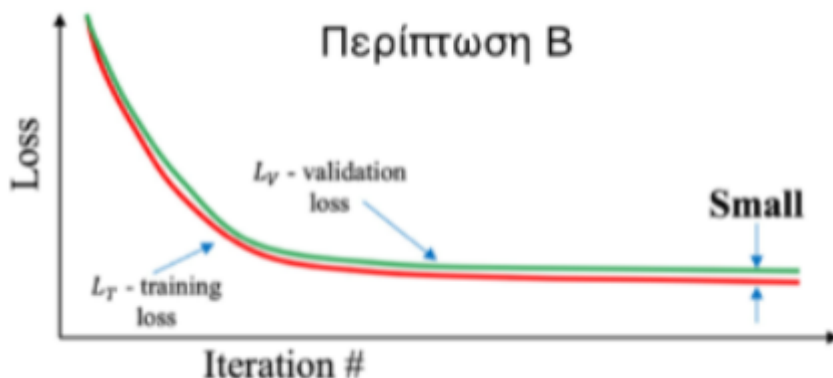
Θεοδώρα Εξάρχου  
ΑΜ:03120865

**Άσκηση 1 (Multi Layer Perceptron - Regularization)**

1) Τι συμπέρασμα βγάξετε από το διάγραμμα για την αρχιτεκτονική του μοντέλου σας για καθεμία από τις περιπτώσεις Α και Β;



**Περίπτωση Α:** Αναφέρεται στο φαινόμενο της υπερπροσαρμογής (overfitting) σε μοντέλα μηχανικής μάθησης. Αν κατανοήσουμε σωστά την ανάλυση σου, η σημαντική και συνεχής μείωση του εκπαιδευτικού σφάλματος (Training loss) σε συνδυασμό με την αύξηση του σφάλματος επικύρωσης (Validation loss) υποδηλώνει ότι το μοντέλο έχει μάθει να προσαρμόζεται υπερβολικά στα δεδομένα εκπαίδευσης, αλλά δυσκολεύεται να γενικεύσει σε νέα δεδομένα. Αυτό συμβαίνει συνήθως όταν το μοντέλο έχει πολύ υψηλή χωρητικότητα, που σημαίνει ότι έχει την ικανότητα να μάθει τα πραγματικά αλλά και τα τυχαία μοτίβα που υπάρχουν στα δεδομένα εκπαίδευσης, συμπεριλαμβανομένου και του θορύβου. Ως αποτέλεσμα, το μοντέλο μπορεί να είναι πολύ ακριβές στο σύνολο δεδομένων εκπαίδευσης, αλλά να παρουσιάζει χαμηλή ακρίβεια σε νέα, μη εκπαιδευτικά δεδομένα.



**Περίπτωση Β:** Η ανάλυση επικεντρώνεται στη συνεχή μείωση της συνάρτησης κόστους και στη σταθεροποίηση των τιμών της τόσο για το σύνολο εκπαίδευσης όσο και για το σύνολο επικύρωσης κατά τη διάρκεια των επαναλήψεων. Αυτό υποδεικνύει ότι το μοντέλο μαθαίνει αποτελεσματικά τα μοτίβα στα δεδομένα εκπαίδευσης χωρίς να προσθέτει περιττό θόρυβο. Η συνεχή μείωση της

συνάρτησης κόστους και η παρόμοια χαμηλή τιμή της για και τα δύο σύνολα υποδεικνύει ότι το μοντέλο δεν υπερπροσαρμόζεται στα δεδομένα εκπαίδευσης και έχει καταφέρει να επιτύχει μια ισορροπημένη γενίκευση. Το μικρό χάσμα μεταξύ της συνάρτησης κόστους για το σύνολο εκπαίδευσης και το σύνολο επικύρωσης υποδεικνύει ότι το μοντέλο είναι σε θέση να γενικεύει καλά σε νέα δεδομένα, δείχνοντας έτσι μια καλή ισορροπία μεταξύ εκπαίδευσης και γενίκευσης.

## **2) Ποια τιμή επαναλήψεων (εποχών) η είναι πιο πιθανό να οδηγεί στο καλύτερο δυνατό μοντέλο μάθησης**

**Περίπτωση Α:** Η καλύτερη τιμή είναι κοντά στο μισό των εποχών που φαίνεται στο διάγραμμα, καθώς μετά έχουμε το overfitting.

**Περίπτωση Β:** Η καλύτερη τιμή είναι κοντά στο μισό των εποχών, καθώς το loss μειώνεται με πολύ αργό ρυθμό. Επομένως, μπορούμε να εξοικονομήσουμε πόρους με μικρή διαφορά στο loss.

## **3) Τεχνικές για την περίπτωση Α που θα μπορούσαν να βελτιώσουν την επίδοση του μοντέλου**

- Εφαρμόζοντας Dropout rate στο μοντέλο
- Αλλάζοντας τις υπερπαραμέτρους

## **4) Εξηγήστε για ποιο λόγο είναι απαραίτητο το σύνολο ελέγχου (testing set) πέρα από τα σύνολα εκπαίδευσης και επικύρωσης.**

Το σύνολο ελέγχου (testing set) είναι απαραίτητο πέρα από τα σύνολα εκπαίδευσης και επικύρωσης για την **Αντικειμενική Αξιολόγηση σε Άγνωστα Δεδομένα**. Το σύνολο ελέγχου βοηθά στον αντικειμενικό έλεγχο του μοντέλου σε άγνωστα δεδομένα, εξασφαλίζοντας ότι η απόδοση που παρατηρούμε είναι αντιπροσωπευτική της πραγματικής ικανότητας του μοντέλου να χειριστεί νέα δεδομένα που δεν έχουν χρησιμοποιηθεί στη διαδικασία εκπαίδευσης και επικύρωσης.

## **Άσκηση 2 (Representation Learning - Autoencoders)**

**(α)** Ποια είναι η διάσταση των χαρακτηριστικών εισόδου  $x_i$  στον auto-encoder?

Τα διανύσματα  $u_i$  και  $v_i$  αντιπροσωπεύουν τις λέξεις στο μοντέλο skipgram. Αυτά τα διανύσματα αποτελούν τις ενσωματώσεις (embeddings) των λέξεων και έχουν διάσταση 256. Όταν εκπαιδεύουμε έναν αυτοκωδικοποιητή, χρησιμοποιούμε αυτά τα διανύσματα ως είσοδο. Άρα, η διάσταση των χαρακτηριστικών εισόδου  $x_i$  στον αυτοκωδικοποιητή θα είναι η ίδια με τη διάσταση των διανυσμάτων ενσωμάτωσης, δηλαδή 256.

**(β)** Ποια είναι η διάσταση των χαρακτηριστικών εξόδου  $y_i$  του auto-encoder?

Συνεπώς, η διάσταση των χαρακτηριστικών εξόδου  $y_i$  θα ταυτίζεται με τη διάσταση της εισόδου, δηλαδή 256.

**(γ)** Ποιά είναι η διάσταση της λανθάνουσας αναπαράστασης (latent representation) του auto-encoder?

Η διάσταση της λανθάνουσας αναπαράστασης (latent representation) ενός αυτοκωδικοποιητή είναι η διάσταση του μικρότερου κρυμμένου στρώματος (hidden layer) στον αυτοκωδικοποιητή. Εδώ, το μικρότερο κρυμμένο στρώμα έχει διάσταση 50. Επομένως, η διάσταση της λανθάνουσας αναπαράστασης του αυτοκωδικοποιητή είναι 50.

## **Άσκηση 3 (Recurrent Neural Networks) Χειρόγραφα Παρακάτω**

## **Άσκηση 4 (Convolutional Neural Networks) Χειρόγραφα Παρακάτω**

## **Άσκηση 5 (Generative models)**

### **α) Variational Autoencoders**

Χρησιμοποιούν κωδικοποιητή και αποκωδικοποιητή για την παραγωγή νέων δεδομένων. Εισάγουν κανονικοποίηση στον λανθάνοντα χώρο για να αποφεύγουν την υπερπροσαρμογή και να διασφαλίζουν τη συνεχή και πλήρη δομή του λανθάνοντα χώρου. Ο κωδικοποιητής συμπιέζει τα δεδομένα στον λανθάνοντα χώρο και ο αποκωδικοποιητής ανακατασκευάζει τα αρχικά δεδομένα από αυτόν τον χώρο.

### **β) Generative Adversarial Networks**

Αποτελούνται από δύο νευρωνικά δίκτυα, μια γεννήτρια και έναν διαχωριστή, που ανταγωνίζονται μεταξύ τους. Η γεννήτρια παράγει ψευδή δεδομένα και ο διαχωριστής προσπαθεί να τα διακρίνει από τα πραγματικά. Η γεννήτρια εκπαιδεύεται για να παράγει δεδομένα που μπορούν να ξεγελάσουν τον διαχωριστή, ενώ ο διαχωριστής εκπαιδεύεται για να αναγνωρίζει τα ψευδή δεδομένα.

### **γ) Diffusion models**

Βασίζονται σε διαδικασίες διάχυσης και αντιστρόφου διάχυσης θορύβου για τη δημιουργία νέων δειγμάτων. Προσθέτουν θόρυβο σταδιακά στα δεδομένα και μαθαίνουν να αντιστρέφουν αυτή τη διαδικασία. Θόρυβος προστίθεται σταδιακά στα δεδομένα σε μια προωθητική διαδικασία και αφαιρείται σε μια αντίστροφη διαδικασία διάχυσης.

## **Διαδικασία Εκπαίδευσης**

### **Variational Autoencoders**

Η εκπαίδευση περιλαμβάνει τη βελτιστοποίηση του κωδικοποιητή και του αποκωδικοποιητή για την ελαχιστοποίηση του σφάλματος ανακατασκευής. Εισάγεται κανονικοποίηση στον λανθάνοντα χώρο για να διατηρείται η δομή και να αποφεύγεται η υπερπροσαρμογή.

### **Generative Adversarial Networks**

Εκπαιδεύουν τη γεννήτρια και τον διαχωριστή σε εναλλασσόμενες περιόδους εκπαίδευσης. Ο διαχωριστής εκπαιδεύεται για να διακρίνει πραγματικά από ψευδή δεδομένα, ενώ η γεννήτρια εκπαιδεύεται για να ξεγελάσει τον διαχωριστή. Η σύγκλιση είναι δύσκολη και απαιτεί προσεκτική ρύθμιση των υπερπαραμέτρων.

### **Diffusion models**

Απαιτούν μεγάλο αριθμό βημάτων εκπαίδευσης (συνήθως μερικές χιλιάδες) για να παράγουν δείγματα υψηλής ποιότητας. Χρησιμοποιούν δυναμικές Langevin και προσδιοριστικές διαφορικές εξισώσεις για τη δημιουργία δειγμάτων.

## **Αποτελεσματικότητα και Απαιτήσεις**

Οι **Variational Autoencoders (VAEs)** απαιτούν λιγότερους πόρους για εκπαίδευση σε σχέση με τα GANs. Μπορεί να παράγουν δείγματα χαμηλότερης ποιότητας σε σύγκριση με τα GANs. Τα **Generative Adversarial Networks (GANs)** απαιτούν περισσότερη μνήμη και χρόνο εκπαίδευσης. Η σύγκλιση είναι δύσκολη και συχνά ασταθής. Μπορούν να παράγουν εξαιρετικά ρεαλιστικά δείγματα όταν εκπαιδεύονται σωστά. Τα **Diffusion Models** παράγουν δείγματα υψηλής ποιότητας. Ωστόσο, η διαδικασία δειγματοληψίας είναι αργή και απαιτεί πολλές αξιολογήσεις του νευρωνικού δικτύου, καθιστώντας τα δύσκολα για εφαρμογές πραγματικού χρόνου.

## Ομοιότητες και Διαφορές

Όλα χρησιμοποιούν νευρωνικά δίκτυα για τη δημιουργία δεδομένων και επιδιώκουν την εκμάθηση κατανομών δεδομένων για τη δημιουργία ρεαλιστικών δειγμάτων. Όσον αφορά τις διαφορές, Variational Autoencoders χρησιμοποιούν κωδικοποίηση και αποκωδικοποίηση με κανονικοποίηση του λανθάνοντα χώρου. Generative Adversarial Networks βασίζονται σε έναν ανταγωνισμό μεταξύ γεννήτριας και διαχωριστή. Diffusion Models προσθέτουν και αφαιρούν θόρυβο στα δεδομένα σε μια διαδικασία διάχυσης.

## Πλεονεκτήματα και Μειονεκτήματα

Τα **Variational Autoencoders (VAEs)** έχουν ως πλεονεκτήματα την εύκολη εκπαίδευση και την καλή κανονικοποίηση του λανθάνοντα χώρου. Ωστόσο, μπορεί να παράγουν δείγματα χαμηλότερης ποιότητας. Τα **Generative Adversarial Networks (GANs)** έχουν την ικανότητα δημιουργίας ρεαλιστικών δειγμάτων, αλλά η εκπαίδευσή τους είναι δύσκολη, απαιτεί μεγάλη υπολογιστική ισχύ και προσεκτική ρύθμιση των υπερπαραμέτρων. Τα **Diffusion Models** παράγουν δείγματα υψηλής ποιότητας και έχουν καλή θεωρητική βάση, αλλά η διαδικασία δειγματοληψίας είναι αργή και η εφαρμογή τους σε πραγματικά προβλήματα είναι δύσκολη λόγω του υψηλού υπολογιστικού κόστους.

## Άσκηση 6 (Graph Neural Networks) Χειρόγραφα Παρακάτω

### Exercice 3

$$a) C = \frac{1}{4} ([4, 0]^T + [0, 4]^T + [0, 0]^T + [0, 0]^T) = [1, 1]^T \\ \Rightarrow c_1 = [1, 1]^T$$

$$x_1 = \text{concat}(\langle \text{start} \rangle, a) = [0, 0, 0, 1, 1]^T$$

$$h_1 = g(W_{hh} \cdot h_0 + W_{hx} \cdot x_1) = \\ = g\left(\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \cdot [0, 0]^T + \begin{bmatrix} 1 & 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 \end{bmatrix} \cdot [0, 0, 0, 1, 1]^T\right) = \\ = g\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix}\right) \stackrel{\text{ReLU}}{=} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$y_1 = W_{gh} \cdot h_1 = \begin{bmatrix} -5 & -5 \\ 0 & 3 \\ 1 & 2 \\ 2 & 2 \\ 3 & -1 \\ 2,9 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -10 \\ 3 \\ 3 \\ 4 \\ 2 \\ 2,9 \end{bmatrix} \quad \text{H cat}$$

$$y_1 = y_{\text{cat}}$$

$$x_2 = \text{concat}(y_1, c_2) = [1, -2, 0, 1, 1]^T$$

$$h_2 = g(W_{hh} \cdot h_1 + W_{hx} \cdot x_2) = \\ = g\left(\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 1 & 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 \end{bmatrix} \cdot [1, -2, 0, 1, 1]^T\right) \\ = g\left(\begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 2 \\ -1 \end{bmatrix}\right) = g\left(\begin{bmatrix} 3 \\ 0 \end{bmatrix}\right) = \begin{bmatrix} 3 \\ 0 \end{bmatrix}$$



$$y_2 = W_{gh} \cdot h_2 = \begin{bmatrix} -5 & -5 \\ 0 & 3 \\ 1 & 2 \\ 2 & 2 \\ 3 & -1 \\ 2.9 & 0 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ 0 \end{bmatrix} = \begin{bmatrix} -15 \\ 0 \\ 3 \\ 6 \\ 9 \\ 8.7 \end{bmatrix} \quad \# \text{staring}$$

$$y_2 = y_{\text{staring}}$$

$$x_3 = \text{concat}(y_2, c_3) = [0, -1, -1, 1, 1]^T$$

$$h_3 = g(W_{hh} \cdot h_2 + W_{hx} \cdot x_3) =$$

$$= g \left( \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix} \cdot [0, -1, -1, 1, 1] \right)$$

$$= g \left( \begin{bmatrix} 0 \\ 3 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix} \right) = g \left( \begin{bmatrix} 0 \\ 2 \end{bmatrix} \right) = \begin{bmatrix} 0 \\ 2 \end{bmatrix}$$

$$y_3 = W_{gh} h_3 = \begin{bmatrix} -5 & -5 \\ 0 & 3 \\ 1 & 2 \\ 2 & 2 \\ 3 & -1 \\ 2.9 & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 2 \end{bmatrix} = \begin{bmatrix} -10 \\ 6 \\ 4 \\ 4 \\ -2 \\ 0 \end{bmatrix} \quad \# \text{<stop>}$$

Ensemble example <start> cat staring <stop>



6) Ο αποκωδικοποιητής RNN έχει προκατα-  
 scaled dot-product για κάθε βήμα το διαυγές  
 επεξεργαστής είναι  $h_{t-1}$  και οι αναρπαστές του κα-  
 ταμπίστη είναι  $i_1, i_2, i_3, i_4$  χρησιμοποιούνται ως  
 κλειδιά (key) και ως τιμές (value)

$$h_0 = 0$$

Για το πρώτο βήμα ( $t=1$ ) υπολογίζουμε τις ανα-  
 νωμίες πιθανότητας (attention probabilities) και το διαυγές  
 είναι που προκύπτει  $\alpha$ .

$$\text{key} = W^k x_i$$

$$\text{value} = W^v x_i$$

$$\text{query} = W^q \textcircled{x_i} \rightarrow \text{embedding matrix}$$

$$Z = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$

→ Διαστολή Διαυγισμός

Attention Scores

$$\text{score} = QK^T = h_0 [i_1, i_2, i_3, i_4]^T = \begin{bmatrix} 0 \\ 0 \end{bmatrix}^T \begin{bmatrix} 4 & 0 \\ 0 & 4 \\ 0 & 6 \\ 0 & 0 \end{bmatrix}$$

$$= [0 \ 0 \ 0 \ 0]$$

$$\text{Scaled score} = \frac{\text{score}}{\sqrt{d_k}} = \frac{[0 \ 0 \ 0 \ 0]}{\sqrt{2}}$$



$$P = \text{softmax}(\text{Scale} + \text{Score})$$

$$P = \frac{\exp(0)}{\sum_{k=1}^4 \exp(0)} = 0,25$$

Για το συνολικό διάνυσμα  $G$

$$G = \sum_{j=1}^4 p \cdot i_j = 0,25 [4,0]^T + 0,25 [0,4]^T + 0,25 [0,0]^T + 0,25 [0,0]^T$$

$$G = 0,25 [4,0]^T + 0,25 [0,4]^T = [1,0]^T + [0,1]^T = [1,1]^T$$

2ο  $G$  είναι όμοιο με του πρώτου ερωτήματος  
 που δεν χρειάζεται να συνεχίσουμε τον έλεγχο  
 Θα έχουμε την ίδια λέξη cat



## Exercice 4

AlexNet

διαστάσεις εισόδου  $227 \times 227 \times 3$  (εξαρτήσεις με RGB channels)

Φίλτρα  $11 \times 11 \times 3$  στο πρώτο convolutional layer.

Δίκτυο συνολικά 96 φίλτρα, stride ίσο με 4 και padding μηδενικό

a) Διαστάσεις στην έξοδο πρώτου convolutional layer.

$$w_2 = \frac{w_1 - F_w + 2P}{S_w} + 1 \quad \xrightarrow{P=0} \quad w_2 = \frac{w_1 - F}{S} + 1$$

$$H_2 = \frac{H_1 - F_h + 2p}{S_h} + 1 \quad \xrightarrow{P=0} \quad H_2 = \frac{H_1 - F}{S} + 1$$

$$P_2 = k P_1$$

$$\left. \begin{aligned} w_2 &= \frac{227 - 11}{4} + 1 \\ H_2 &= \frac{227 - 11}{4} + 1 \\ P_2 &= 1 \text{ για πρώτο φίλτρο} \end{aligned} \right\} \begin{aligned} w_2 &= 55 \\ H_2 &= 55 \\ p_2 &= 1 \end{aligned}$$



Διαστάσεις  $55 \times 55 \times 96$   
 και για  $k=96$  φίλτρα έχουμε  $P_2=961=96$

β)  $h_i$   $h_o$  units στο  $i$   $h_o$  convolutional layer

$$5.5 \cdot 55 = 3025 \text{ units σε ένα filter}$$

Για 96  $h_o$  έχουμε 290400 units

γ) Εκπαίδευση παραμέτρων του 1<sup>o</sup> conv layer με διαφορετικές βάσεις

$$[width \times height \times channels + 1] \times filters =$$

$$(11 \cdot 11 \cdot 3 + 1) \cdot 96 = 34.944$$

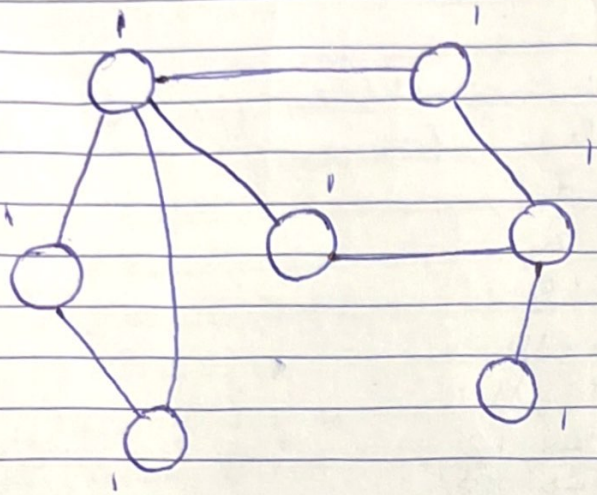
δ) Feedforward layer. 256 units

$$\begin{aligned} \text{Numb. of par.} &= (\text{Num. Input} + \text{Units}) \cdot \text{Units} \\ \text{Numbers of parameters} &= (227 \cdot 227 \cdot 3 \cdot 256) + 256 \\ &= 39574528 \end{aligned}$$

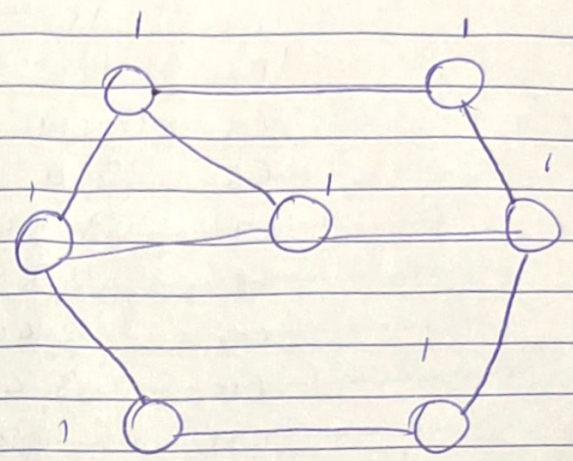


# Lesson 6

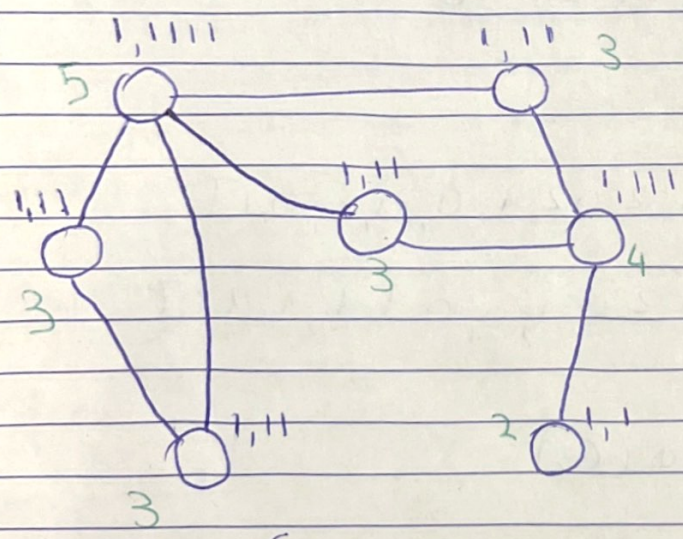
G<sub>1</sub>



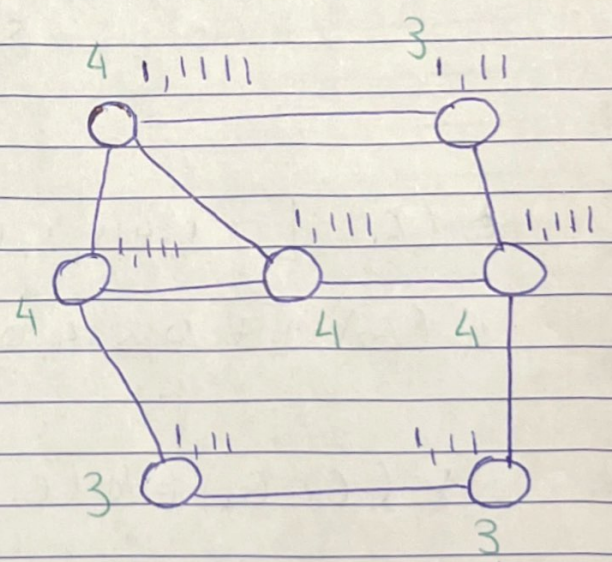
G<sub>2</sub>



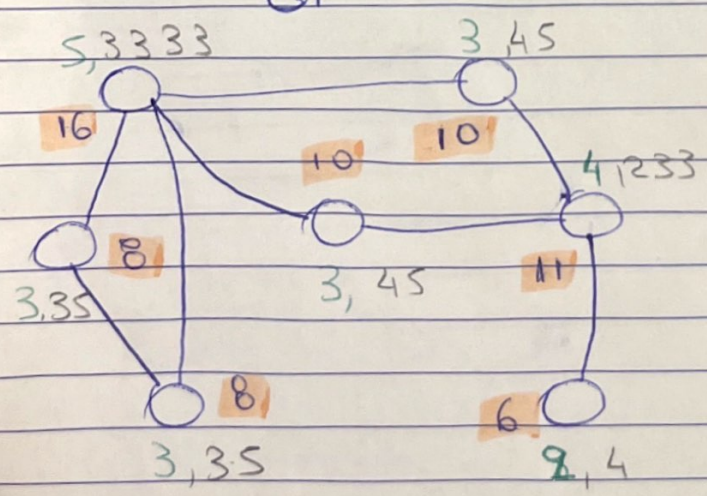
G<sub>1</sub>



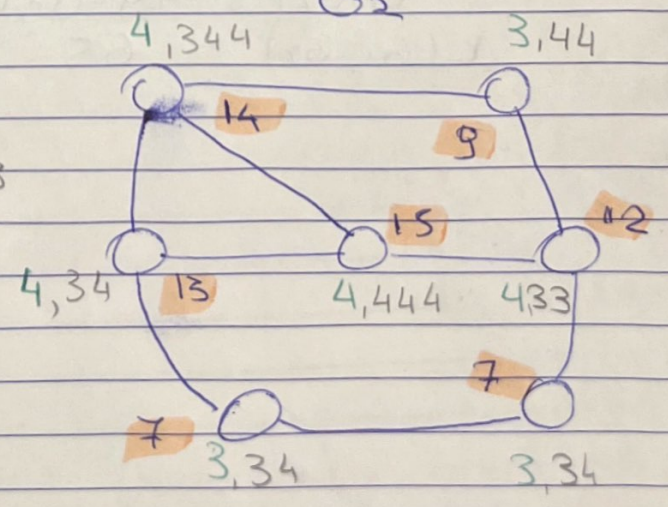
G<sub>2</sub>



G<sub>1</sub>



G<sub>2</sub>





hash

$c_2$	1, 1	$\rightarrow 2$
$c_3$	1, 11	$\rightarrow 3$
$c_4$	1, 111	$\rightarrow 4$
$c_5$	1, 1111	$\rightarrow 5$
$c_6$	2, 4	$\rightarrow 6$
$c_7$	3, 34	$\rightarrow 7$
$c_8$	3, 35	$\rightarrow 8$
$c_9$	3, 44	$\rightarrow 9$
$c_{10}$	3, 45	$\rightarrow 10$
$c_{11}$	4, 233	$\rightarrow 11$
$c_{12}$	4, 33	$\rightarrow 12$
$c_{13}$	4, 34	$\rightarrow 13$
$c_{14}$	4, 344	$\rightarrow 14$
$c_{15}$	4, 444	$\rightarrow 15$
$c_{16}$	5, 333	$\rightarrow 16$

$$\phi(G_1) = [7, 1, 4, 1, 1, 1, 0, 2, 0, 2, 1, 0, 0, 0, 0, 1]$$

$$\phi(G_2) = [7, 0, 3, 4, 0, 0, 2, 0, 1, 0, 0, 1, 1, 1, 1, 0]$$

$$\chi(G_1, G_2) = \phi(G_1) \cdot \phi(G_2) =$$

$$\chi(G_1, G_2) = 49 + 0 + 12 + 4 + 0$$

$$\chi(G_1, G_2) = 65$$