

ΤΜΗΜΑ ΠΛΗΡΟΦΟΡΙΚΗΣ  
ΠΜΣ ΕΠΙΣΤΗΜΗΣ ΥΠΟΛΟΓΙΣΤΩΝ

Διπλωματική Εργασία  
Μεταπτυχιακού Διπλώματος Ειδίκευσης

---

# Variational Autoencoders & Applications

---

Φοιτητής:  
Χρήστος ΚΟΡΜΑΡΗΣ  
ΕΥ1617

Επιβλέπων Καθηγητής:  
Δρ. Μιχάλης ΤΙΤΣΙΑΣ

ΑΘΗΝΑ, 3 ΜΑΙΟΥ 2018

# Δήλωση Συγγραφής

Εγώ, ο Χρήστος Κορμαρής, δηλώνω ότι αυτή η διατριβή με τον τίτλο "Variational Autoencoders" που παρουσιάζεται, είναι δικό μου έργο και επιβεβαιώνω ότι:

- Το έργο αυτό έγινε εξ' ολοκλήρου κατά τη διαδικασία υποψηφιότητας μου, για ερευνητικό πτυχίο στο Οικονομικό Πανεπιστήμιο Αθηνών.
- Οποιοδήποτε σημείο της παρούσας διατριβής έχει χρησιμοποιηθεί προηγουμένως για παρόμοιο πτυχίο στο Οικονομικό Πανεπιστήμιο Αθηνών ή σε οποιοδήποτε άλλο ίδρυμα, αυτό έχει αναφερθεί σαφώς.
- Στα σημεία που έχω συμβουλευτεί τη δημοσιευμένη εργασία άλλων ερευνητών, αυτό διευκρινίζεται πάντα σαφώς.
- Όπου έχω αναφερθεί σε δημοσιεύσεις άλλων ερευνητών, η πηγή αναφέρεται με ευκρίνεια. Με εξαίρεση τις αναφορές αυτές, η εργασία είναι εξ' ολοκλήρου δικό μου έργο.
- Έχω αναφέρει όλες τις κύριες πηγές βοήθειας.
- Όπου το έργο μου βασίζεται σε συνεργασία που έχω κάνει από κοινού με άλλους, αυτό αναφέρεται με ακρίβεια.

## Αρχαίο Κείμενο

“ ξεῖν’, ἥ τοι μὲν ὄνειροι ἀμήχανοι ἀκριτόμυθοι  
γίγνοντ’, οὐδέ τι πάντα τελείεται ἀνθρώποισι.  
δοιαί γάρ τε πύλαι ἀμενηνῶν εἰσὶν ὀνείρων·  
αἱ μὲν γὰρ κεράεσσι τετεύχεται, αἱ δ’ ἐλέφαντι·  
τῶν οἷ μὲν κ’ ἔλθωσι διὰ πριστοῦ ἐλέφαντος,  
οἷ ῥ’ ἐλεφαίρονται, ἔπε’ ἀκράαντα φέροντες·  
οἷ δὲ διὰ ξεστῶν κεράων ἔλθωσι θύραζε,  
οἷ ῥ’ ἔτυμα κραίνουσι, βροτῶν ὅτε κέν τις ἴδῃται. “

Ὅμηρου Ὀδύσσεια, Ραψωδία τ, στ. 562 (Πρωτότυπο)

## Απόδοση στην Καθαρεύουσα

“Ξένε, τὰ ὄνειρα βέβαια ἀμήχανα καὶ ἀκριτόμυθοι εἶναι, οὐδέ πάντα τελοῦνται στοὺς ἀνθρώπους.  
Γιατί δύο οἱ πύλες των ἀδυνάμων εἶναι ὀνείρων· ἡ μία ἀπὸ κέρατα ἔχει φτιαχθῆ, κι ἡ ἄλλη ἀπὸ  
ἐλεφαντόδοντο. Καὶ ὅσα διέλθουν ἀπὸ το πριονιστό ἐλεφαντόδοντο μᾶς ἀπατοῦν, ἔπη ἀνεκτέλεστα  
φέροντα· καὶ ὅσα διὰ ξυστῶν κεράτων ἔλθουν πρὸς τὰ ἔξω, αὐτὰ τ’ ἀληθινὰ κραίνουν, σέ ὅποιον ἀπὸ  
τοὺς βροτούς τὰ ἰδῇ.”

Ὅμηρου Ὀδύσσεια, Ραψωδία τ, στ. 562 (Μετάφραση Κώστα Δούκα)

## Προοίμιο (προ + οἶμος = δρόμος)

Σχολή Επιστημών & Τεχνολογίας της Πληροφορίας  
Τμήμα Πληροφορικής

Πρόγραμμα Μεταπτυχιακών Σπουδών στην Επιστήμη των Υπολογιστών

### Variational Autoencoders & Applications

Φοιτητής: Χρήστος Κορμαρής. Επιβλέπων Καθηγητής: Διδάκτωρ Μιχάλης Τίτσιας

Ο variational autoencoder είναι μια μέθοδος που μπορεί να παράγει τεχνητά δεδομένα που θα μοιάζουν με ένα δοσμένο σύνολο πραγματικών δεδομένων. Για παράδειγμα, αν θέλουμε να δημιουργήσουμε καινούργιες τεχνητές φωτογραφίες από γάτες, μπορούμε να χρησιμοποιήσουμε ένα variational autoencoder αλγόριθμο για να τα καταφέρουμε, αφού τον εκπαιδεύσουμε σε ένα μεγάλο σύνολο δεδομένων από γάτες. Τα δεδομένα εισόδου είναι μη κατηγοριοποιημένα διότι δεν ενδιαφερόμαστε να τα κατηγοριοποιήσουμε σε κάποια κλάση, αλλά επιθυμούμε να είμαστε σε θέση να μάθουμε τα πιο σημαντικά χαρακτηριστικά ή τις ομοιότητες μεταξύ τους. Εφόσον τα δεδομένα μας δεν έχουν ετικέτες-κατηγορίες, ο variational autoencoder χαρακτηρίζεται ως αλγόριθμος μη επιβλεπόμενης μάθησης. Όσον αφορά το παράδειγμα των εικόνων από γάτες, ο αλγόριθμος μπορεί να μάθει να αναγνωρίζει ότι μια γάτα πρέπει να έχει δύο αυτιά, μια μύτη, μουστάκια, τέσσερα πόδια, μια ουρά και μια ποικιλία από χρώματα. Ο αλγόριθμος χρησιμοποιεί δύο νευρωνικά δίκτυα, έναν κωδικοποιητή και έναν αποκωδικοποιητή, τα οποία εκπαιδεύονται ταυτόχρονα. Ένας variational autoencoder θα έχει καλές εφαρμογές σε περιπτώσεις όπου θα επιθυμούσαμε να παράγουμε μεγαλύτερο σύνολο δεδομένων, για καλύτερη εκπαίδευση σε διάφορα νευρωνικά δίκτυα. Επίσης, κάνει ελάττωση διάστασης των δεδομένων, συμπιέζοντάς τα σε λανθάνουσες μεταβλητές. Τρέχουμε υλοποιήσεις διαφορετικών variational autoencoder σε διάφορα dataset, όπως MNIST, Binarized MNIST, CIFAR-10, OMNIGLOT, YALE Faces, ORL Face Database, MovieLens, υλοποιημένες σε Python 3, χρησιμοποιώντας τρεις διαφορετικές βιβλιοθήκες κάθε φορά, TensorFlow, PyTorch και Keras και παρουσιάζουμε τα αποτελέσματα. Συστήνουμε έναν απλό αλγόριθμο συμπλήρωσης ελλিপών τιμών που χρησιμοποιεί K-NN συνεργατικό φιλτράρισμα για να κάνει προβλέψεις (πχ. για ελλιπή pixel). Τέλος, χρησιμοποιούμε το variational autoencoder σε αλγόριθμο συμπλήρωσης ελλিপών τιμών για να προβλέψουμε τις ελλιπείς τιμές σε διάφορα dataset. Ο K-NN αλγόριθμος τα πήγε σπουδαία στις προβλέψεις και το variational autoencoder σύστημα συμπλήρωσης επέφερε πολύ ικανοποιητικά αποτελέσματα. Έχει επίσης υλοποιηθεί GUI.

## *Ευχαριστίες*

Ευχαριστώ τον καθηγητή Δρ. *Μιχάλη Τίτσια* για τη βοήθεια και τις πολύτιμες συμβουλές του. Επίσης, θέλω να ευχαριστήσω τον υπεύθυνο του προγράμματος, τον κύριο καθηγητή Δρ. *Γεώργιο Πολύζο*, για την κατανόηση και την ενθάρρυνσή του κατά τη διάρκεια του μεταπτυχιακού. Τέλος, θέλω να ευχαριστήσω ξεχωριστά τους καθηγητές του μεταπτυχιακού, κύριο Δρ. *Γεώργιο Σταμούλη*, Δρ. *Ιωάννη Κωτίδη*, Δρ. *Βάνα Καλογεράκη*, Δρ. *Ευάγγελο Μαρκάκη*, Δρ. *Σταύρο Τουμπή*, Δρ. *Γιων Ανδρουτσόπουλο*, Δρ. *Γεώργιο Παπαϊωάννου*, Δρ. *Μιχάλη Βαζιργιάννη*, Δρ. *Βασίλειο Σύρη* και Δρ. *Σοφία Δημέλη* για όλη τη γνώση που μου προσέφεραν.

# Περιεχόμενα

	Σελίδα
Δήλωση Συγγραφής	i
Προοίμιο	iii
Ευχαριστίες	iv
Περιεχόμενα	v
Κατάλογος Εικόνων	viii
Κατάλογος Πινάκων	x
Κατάλογος Κώδικα	xi
Περιεχόμενα	xi
Συντομογραφίες	xii
Σύμβολα	xiii
1 Εισαγωγή	1
1.1 Προκαταρκτικά	1
1.2 Αναπαράσταση του Μοντέλου	1
1.3 Variational Συμπέρασμα (Variational Inference)	2
1.3.1 Πλεονεκτήματα του Variational Inference	3
1.3.2 Μειονεκτήματα του Variational Inference	3
1.3.3 Κύρια Βήματα του Variational Autoencoder	4
2 Υπολογίζοντας το Variational Ελάχιστο Φράγμα (ELBO)	5
2.1 Χρησιμοποιώντας την Ανισότητα του Jensen για να Υπολογίσουμε το Variational Ελάχιστο Φράγμα (ELBO)	5

2.2	Χρησιμοποιώντας την KL Απόκλιση για να Υπολογίσουμε το ELBO	6
<b>3</b>	<b>Αλγόριθμος Προς τα Πίσω Διάδοσης (Back-Propagation)</b>	<b>8</b>
3.1	Βελτιστοποιώντας τη Συνάρτηση Κόστους (ELBO)	8
3.2	Το Reparameterization Trick	10
3.3	Υπολογίζοντας τους κανόνες ενημέρωσης βαρών του κωδικοποιητή και του αποκωδικοποιητή	12
<b>4</b>	<b>Δομή του Variational Autoencoder &amp; Αποτελέσματα Πειραμάτων</b>	<b>14</b>
4.1	Δομή του Variational Autoencoder	14
4.2	Το TensorBoard	16
4.3	Σύνολα Δεδομένων (Datasets)	18
4.4	Αποτελέσματα Πειραμάτων	19
<b>5</b>	<b>Αλγόριθμοι Συμπλήρωσης Ελλιπών Τιμών &amp; Variational Autoencoders</b>	<b>27</b>
5.1	Ένας Απλός Αλγόριθμος Συμπλήρωσης Ελλιπών Τιμών Χρησιμοποιώντας K-NN Συνεργατικό Φιλτράρισμα (CF)	27
5.2	Ένας Προτεινόμενος Αλγόριθμος Συμπλήρωσης Ελλιπών Τιμών Χρησιμοποιώντας Variational Autoencoder	30
5.3	Αποτελέσματα Πειραμάτων του K-NN Αλγορίθμου Συμπλήρωσης Ελλιπών Τιμών και Αλγορίθμων Συμπλήρωσης Ελλιπών Τιμών με Χρήση VAE	32
5.4	K-NN vs VAE Αλγόριθμος Συμπλήρωσης Ελλιπών Τιμών στο MovieLens Dataset	36
<b>6</b>	<b>Γραφική Διεπαφή (GUI) για Variational Autoencoder &amp; Αλγορίθμους Συμπλήρωσης Ελλιπών Τιμών</b>	<b>37</b>
<b>7</b>	<b>Περαιτέρω Εφαρμογές των Variational Autoencoder</b>	<b>47</b>
<b>A</b>	<b>Θεωρητικό Υπόβαθρο</b>	<b>48</b>
A.1	Κανόνας του Bayes	48
A.2	Συνάρτηση Softmax	48

A.3	Εντροπία Θεωρίας Πληροφοριών (Entropy)	49
A.4	Cross-Entropy	49
A.5	Ανισότητα του Jensen	50
A.6	Kullback-Leibler (KL) Απόκλιση	51
A.7	Mean Absolute Error (MAE)	53
A.8	Mean Squared Error (MSE)	53
A.9	Το Ξυράφι του Occam	53
A.10	Root Mean Squared Error (RMSE)	53
A.11	Κανόνες Ενημέρωσης Βαρών των Νευρωνικών Δικτύων	54
<b>B</b>	<b>Πιθανοτικές Κατανομές</b>	<b>55</b>
B.1	Κατανομή Bernoulli	55
B.2	Διωνυμική Κατανομή	56
B.3	Κατανομή Gauss (ή Κανονική Κατανομή)	57
B.4	Νόμος Μεγάλος Αριθμών (N.M.A.)	59
B.5	Κεντρικό Οριακό Θεώρημα (Κ.Ο.Θ.)	60
<b>C</b>	<b>Προγραμματιστικές Υλοποιήσεις για Variational Autoencoder σε Python</b>	<b>61</b>
C.1	Υλοποίηση VAE με χρήση βιβλιοθήκης TensorFlow	61
C.2	Υλοποίηση VAE με χρήση βιβλιοθήκης PyTorch	67
C.3	Υλοποίηση VAE με χρήση βιβλιοθήκης Keras	72
C.4	K-NN Αλγόριθμος Συμπλήρωσης Ελλιπών Τιμών σε Python	74
C.5	VAE Αλγόριθμος Συμπλήρωσης Ελλιπών Τιμών σε PyTorch	76
	<b>Βιβλιογραφία</b>	<b>78</b>
	<b>Σύνδεσμοι</b>	<b>79</b>



# Κατάλογος Εικόνων

Εικόνα 1.1 Παράδειγμα δεδομένων εισόδου και εξόδου ενός VAE. Όπως παρατηρούμε, τα τελικά αποτελέσματα είναι μερικώς θολά. ....	4
Εικόνα 1.2 Σχήμα Variational Autoencoder.....	4
Εικόνα 2.1 Το Elbo Room είναι ένα bar που βρίσκεται στην ιστορική περιοχή Mission District του San Francisco στην οδό Valencia Street μεταξύ 17th Street και 18th Street. ....	7
Εικόνα 3.1 Το δεξιά δίκτυο χρησιμοποιεί το reparameterization trick. Το back-propagation μπορεί να εφαρμοστεί μόνο σε αυτό το δίκτυο. ....	11
Εικόνα 4.1 Ο autoencoder ονομάζεται και δίκτυο Diabolo λόγω της εμφάνισης της δομής του. (Rumelhart et al., 1986a; Bourlard & Kamp, 1988; Hinton & Zemel, 1994; Schwenk & Milgram, 1995; Japkowicz, Hanson, & Gluck, 2000, Yoshua Bengio 2009. [1]) .....	15
Εικόνα 4.2 Αυτή είναι απεικόνιση της υλοποίησης σε TensorFlow, με το TensorBoard. ....	16
Εικόνα 4.3 Αυτή είναι διευρυμένη απεικόνιση της υλοποίησης σε TensorFlow, με το TensorBoard. ....	17
Εικόνα 4.4 MNIST VAE σε TensorFlow. root mean squared error: 0.142598 .....	19
Εικόνα 4.5 MNIST VAE σε TensorFlow.....	21
Εικόνα 4.6 OMNIGLOT Αγγλικό αλφάβητο, χαρακτήρες 1-10, αρχικοί και ανακατασκευασμένοι. ....	22
Εικόνα 4.7 OMNIGLOT Ελληνικό αλφάβητο, χαρακτήρες 1-10, αρχικοί και ανακατασκευασμένοι. ....	22
Εικόνα 4.8 Binarized MNIST και (Real-valued) MNIST ψηφία, αρχικά και ανακατασκευασμένα.....	24
Εικόνα 4.9 CIFAR-10 εικόνες, RGB και grayscale, αρχικές και ανακατασκευασμένες. ....	24
Εικόνα 4.10 OMNIGLOT Αγγλικοί και Ελληνικοί χαρακτήρες, αρχικοί και ανακατασκευασμένοι. ....	25
Εικόνα 5.1 Αρχικά MNIST Δεδομένα Ελέγχου (Test).....	32

Εικόνα 5.2 <b>MNIST 1-NN, 100-NN &amp; VAE</b> αλγόριθμοι ελλিপών τιμών.....	33
Εικόνα 5.3 <b>1-NN &amp; VAE</b> αλγόριθμοι ελλিপών τιμών στο OMNIGLOT Αγγλικό αλφάβητο, χαρακτήρες 1-10, αρχικές, τυχαίες ελλιπείς τιμές και ανακατασκευασμένες.....	34
Εικόνα 5.4 <b>1-NN &amp; VAE</b> αλγόριθμοι ελλিপών τιμών στο OMNIGLOT Ελληνικό αλφάβητο, χαρακτήρες 1-10, αρχικές, τυχαίες ελλιπείς τιμές και ανακατασκευασμένες.....	35
Εικόνα 6.1 .....	38
Εικόνα 6.2 <b>Dropdown menus</b> .....	38
Εικόνα 6.3 <b>GUI VAE σε TensorFlow, MNIST dataset</b> .....	39
Εικόνα 6.4 <b>GUI VAE σε TensorFlow, CIFAR-10 dataset</b> .....	40
Εικόνα 6.5 <b>GUI VAE σε TensorFlow, OMNIGLOT dataset</b> .....	41
Εικόνα 6.6 <b>GUI K-NN αλγόριθμος Ελλিপών Τιμών, MNIST dataset</b> .....	42
Εικόνα 6.7 <b>GUI K-NN αλγόριθμος Ελλিপών Τιμών, CIFAR-10 dataset</b> .....	43
Εικόνα 6.8 <b>GUI K-NN αλγόριθμος Ελλিপών Τιμών, OMNIGLOT dataset</b> .....	44
Εικόνα 6.9 <b>GUI λεπτομέρειες των dataset</b> .....	45
Εικόνα 6.10 <b>GUI σχετικά</b> .....	46

# Κατάλογος Πινάκων

Πίνακας 4.1	Λεπτομέρειες των dataset.	18
Πίνακας 4.2	Πλήθος κατηγοριών, διαστάσεις και είδη τιμών των dataset.	18
Πίνακας 4.3	Σύνδεσμοι για τα dataset.	18
Πίνακας 4.4	VAE στο MNIST dataset σε TensorFlow.	20
Πίνακας 4.5	VAE στο Binarized MNIST dataset σε TensorFlow.	21
Πίνακας 4.6	VAE στο OMNIGLOT Αγγλικό dataset σε PyTorch.	23
Πίνακας 4.7	VAE στο OMNIGLOT Ελληνικό dataset σε PyTorch.	23
Πίνακας 4.8	VAE στο Binarized MNIST dataset, σε Keras.	25
Πίνακας 4.9	VAE στο MNIST dataset, σε Keras.	25
Πίνακας 4.10	VAE στο CIFAR-10 RGB dataset, σε Keras.	25
Πίνακας 4.11	VAE στο CIFAR-10 Grayscale dataset, σε Keras.	26
Πίνακας 4.12	VAE στο OMNIGLOT Αγγλικό dataset, σε Keras.	26
Πίνακας 4.13	VAE στο OMNIGLOT Ελληνικό dataset, σε Keras.	26
Πίνακας 5.1	Αλγόριθμοι συμπλήρωσης ελλιπών τιμών στο MNIST dataset.	33
Πίνακας 5.2	Αλγόριθμος συμπλήρωσης ελλιπών τιμών στο OMNIGLOT English dataset.	34
Πίνακας 5.3	Αλγόριθμος συμπλήρωσης ελλιπών τιμών στο OMNIGLOT Ελληνικό dataset.	35
Πίνακας 5.4	Σύγκριση μεταξύ των 10-NN και VAE σε TensorFlow αλγορίθμων συμπλήρωσης ελλιπών τιμών.	36
Πίνακας B.1	Κατανομή Bernoulli, συνάρτηση μάζας πιθανότητας (pmf), μέση τιμή και διασπορά.	55
Πίνακας B.2	Διωνυμική κατανομή, συνάρτηση μάζας πιθανότητας (pmf), μέση τιμή και διασπορά.	56
Πίνακας B.3	Κατανομή Gauss, συνάρτηση μάζας πιθανότητας (pmf), μέση τιμή και διασπορά.	59

# Κατάλογος Κώδικα

3.1	προς τα πίσω διάδοση σε TensorFlow . . . . .	12
3.2	προς τα πίσω διάδοση σε PyTorch . . . . .	12
4.1	εντολή για έναρξη του TensorBoard . . . . .	16
5.1	X_train_common . . . . .	28
5.2	X_test_common . . . . .	28
5.3	ψευδοκώδικας VAE αλγορίθμου συμπλήρωσης ελλιπών τιμών . . . . .	31
6.1	εντολή για εγκατάσταση Python εξαρτήσεων . . . . .	37
6.2	εντολή για εκτέλεση του GUI . . . . .	37
6.3	εντολή για δημιουργία εκτελέσιμου αρχείου για το GUI . . . . .	37
C.1	Κύριος Βρόχος του VAE σε TensorFlow . . . . .	61
C.2	Συνάρτηση VAE σε TensorFlow . . . . .	63
C.3	Κύριος Βρόχος του VAE σε PyTorch . . . . .	67
C.4	Αρχικοποίηση των Βαρών του VAE σε PyTorch . . . . .	68
C.5	Συνάρτηση Εκπαίδευσης του VAE σε PyTorch . . . . .	70
C.6	Κύριος Βρόχος του VAE σε Keras . . . . .	72
C.7	K-NN αλγόριθμος συμπλήρωσης ελλιπών τιμών σε Python . . . . .	74
C.8	VAE αλγόριθμος συμπλήρωσης ελλιπών τιμών σε PyTorch . . . . .	76

# Συντομογραφίες

<b>CF</b>	<b>C</b> ollaborative <b>F</b> iltering (Συνεργατικό Φιλτράρισμα)
<b>CIFAR</b> dataset	<b>C</b> anadian <b>I</b> nstitute for <b>A</b> dvanced <b>R</b> esearch dataset
<b>ELBO</b>	<b>E</b> vidence <b>L</b> ower <b>B</b> ound (Ελάχιστο Φράγμα Απόδειξης)
<b>EM</b> algorithm	<b>E</b> xpectation <b>M</b> aximization algorithm (αλγόριθμος Μεγιστοποίησης της Προσδοκίας)
<b>KL</b> divergence	<b>K</b> ullback- <b>L</b> eibler divergence (KL απόκλιση)
<b>K-NN</b>	<b>K</b> Nearest <b>N</b> eighbors (K Κοντινότεροι Γείτονες)
<b>MNIST</b> dataset	<b>M</b> odified <b>N</b> ational <b>I</b> nstitute of <b>S</b> tandards & <b>T</b> echnology dataset
<b>VAE</b>	<b>V</b> ariational <b>A</b> utoencoder
<b>VB</b> methods	<b>V</b> ariational <b>B</b> ayes methods (μέθοδοι Variational Bayes)
<b>TM</b>	<b>T</b> υχαία <b>M</b> εταβλητή

# Σύμβολα

$D_{KL}(P    Q)$	Kullback-Leibler Απόκλιση μεταξύ δύο κατανομών $\mathbf{P}$ και $\mathbf{Q}$
$E$	Μέση τιμή
$L$	Μεταβλητό Ελάχιστο Φράγμα (ELBO)
$N(\mu, \sigma^2)$	Κανονική (Normal) κατανομή με Μέση τιμή $\mu$ και Διασπορά $\sigma^2$
$P(X)$	Κατανομή πιθανότητας τυχαίας μεταβλητής $X$
$x$	Ένα δεδομένο εισόδου
$X$	Πολλά δεδομένα εισόδου
$z$	Μία λανθάνουσα μεταβλητή
$Z$	Πολλές λανθάνουσες μεταβλητές
$D$	Διάσταση δεδομένων εισόδου / Αριθμός των pixel
$M_1$	Αριθμός νευρώνων στον κωδικοποιητή
$M_2$	Αριθμός νευρώνων στον αποκωδικοποιητή
$Z_{dim}$	Διάσταση λανθανουσών μεταβλητών
$\epsilon \sim N(\mu, \sigma^2)$	Η T.M. (τυχαία μεταβλητή) $\epsilon$ ακολουθεί κανονική κατανομή
$\theta$	Παράμετροι του αποκωδικοποιητή (decoder)
$\mu$	Μία μεταβλητή μέσης τιμής
$M$	Πολλές μεταβλητές μέσης τιμής
$\sigma$	Μία μεταβλητή τυπικής απόκλισης
$\Sigma$	Πολλές μεταβλητές τυπικής απόκλισης
$\sigma^2$	Διασπορά
$\Sigma^2$	Διασπορές
$\phi$	Παράμετροι του κωδικοποιητή (encoder)
@	Τελεστής πολλαπλασιασμού πινάκων σε PyTorch
.*	Τελεστής πολλαπλασιασμού πινάκων στοιχείο προς στοιχείο (element-wise)

*Αφιερώνω τη διπλωματική μου εργασία στους γονείς μου.*

# Κεφάλαιο 1

## Εισαγωγή

### 1.1 Προκαταρκτικά

Έστω  $X$  τα δεδομένα εκπαίδευσης. Σκοπεύουμε να μεγιστοποιήσουμε την πιθανότητα κάθε δεδομένου εκπαίδευσης  $x$ , από τα δεδομένα εκπαίδευσης  $X$ , σύμφωνα με τον τύπο:

$$P(X) = \int P(X, z) dz = \int P(X | z, \phi) \cdot P(z) dz$$

όπου  $Z$  είναι συνεχής και ΟΧΙ διακριτή κατανομή και κάθε  $z$  είναι δεδομένο του  $Z$ . Συνεπώς, παίρνουμε ολοκλήρωμα από κοινού κατανομών και όχι άθροισμα, ώστε να υπολογίσουμε την κατανομή  $X$ . Η κατανομή  $X$  μπορεί να είναι είτε συνεχής είτε διακριτή. Για παράδειγμα, αν το dataset περιλαμβάνει εικόνες, η  $X$  μπορεί να είναι μια Bernoulli διακριτή κατανομή με διακριτές τιμές, 1 ή 0, 1 για άσπρα pixel και 0 για μαύρα pixel. Αντιθέτως, αν η  $X$  παίρνει πραγματικές τιμές μεταξύ του διαστήματος  $[0, 1]$ , τότε είναι συνεχής κατανομή. Με μια συνεχή κατανομή, το dataset μπορεί να αναπαραστήσει παραπάνω από δύο χρώματα. (Carl Doersch 2016. [2])

### 1.2 Αναπαράσταση του Μοντέλου

Τα δεδομένα εισόδου του μοντέλου μας είναι κυρίως εικόνες. Οι εικόνες εισόδου αναπαρίστανται ως μια μεταβλητή μεγέθους  $N \times D$ . Το  $N$  συμβολίζεται ως το μέγεθος των εικόνων και το  $D$  συμβολίζει τον αριθμό των pixel. Ένα καλό παράδειγμα είναι το *MNIST* dataset, που περιλαμβάνει εικόνες από μονοψήφιους αριθμούς, από το 0-9. Άλλα παραδείγματα από dataset με εικόνες είναι το *Binarized MNIST* dataset, τα *CIFAR-10* / *CIFAR-100* dataset, το *OMNIGLOT* dataset, το *YALE Faces* dataset και το *The Database of Faces* dataset. Επίσης, θα εξετάσουμε το *MovieLens* dataset, το οποίο περιλαμβάνει βαθμολογίες που έδωσαν πολλοί χρήστες σε συγκεκριμένες ταινίες.

Ο variational autoencoder είναι μια ακολουθία από δύο νευρωνικά δίκτυα, το ένα μετά το άλλο. Το πρώτο νευρωνικό δίκτυο ονομάζεται κωδικοποιητής (encoder) και το δεύτερο νευρωνικό δίκτυο ονομάζεται αποκωδικοποιητής (decoder). Επίσης, ο κωδικοποιητής και ο αποκωδικοποιητής εκπαιδεύονται ταυτόχρονα. Ο ρόλος του κωδικοποιητή είναι να μάθει πως να αναπαριστά τα κρυμμένα



χαρακτηριστικά του dataset, αποθηκεύοντάς τα σε λανθάνουσες μεταβλητές ελαττωμένης διάστασης. Από την άλλη πλευρά, ο αποκωδικοποιητής, κατασκευάζει τεχνητά δεδομένα, από τις λανθάνουσες μεταβλητές που έχουμε μάθει. Τα τεχνητά δεδομένα πρέπει να είναι παρόμοια με τα πρωτότυπα δεδομένα εισόδου, αλλά ΟΧΙ πανομοιότυπα. Διαφορετικά, η διαδικασία απέτυχε. Όταν ο αποκωδικοποιητής τελειώσει, ο στόχος μας έχει ολοκληρωθεί.

### 1.3 Variational Συμπέρασμα (Variational Inference)

Πρώτα θέλουμε να υπολογίσουμε τον κωδικοποιητή, δηλαδή θέλουμε να υπολογίσουμε:

$$P(Z | X) = \frac{P(X | Z) \cdot P(Z)}{P(X)} = \frac{P(X | Z) \cdot P(Z)}{\int P(X, Z) dz}$$

Ο κάθε όρος μπορεί να αναπαρασταθεί ως:

$$\text{posterior} = \frac{\text{likelihood} \cdot \text{prior}}{\text{normalizing constant}}$$

Όπου posterior:  $P(Z|X)$ , likelihood:  $P(X|Z)$ , prior:  $P(Z)$ , normalizing constant:  $P(X)$ .

Ο παρανομαστής  $P(X)$ , εκτός από σταθερά κανονικοποίησης (*normalizing constant*), ονομάζεται και απόδειξη (*evidence*). Μπορούμε να τον υπολογίσουμε περιθωριοποιώντας ως προς τις λανθάνουσες μεταβλητές  $Z$ , ως εξής:

$$P(X) = \int P(X | Z, \theta) \cdot P(Z) dz = \int P(X, Z, \theta) dz$$

Ωστόσο, υπολογίζοντας αυτό το ολοκλήρωμα απαιτεί εκθετικό χρόνο, επειδή η κατανομή των λανθανουσών μεταβλητών  $Z$  είναι συνεχής. Ο όρος  $P(X | Z, \theta)$  είναι μια πολύπλοκη συνάρτηση πιθανοφάνειας, λόγω της μη γραμμικότητας των κρυφών επιπέδων. Το πρόβλημα μεγιστοποίησης του όρου  $\log P(Z | X)$ , μέσω του κανόνα του Bayes, ανάγεται στο ακόλουθο:

$$\log P(Z | X) = \log \frac{P(X | Z) \cdot P(Z)}{P(X)} = \log P(X | Z) + \log P(Z) - \log P(X)$$

Εφόσον ο όρος  $P(X)$  είναι μη υπολογίσιμος, ο όρος  $P(Z | X)$  είναι επίσης μη υπολογίσιμος, μέσω του κανόνα του Bayes. Για να υπολογίσουμε το  $P(Z | X)$ , θα χρησιμοποιήσουμε μια άλλη μέθοδο που λέγεται **variational συμπέρασμα (variational inference)**, χρησιμοποιώντας μια οικογένεια κατανομών που ονομάζουμε  $Q_\phi(Z | X)$ , όπου  $\phi$  είναι μια παράμετρος. Μαθαίνουμε την παράμετρο  $\phi$  μέσω stochastic (ή mini-batch) αλγορίθμου ανοδικής (ή καθοδικής) κλίσης. Σε κάθε επανάληψη,

υπολογίζουμε τη συνάρτηση κόστους ή την πιθανοφάνεια, η οποία είναι το ελάχιστο φράγμα του όρου  $\log P(X)$ . Θέλουμε να μεγιστοποιήσουμε αυτόν τον όρο, επομένως θέλουμε να μεγιστοποιήσουμε το ελάχιστο φράγμα. Θα αναλύσουμε αυτή τη διαδικασία περαιτέρω στο [Κεφάλαιο 2](#).

Χρησιμοποιώντας το variational συμπίεσμα (variational inference) καταφέραμε να κάνουμε τον υπολογισμό του όρου  $P(Z|X)$  εφικτό. Στο δεύτερο σκέλος του variational autoencoder, τον αποκωδικοποιητή, θέλουμε να υπολογίσουμε τον όρο  $P_\theta(X|Z)$ . Θα χρησιμοποιήσουμε stochastic (ή mini-batch) αλγόριθμο ανοδικής (ή καθοδικής) κλίσης για να μάθουμε τις παραμέτρους  $\theta$ .

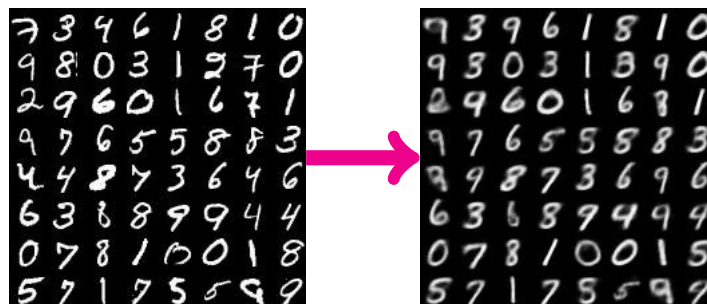
### 1.3.1 Πλεονεκτήματα του Variational Inference

(Max Welling, Diederik P. Kingma 2013. [3])

1. **Διαχείριση Μη Υπολογισιμότητας (Intractability):** Όπως εξηγήσαμε παραπάνω, ο υπολογισμός του όρου  $P(X)$  είναι αδύνατος και ως αποτέλεσμα ισχύει το ίδιο και για τον όρο  $P(Z|X)$ . Συνεπώς, ο αλγόριθμος **Expectation-Maximization (EM)** και άλλοι αλγόριθμοι τύπου Variational Bayes δεν μπορούν να χρησιμοποιηθούν. Το variational συμπίεσμα (variational inference) όμως, που δεν απαιτεί το υπολογισμό αυτού του όρου, μπορεί να δώσει τη λύση.
2. **Large datasets:** Αλγόριθμοι βελτιστοποίησης που εκπαιδεύονται σε όλο το σύνολο των δεδομένων σε κάθε επανάληψη (πχ. batch gradient descent) είναι πολύ δαπανηροί. Χάρη στο variational inference οι παράμετροι ενημερώνονται χρησιμοποιώντας μικρές δέσμες δεδομένων (mini-batches) ή ακόμα μεμονομένα στιγμιότυπα (πχ. stochastic gradient descent). Λύσεις που βασίζονται σε δειγματοληψία, όπως Monte Carlo EM είναι πολύ αργές γιατί εμπλέκουν ακριβούς βρόχους δειγματοληψίας για κάθε δεδομένο ξεχωριστά.

### 1.3.2 Μειονεκτήματα του Variational Inference

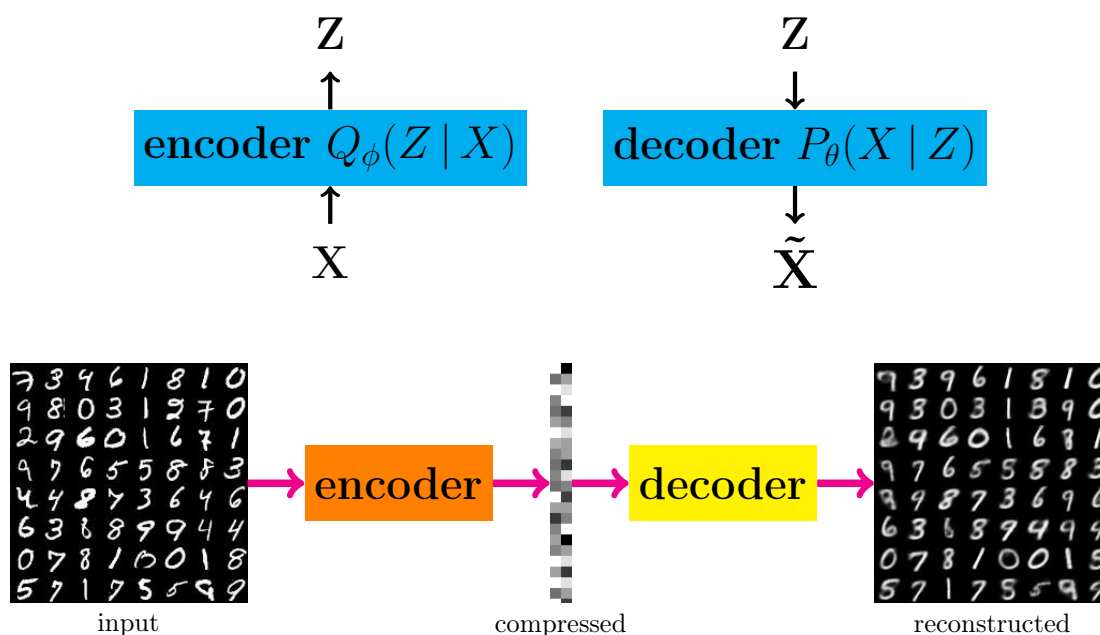
(Aaron Courville, Ian Goodfellow, Yoshua Bengio 2016. [4]) Η προσέγγιση του variational autoencoder είναι κομψή, θεωρητικά ευχάριστη και απλή στην υλοποίηση. Επίσης, επιφέρει σχεδόν τέλεια αποτελέσματα και είναι μέσα στις τεχνολογίες τελευταίας λέξης που αφορούν προσεγγίσεις σε γεννητική μοντελοποίηση (generative modelling). Το κύριο μειονέκτημα είναι ότι δείγματα από variational autoencoder εκπαιδευμένα σε εικόνες τείνουν να είναι μερικώς θολά. Οι αιτίες αυτού του φαινομένου δεν είναι ακόμα γνωστές και παραμένει να βρεθούν. Μια πιθανή εξήγηση είναι ότι η θολότητα είναι ένα εγγενές αποτέλεσμα της μέγιστης πιθανοφάνειας, που ελαχιστοποιεί την απόκλιση  $D_{KL}(P(Z|X) || Q(Z|X))$ . Αυτό σημαίνει ότι το μοντέλο θα αναθέσει μεγάλες πιθανότητες σε δεδομένα που ανήκουν στο σύνολο εκπαίδευσης, αλλά μπορεί επίσης να αναθέσει μεγάλες πιθανότητες και σε υπόλοιπα δεδομένα. Τα υπόλοιπα δεδομένα μπορεί να περιλαμβάνουν θολές εικόνες.



ΕΙΚΟΝΑ 1.1: Παράδειγμα δεδομένων εισόδου και εξόδου ενός VAE. Όπως παρατηρούμε, τα τελικά αποτελέσματα είναι μερικώς θολά.

### 1.3.3 Κύρια Βήματα του Variational Autoencoder

1. Λαμβάνουμε τα δεδομένα (εικόνες) εισόδου  $X$ , ως μεταβλητή μεγέθους  $N \times D$ .
2. Εκπαιδεύουμε τον κωδικοποιητή, που συμβολίζεται ως  $Q_\phi(Z | X)$ , χρησιμοποιώντας batch gradient descent (όπου  $\phi$  είναι οι παράμετροι του κωδικοποιητή και  $Z$  είναι οι λανθάνουσες μεταβλητές ελαττωμένης διάστασης).
3. Ταυτόχρονα με τον κωδικοποιητή, εκπαιδεύουμε τον αποκωδικοποιητή, που συμβολίζεται ως  $P_\theta(X | Z)$ , χρησιμοποιώντας batch gradient descent (όπου  $\theta$  είναι οι παράμετροι του αποκωδικοποιητή).
4. Είμαστε έτοιμοι να δείξουμε και να χρησιμοποιήσουμε τα καινούρια μας τεχνητά δεδομένα (εικόνες)  $\tilde{X}$ , τα οποία πρέπει να μοιάζουν με τα αρχικά δεδομένα.



ΕΙΚΟΝΑ 1.2: Σχήμα Variational Autoencoder.

## Κεφάλαιο 2

# Υπολογίζοντας το Variational Ελάχιστο Φράγμα (ELBO)

### 2.1 Χρησιμοποιώντας την Ανισότητα του Jensen για να Υπολογίσουμε το Variational Ελάχιστο Φράγμα (ELBO)

Εφόσον ο όρος  $P(Z|X)$  είναι μη υπολογίσιμος, όπως εξηγήσαμε στο [Κεφάλαιο 1](#), χρησιμοποιούμε μια αυθαίρετη κατανομή  $Q(Z)$  για να προσεγγίσουμε την πραγματική εκ των υστέρων κατανομή  $P(Z|X)$ .

$$\log P(X) = \log \int P(X, z) dz = \log \int P(X, z) \cdot \frac{Q(Z)}{Q(Z)} dz = \log(E_q[\frac{P(X, z)}{Q(Z)}])$$

Από την ανισότητα του Jensen για κοίλες συναρτήσεις έχουμε  $f(E[X]) \geq E[f(X)]$ . Εφόσον η λογαριθμική συνάρτηση είναι κοίλη, έχουμε:

$$\log P(X) \geq E_q[\log P(X, z)] - E_q[\log Q(Z)]$$

Ορίζουμε το ELBO ως εξής:

$$ELBO = L(X, Q) = E_q[\log P(X, z)] - E_q[\log Q(Z)]$$

Συμπερασματικά, είναι προφανές ότι το ELBO είναι το ελάχιστο φράγμα των λογαριθμικών πιθανοτήτων των παρατηρήσεων ( $\log P(X)$ ). Συνεπώς, αν σε ορισμένες περιπτώσεις επιθυμούμε να μεγιστοποιήσουμε την περιθώρια πιθανότητα, μπορούμε εναλλακτικά να μεγιστοποιήσουμε το variational ελάχιστο φράγμα (ELBO).

## 2.2 Χρησιμοποιώντας την KL Απόκλιση για να Υπολογίσουμε το ELBO

(Carl Doersch 2016. [2]) Στην περίπτωση των Variational Autoencoder, η KL-απόκλιση (Kullback-Leibler απόκλιση) είναι η πιθανοφάνεια μεταξύ της πραγματικής κατανομής των λανθανουσών μεταβλητών  $Z$ , δεδομένου των  $X$ ,  $P(Z|X)$ , και της προσεγγιστικής κατανομής των λανθανουσών μεταβλητών  $Z$ , δεδομένου των  $X$ ,  $Q_\phi(Z|X)$ . Για το δεύτερο όρο,  $Q_\phi(Z|X)$ , ισχύει η ακόλουθη εξίσωση:

$$Q_\phi(Z|X) \approx Q_\phi(Z)$$

Η KL απόκλιση μεταξύ δύο κατανομών, παίρνει την ακόλουθη μορφή:

$$\begin{aligned} D_{KL}[Q_\phi(Z) || P(Z|X)] &= E_{Z \sim Q}[\frac{\log Q_\phi(Z)}{\log P(Z|X)}] \Rightarrow \\ D_{KL}[Q_\phi(Z) || P(Z|X)] &= E_{Z \sim Q}[\log Q_\phi(Z) - \log P(Z|X)] \Rightarrow \end{aligned}$$

όπου το  $D$  συμβολίζει την KL-απόκλιση μεταξύ δύο κατανομών.

Αφού εφαρμόσουμε τον κανόνα του Bayes στον δεύτερο όρο, έχουμε:

$$\begin{aligned} D_{KL}[Q_\phi(Z) || P(Z|X)] &= E_{Z \sim Q}[\log Q_\phi(Z) - \log \frac{P_\theta(X|Z) \cdot P(Z)}{P(X)}] \Rightarrow \\ D_{KL}[Q_\phi(Z) || P(Z|X)] &= E_{Z \sim Q}[\log Q_\phi(Z) - \log P_\theta(X|Z) - \log P(Z) + \log P(X)] \Rightarrow \\ \log P(X) - D_{KL}[Q_\phi(Z) || P(Z|X)] &= E_{Z \sim Q}[\log P_\theta(X|Z)] - D_{KL}[Q_\phi(Z) || P(Z)] \Rightarrow \\ \log P(X) - D_{KL}[Q_\phi(Z|X) || P(Z|X)] &= E_{Z \sim Q}[\log P_\theta(X|Z)] - D_{KL}[Q_\phi(Z|X) || P(Z)] \end{aligned}$$

Η τελευταία εξίσωση είναι το variational ελάχιστο φράγμα, το οποίο ονομάζουμε **ELBO** από δω και στο εξής. Το αριστερό σκέλος της εξίσωσης έχει τον όρο που θέλουμε να μεγιστοποιήσουμε,  $P(X)$ , συν έναν όρο σφάλματος. Ο όρος σφάλματος είναι η Kullback-Leibler απόκλιση μεταξύ των  $Q_\phi(Z|X) \approx Q_\phi(Z)$  και  $P(Z|X)$ , που οδηγεί την κατανομή  $Q$  να παράγει λανθάνουσες μεταβλητές  $Z$ , δεδομένου των μεταβλητών εισόδου  $X$ . Θέλουμε να ελαχιστοποιήσουμε την Kullback-Leibler απόκλιση μεταξύ των δύο κατανομών. Το πρόβλημα ανάγεται στη μεγιστοποίηση του όρου ELBO. Αν η κατανομή  $Q$  έχει προσεγγιστεί με υψηλή ακρίβεια, τότε ο όρος του σφάλματος γίνεται μικρός.

Συνοπτικά, το ELBO προκύπτει από τον παρακάτω τύπο:

$$ELBO = L(X, Q) = \log P(X) - D_{KL}[Q_\phi(Z|X) || P(Z|X)] \Rightarrow$$

$$\log P(X) = L(X, Q) + D_{KL}[Q_\phi(Z|X) || P(Z|X)]$$

εφόσον η KL-απόκλιση είναι μη αρνητική, έχουμε:

$$\log P(X) \geq L(X, Q)$$

γι' αυτό ονομάζουμε το L, ELBO ή variational ελάχιστο φράγμα (Max Welling, Diederik P. Kingma 2013. [3]). Επίσης, το ELBO είναι ίσο με:

$$ELBO = L(X, Q) = E_{Z \sim Q}[\log P_\theta(X|Z)] - D_{KL}[Q_\phi(Z|X) || P(Z)] \Rightarrow_{Q_\phi(Z|X) \approx Q_\phi(Z)}$$

$$\mathbf{L(X, Q)} = \mathbf{E_{Z \sim Q}[\log P_\theta(X|Z)]} - \mathbf{D_{KL}[Q_\phi(Z) || P(Z)]}$$

Ο όρος  $E_{Z \sim Q}[\log P_\theta(X|Z)]$  ονομάζεται **κόστος ανακατασκευής** (reconstruction cost).

Ο όρος  $D_{KL}[Q_\phi(Z)||P(Z)]$  ονομάζεται **ποινή** (penalty) ή **όρος κανονικοποίησης** (regularization term). Ο όρος ποινής εξασφαλίζει ότι η εξήγηση των δεδομένων,  $Q_\phi(Z|X) \approx Q_\phi(Z)$  δεν αποκλίνει πολύ από τον όρο των παρατηρήσεων  $P(Z)$ . Επιπλέον, ο όρος ποινής μας βοηθά να εφαρμόσουμε το **Ξυράφι του Occam** στον κωδικοποιητή μας (inference model). Ο όρος αυτός είναι πάντα μεγαλύτερος ή ίσος του 0 και άρα μπορεί να παραλειφθεί.

Μπορούμε να εκπαιδεύσουμε ταυτόχρονα τον κωδικοποιητή και τον αποκωδικοποιητή.



ΕΙΚΟΝΑ 2.1: Το Elbo Room είναι ένα bar που βρίσκεται στην ιστορική περιοχή Mission District του San Francisco στην οδό Valencia Street μεταξύ 17th Street και 18th Street.

## Κεφάλαιο 3

# Αλγόριθμος Προς τα Πίσω Διάδοσης (Back-Propagation)

### 3.1 Βελτιστοποιώντας τη Συνάρτηση Κόστους (ELBO)

Όπως αναφέραμε στο προηγούμενο κεφάλαιο, εκπαιδεύουμε ταυτόχρονα τον κωδικοποιητή - (inference model)  $Q_\phi(Z|X)$ , και τον αποκωδικοποιητή - (generative model),  $P_\theta(X|Z)$ , βελτιστοποιώντας το variational lower bound (ELBO), χρησιμοποιώντας αλγόριθμο **προς τα πίσω διάδοσης** (gradient back-propagation). Καταλήξαμε λοιπόν στον ακόλουθο τύπο για το ELBO:

$$L(X, Q) = E_{Z \sim Q}[\log P_\theta(X|Z)] - D_{KL}[Q_\phi(Z) || P(Z)]$$

Οι κανόνες ενημέρωσης καθορίζονται βασισμένοι στον επιλεγμένο αλγόριθμο προς τα πίσω διάδοσης (back-propagation).

Τώρα, θα προσπαθήσουμε να βρούμε έναν κατάλληλο τύπο, για την KL-απόκλιση μεταξύ της κατανομής  $P(Z|X)$  και της κατανομής  $Q(Z|X)$ . (Carl Doersch 2016. [2])

έχουμε:

$$\mathbf{Q}_\phi(\mathbf{Z}) = N_1 = N(Z|\mu_1, \sigma_1^2) = N(Z|M, \Sigma^2), \text{ όπου: } \mu_1 = M \text{ και } \sigma_1 = \Sigma$$

$$\mathbf{P}(\mathbf{Z}) = N_2 = N(Z|\mu_2, \sigma_2^2) = N(Z|0, I), \text{ όπου: } \mu_2 = 0 \text{ και } \sigma_2 = I$$

ακόμα έχουμε:

$$\int \mathbf{Q}_\phi(\mathbf{Z}) \cdot \log \mathbf{P}(\mathbf{Z}) \, d\mathbf{z} = \int N(Z|M, \Sigma^2) \cdot \log N(Z|0, I) \, dz = -\frac{J}{2} \log 2\pi - \frac{1}{2} \cdot \sum_{j=1}^J (\mu_j^2 + \sigma_j^2) \Rightarrow$$

$$\int \mathbf{Q}_\phi(\mathbf{Z}) \cdot \log \mathbf{P}(\mathbf{Z}) \, d\mathbf{z} = -\frac{J}{2} \log 2\pi - \frac{1}{2} \cdot (M^2 + \Sigma^2) \quad (1)$$

και:

$$\int \mathbf{Q}_\phi(\mathbf{Z}) \cdot \log \mathbf{Q}_\phi(\mathbf{Z}) \, d\mathbf{z} = \int N(Z|M, \Sigma^2) \cdot \log N(Z|M, \Sigma^2) \, dz = -\frac{J}{2} \log 2\pi - \frac{1}{2} \cdot \sum_{j=1}^J (1 + \log \sigma_j^2) \Rightarrow$$

$$\int \mathbf{Q}_\phi(\mathbf{Z}) \cdot \log \mathbf{Q}_\phi(\mathbf{Z}) \, d\mathbf{z} = -\frac{J}{2} \log 2\pi - \frac{1}{2} \cdot (1 + \log \Sigma^2) \quad (2)$$

όπου  $J$  είναι η διάσταση των λανθανουσών μεταβλητών  $Z$ . Οι μέσες τιμές  $M$  και οι διασπορές  $\Sigma$  ορίζονται ως εξής:

$$M = \begin{bmatrix} M_{11} & M_{12} & M_{13} & \dots & M_{1J} \\ M_{21} & M_{22} & M_{23} & \dots & M_{2J} \\ \dots & \dots & \dots & \dots & \dots \\ M_{N1} & M_{N2} & M_{N3} & \dots & M_{NJ} \end{bmatrix}, \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} & \Sigma_{13} & \dots & \Sigma_{1J} \\ \Sigma_{21} & \Sigma_{22} & \Sigma_{23} & \dots & \Sigma_{2J} \\ \dots & \dots & \dots & \dots & \dots \\ \Sigma_{N1} & \Sigma_{N2} & \Sigma_{N3} & \dots & \Sigma_{NJ} \end{bmatrix}$$

όπου  $N$  είναι ο αριθμός των pixel

Συνεχίζοντας, μπορούμε να υπολογίσουμε την KL-απόκλιση μεταξύ των κατανομών  $P$  και  $Q$  από τον τύπο του ELBO, ως εξής:

$$D_{KL}[Q_\phi(Z|X) || P(Z|X)] = D_{KL}[N_1 || N_2] = D_{KL}[N(Z|\mu_1, \sigma_1^2) || N(Z|\mu_2, \sigma_2^2)] \Rightarrow$$

$$D_{KL}[Q_\phi(Z|X) || P(Z|X)] = D_{KL}[N(Z|\mu_1, \sigma_1^2) || N(Z|0, I)] \Rightarrow$$

$$D_{KL}[Q_\phi(Z|X) || P(Z|X)] = \int Q_\phi(Z) \cdot \log \frac{P(Z)}{Q_\phi(Z)} \, dz \Rightarrow$$

$$D_{KL}[Q_\phi(Z|X) || P(Z|X)] = \int Q_\phi(Z) \cdot (\log P(Z) - \log Q_\phi(Z)) \, dz \Rightarrow$$

$$D_{KL}[Q_\phi(Z|X) || P(Z|X)] = \int Q_\phi(Z) \cdot \log P(Z) - Q_\phi(Z) \cdot \log Q_\phi \, dz \Rightarrow_{(1),(2)}$$

$$D_{KL}[Q_\phi(Z|X) || P(Z|X)] = -\frac{J}{2} \log 2 \cdot \pi - \frac{1}{2} \cdot \sum_{j=1}^J (\mu_j^2 + \sigma_j^2) - (-\frac{J}{2} \log 2 \cdot \pi - \frac{1}{2} \cdot \sum_{j=1}^J (1 + \log \sigma_j^2)) \Rightarrow$$



$$D_{KL}[Q_\phi(Z|X) || P(Z|X)] = -\frac{J}{2} \log 2 \cdot \pi + \frac{J}{2} \cdot \log 2 \cdot \pi - \frac{1}{2} \cdot \sum_{j=1}^J (\mu_j^2 + \sigma_j^2) + \frac{1}{2} \cdot \sum_{j=1}^J (1 + \log \sigma_j^2) \Rightarrow$$

$$D_{KL}[Q_\phi(Z|X) || P(Z|X)] = \frac{1}{2} \cdot \sum_{j=1}^J (1 + \log \sigma_j^2 - \mu_j^2 - \sigma_j^2) \Rightarrow$$

$$D_{KL}[Q_\phi(Z|X) || P(Z|X)] = \frac{1}{2} \cdot (J + \log \Sigma^2 - M^2 - \Sigma^2)$$

αν η διάσταση της παραμέτρου  $J = 1$ , των λανθανουσών μεταβλητών  $Z$ , αυτό σημαίνει ότι έχουμε ενιαίες (univariate) κατανομές Gauss και καταλήγουμε στον εξής τύπο:

$$\mathbf{D}_{KL}[\mathbf{Q}_\phi(\mathbf{Z}|\mathbf{X}) || \mathbf{P}(\mathbf{Z}|\mathbf{X})] = \frac{1}{2} \cdot (1 + \log \Sigma^2 - M^2 - \Sigma^2)$$

Ας θυμηθούμε ότι ο όρος της KL-απόκλισης έχει αρνητικό πρόσημο στον τύπο του variational lower bound (ELBO), επομένως θέλουμε να τον ελαχιστοποιήσουμε.

### 3.2 Το Reparameterization Trick

(Carl Doersch 2016. [2]) Έχοντας καταλήξει σε έναν κατάλληλο τύπο για τον όρο της KL-απόκλισης, μένει ακόμα να τρέξουμε τον αλγόριθμο στοχαστικής (ή mini-batch) καθοδικής κλίσης (gradient descent) για διάφορες τιμές του  $X$ , παίρνοντας δείγματα από ένα dataset,  $D$ . Η πλήρης εξίσωση που θέλουμε να βελτιστοποιήσουμε είναι η εξής:

$$E_{X \sim D}[E_{Z \sim Q}[\log P_\theta(X|Z)] - D_{KL}[Q_\phi(Z) || P(Z)]] \quad (1)$$

Θέλουμε να βρούμε την παράγωγο αυτής της εξίσωσης. Το σύμβολο της παραγώγου μπορεί να μετακινηθεί εντός των μέσων τιμών. Επομένως, μπορούμε να δειγματοληπτήσουμε μία μόνο τιμή του  $X$  και μία μόνο τιμή του  $Z$  από την κατανομή  $Q(Z|X)$ , και να υπολογίσουμε την παράγωγο της εξίσωσης:

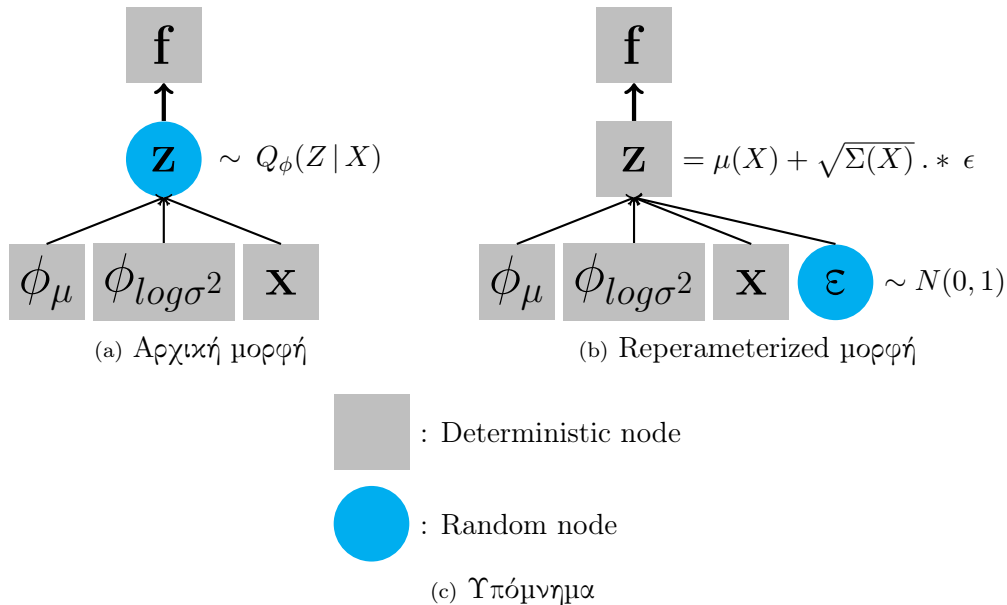
$$\log P_\theta(X|Z) - D_{KL}[Q_\phi(Z) || P(Z)] \quad (2)$$

Στη συνέχεια, μπορούμε να πάρουμε τη μέση τιμή της παραγώγου αυτής της συνάρτησης, για αυθαίρετα πολλά δείγματα  $X$  και  $Z$  και το αποτέλεσμα θα συγκλίνει στην παράγωγο της Συνάρτησης 1. Ωστόσο, υπάρχει ένα σημαντικό πρόβλημα με την Εξίσωση 2. Ο όρος  $E_{Z \sim Q}[\log P_\theta(X|Z)]$  δεν εξαρτάται μόνο

από τις παραμέτρους  $P$ , αλλά και από τις παραμέτρους  $Q$ . Εν τούτοις, στην Εξίσωση 2, η εξάρτηση αυτή έχει εξαφανιστεί! Για να δουλέψουν οι VAE, είναι ουσιώδες να οδηγήσουμε την κατανομή  $Q$  να παράγει κωδικοποιήσεις για την  $X$ , τις οποίες η  $P$  μπορεί αξιόπιστα να αποκωδικοποιήσει. Το προς τα εμπρός πέρασμα του δικτύου (forward pass) δουλεύει σωστά και παράγει τη σωστή μέση τιμή, αν η έξοδος (output) υπολογίζεται κατά μέσο όρο πολλών δειγμάτων  $X$  και  $Z$ . Ωστόσο, πρέπει να προωθήσουμε προς τα πίσω (back-propagate) το σφάλμα μέσω ενός επιπέδου που δειγματοληπτεί τη  $Z$  μέσω της κατανομής  $Q(Z|X)$ , η οποία είναι μη συνεχής διαδικασία και δεν έχει παράγωγο. Ο στοχαστικός αλγόριθμος καθοδικής κλίσης (stochastic gradient descent) μέσω της προς τα πίσω διάδοσης (back-propagation) μπορεί να χειριστεί στοχαστικές εισόδους (inputs), αλλά δεν μπορεί να χειριστεί μονάδες εντός του επιπέδου εισόδου (input layer). Δεδομένου των  $\mu(X)$  και  $\Sigma(X)$  - η μέση τιμή και η συνδιακύμανση της κατανομής  $Q(Z|X)$  - μπορούμε να δειγματοληπτήσουμε από την κανονική κατανομή  $N(\mu(X), \Sigma(X))$ , δειγματοληπτώντας πρώτα από την  $\epsilon \sim N(0, I)$ . Τελικά υπολογίζουμε τη μεταβλητή  $Z = \mu(X) + \sqrt{\Sigma(X)} \cdot \epsilon$ , που ακολουθεί κανονική κατανομή,  $Z \sim N(\mu(X), \Sigma(X))$ , εφόσον, κάθε γραμμικός μετασχηματισμός μιας Gaussian T.M. (τυχαίας μεταβλητής) είναι πάλι Gaussian (συμβουλευτείτε το Παράρτημα Β, [Κατανομή Gauss \(ή Κανονική Κατανομή\)](#)). Επομένως, η εξίσωση της οποίας υπολογίζουμε την παράγωγο είναι η εξής:

$$E_{X \sim D}[E_{Z \sim Q}[\log P_{\theta}(X|Z = \mu(X) + \sqrt{\Sigma(X)} \cdot \epsilon)] - D_{KL}[Q_{\phi}(Z) || P(Z)]]$$

Είναι σημαντικό να προσέξουμε ότι η κατανομή  $Q(Z|X)$  (συνεπώς και η  $P(Z)$ ) πρέπει να είναι συνεχής! Το reparameterization trick μας επιτρέπει να υπολογίσουμε την παράγωγο της μέσης τιμής του ELBO, έτσι ώστε το back-propagation να μπορεί να εφαρμοστεί.



ΕΙΚΟΝΑ 3.1: Το δεξιό δίκτυο χρησιμοποιεί το reparameterization trick. Το back-propagation μπορεί να εφαρμοστεί μόνο σε αυτό το δίκτυο.

### 3.3 Υπολογίζοντας τους κανόνες ενημέρωσης βαρών του κωδικοποιητή και του αποκωδικοποιητή

Στις υλοποιήσεις που παρουσιάζονται σε αυτή τη διατριβή, οι κανόνες ενημέρωσης των βαρών (Phis) του κωδικοποιητή και τα βάρη (Thetas) του αποκωδικοποιητή υπολογίζονται αυτόματα. Οι βιβλιοθήκες TensorFlow και PyTorch περιλαμβάνουν ενσωματωμένους αλγόριθμους προς τα πίσω διάδοσης (back-propagation) που υπολογίζουν τις παραγώγους για τα βάρη που χρησιμοποιούνται για τη διαδικασία της μάθησης. Οι κανόνες ενημέρωσης, επίσης περιλαμβάνονται ως ενσωματωμένες συναρτήσεις και εφαρμόζονται στη διαδικασία χρησιμοποιώντας τις παραγώγους που υπολογίστηκαν προηγουμένως.

Για να κάνουμε τη διαδικασία πιο ξεκάθαρη, ακολουθεί παράδειγμα κλήσης του ενσωματωμένου αλγόριθμου προς τα πίσω διάδοσης και της συνάρτησης ενημέρωσης βαρών σε TensorFlow:

```
var_list # the weights and the biases of the variational autoencoder
elbo # the loss function of the variational autoencoder to be minimized
lr # learning rate for the weights and the biases updates

# Adam Optimizer #
grads_and_vars = tf.train.AdamOptimizer(learning_rate=lr).
    compute_gradients(loss=elbo, var_list=var_list)
apply_updates = tf.train.AdamOptimizer(learning_rate=lr).
    apply_gradients(grads_and_vars=grads_and_vars)
```

ΚΩΔΙΚΑΣ 3.1: προς τα πίσω διάδοση σε TensorFlow

Ακολουθεί παράδειγμα κλήσης του ενσωματωμένου αλγόριθμου προς τα πίσω διάδοσης και της συνάρτησης ενημέρωσης βαρών σε PyTorch. Ας σημειώσουμε ότι και οι δύο διεργασίες εκτελούνται με την ίδια εντολή:

```
params # the weights and the biases of the variational autoencoder
lr # learning rate for the weights and the biases updates
solver = optim.Adam(params, lr=lr) # Adam Optimizer #
elbo_loss.backward() # Backward #
solver.step() # Update #
for p in params: # Housekeeping #
    # initialize the parameter gradients of the next epoch as zeros
    p.grad.data.zero_()
```

ΚΩΔΙΚΑΣ 3.2: προς τα πίσω διάδοση σε PyTorch

Για οδηγίες σχετικά με το πως εφαρμόζονται οι κανόνες ενημέρωσης, συμβουλευτείτε το Παράρτημα Α, στην ενότητα [Α.11 Κανόνες Ενημέρωσης Βαρών των Νευρωνικών Δικτύων](#).

**ΣΗΜΕΙΩΣΗ:** Όπως εξηγήσαμε στο προηγούμενο κεφάλαιο, θέλουμε να μεγιστοποιήσουμε το variational ελάχιστο φράγμα (ELBO). Ωστόσο, στις υλοποιήσεις που παρουσιάζονται σε αυτή τη διπλωματική εργασία, το ELBO γίνεται όλο και μικρότερο σε κάθε εποχή (epoch). Ο λόγος που αυτό συμβαίνει είναι επειδή χρησιμοποιούμε αλγόριθμο καθοδικής κλίσης (gradient descent) αντί για αλγόριθμο ανοδικής κλίσης (gradient ascent).

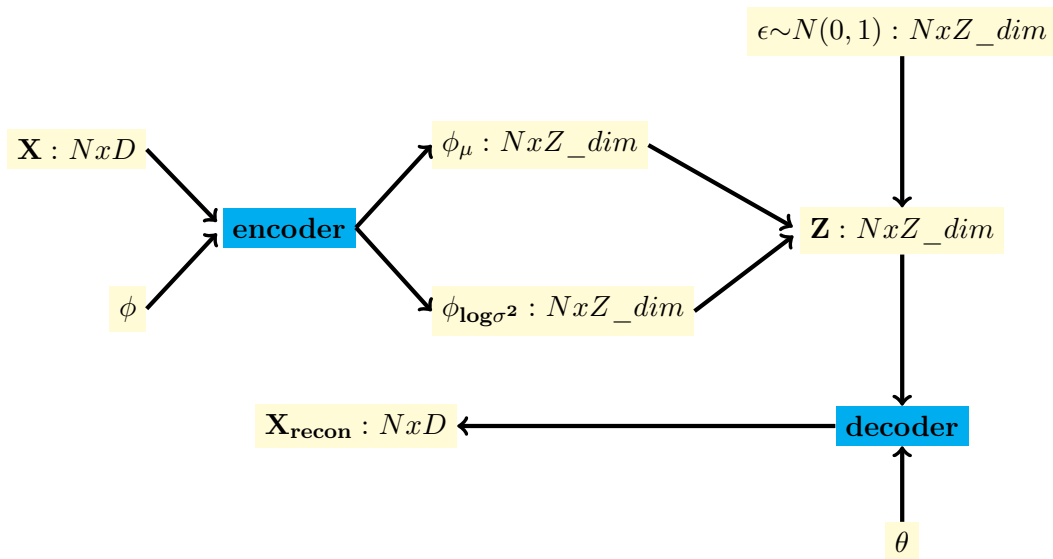
## Κεφάλαιο 4

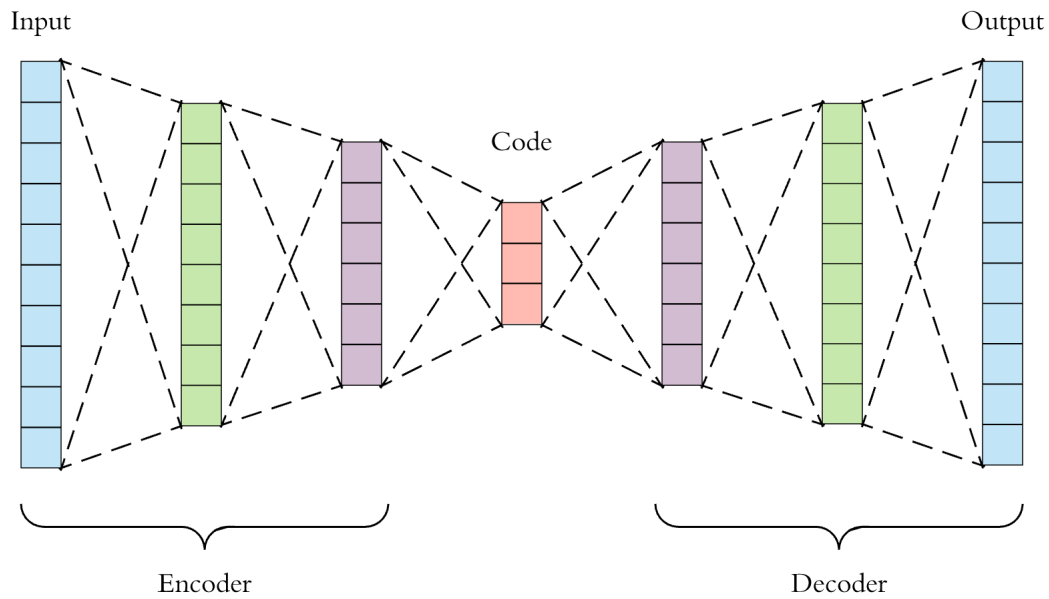
# Δομή του Variational Autoencoder & Αποτελέσματα Πειραμάτων

### 4.1 Δομή του Variational Autoencoder

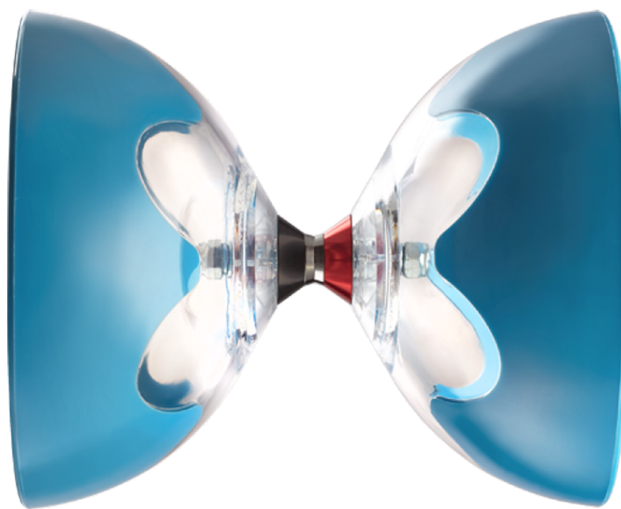
Συνήθως, ο variational autoencoder χρειάζεται πολλές εποχές (epochs) για να εκπαιδευτεί. Ωστόσο, υπάρχουν περιπτώσεις όπου πάρα πολλές εποχές μπορεί να επιφέρουν μη επιθυμητά αποτελέσματα, πχ. θολές εικόνες. Όταν το ελάχιστο φράγμα (ELBO) γίνει πολύ μικρό, ο VAE θα πρέπει να σταματήσει την εκπαίδευση, γιατί στην επόμενη εποχή ή επανάληψη, η τιμή του ELBO θα γίνει τελικά πολύ μεγάλη. Σε έναν autoencoder, ο αριθμός των νευρώνων του κωδικοποιητή (encoder), τον οποίο συμβολίζουμε ως  $M_1$ , πρέπει να είναι ίσος με τον αριθμό των νευρώνων του αποκωδικοποιητή (decoder), τον οποίο συμβολίζουμε ως  $M_2$ .

Το παρακάτω διάγραμμα απεικονίζει τα βήματα που ο variational autoencoder ακολουθεί για να κατασκευάσει τα τεχνητά δεδομένα,  $\mathbf{X}_{\text{recon}}$ , από τα πρωτότυπα δεδομένα,  $\mathbf{X}$ . Μπορούμε να εξετάσουμε πώς τα αρχικά δεδομένα,  $\mathbf{X}$ , μετατρέπονται στα λανθάνοντα δεδομένα,  $\mathbf{Z}$ , τα οποία είναι μια αναπαράσταση των μεταβλητών  $\mathbf{X}$  στην ελαττωμένη διάσταση  $\mathbf{Z}_{\text{dim}}$ :





(a) Δομή του autoencoder



(b) Ένα ζογκλερικό παιχνίδι diabolo

ΕΙΚΟΝΑ 4.1: Ο autoencoder ονομάζεται και δίκτυο Diabolo λόγω της εμφάνισης της δομής του. (Rumelhart et al., 1986a; Bourlard & Kamp, 1988; Hinton & Zemel, 1994; Schwenk & Milgram, 1995; Japkowicz, Hanson, & Gluck, 2000, Yoshua Bengio 2009. [1])

## 4.2 Το TensorBoard

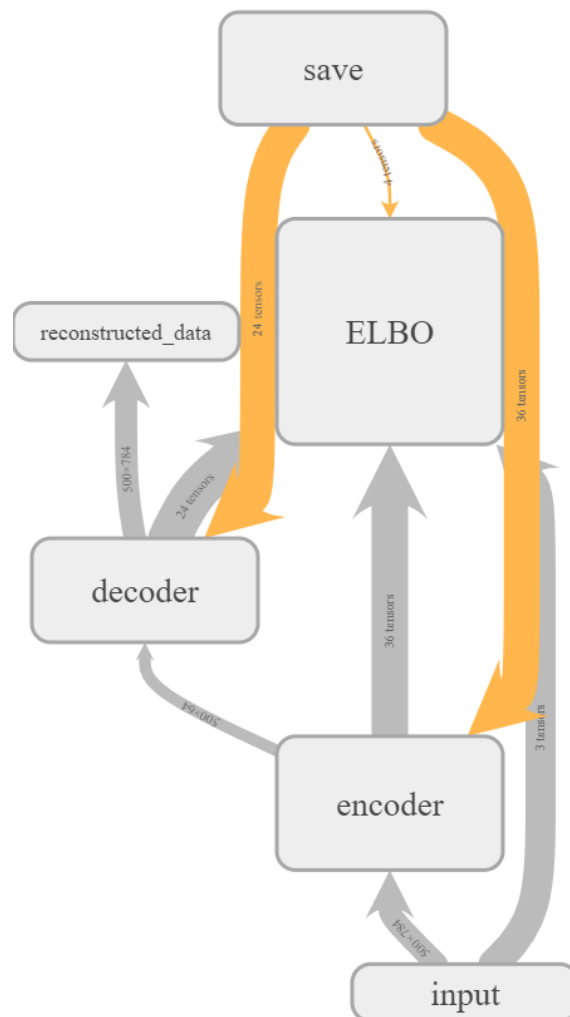
Για προγραμματιστικές υλοποιήσεις των variational autoencoder σε TensorFlow, PyTorch και Keras δείτε το [Παράρτημα Γ](#). Ανοίξτε ένα τερματικό (terminal), για Unix/Linux, ή τη γραμμή εντολών (command prompt), για Windows, και τρέξτε την ακόλουθη εντολή:

```
tensorboard --logdir = "path_to_the_graph"
```

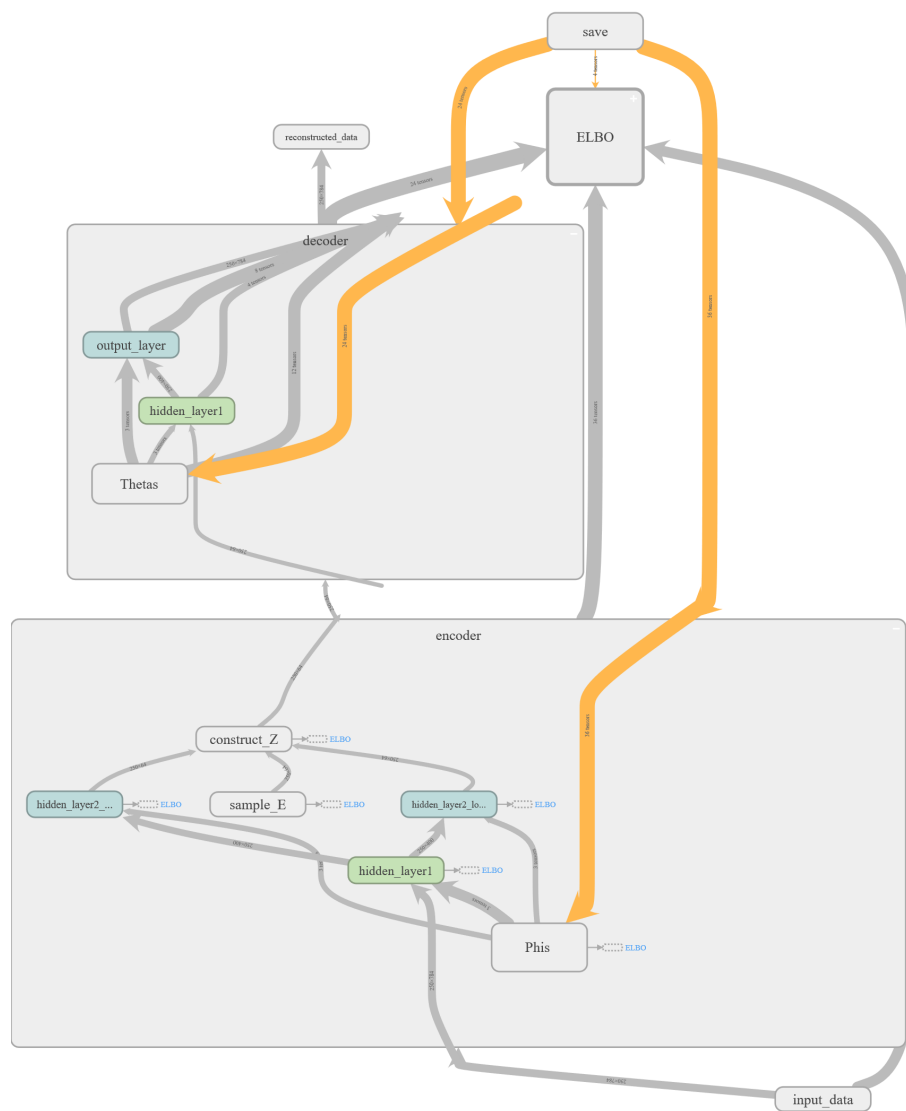
ΚΩΔΙΚΑΣ 4.1: εντολή για έναρξη του TensorBoard

Μετά, ένα μήνυμα θα πρέπει να εμφανιστεί που θα αναφέρεται στο ακόλουθο URL:

<http://localhost:6006>. Ανοίξτε ένα φυλλομετρητή (πχ. Firefox) και πλοηγηθείτε στο URL. Το TensorFlow τώρα τρέχει ως εφαρμογή δικτύου (web), από προεπιλογή στο port 6006.



ΕΙΚΟΝΑ 4.2: Αυτή είναι απεικόνιση της υλοποίησης σε TensorFlow, με το TensorBoard.



ΕΙΚΟΝΑ 4.3: Αυτή είναι διευρυμένη απεικόνιση της υλοποίησης σε TensorFlow, με το TensorBoard.

**ΣΗΜΕΙΩΣΗ:** Το port "6006" στο URL <http://localhost:6006> είναι το "goog" ανάποδα.



### 4.3 Σύνολα Δεδομένων (Datasets)

Ακολουθεί η κατάσταση όλων των συνόλων δεδομένων (dataset) που χρησιμοποιήθηκαν σε αυτή τη διπλωματική διατριβή:

Dataset	# TRAIN	# TEST	# VALIDATION
MNIST [5]	55000	10000	5000
Binarized MNIST	50000	10000	5000
CIFAR-10 [6]	50000	10000	-
OMNIGLOT	390 [En] / 360 [Gr]	130 [En] / 120 [Gr]	-
Cropped YALE Faces [7]	2442	-	-
ORL Face Database	400	-	-
MovieLens 100k	90570	9430	-

ΠΙΝΑΚΑΣ 4.1: Λεπτομέρειες των dataset.

Dataset	# Classes	Dimensions	Type of Values
MNIST [5]	10	28x28 pixels	real values in $[0, 1]$
Binarized MNIST	10	28x28 pixels	$\{0, 1\}$
CIFAR-10 [6]	10	32x32x3 [RGB] / 32x32x1 [Grayscaled] pixels	real values in $[0, 1]$
OMNIGLOT	26 [En] / 24 [Gr]	32x32 pixels	real values in $[0, 1]$
Cropped YALE Faces [7]	38	168x192 pixels	real values in $[0, 255]$
ORL Face Database	40	92x112 pixels	real values in $[0, 255]$
MovieLens 100k	943 users	1682 movies	$\{1, 2, 3, 4, 5\}$

ΠΙΝΑΚΑΣ 4.2: Πλήθος κατηγοριών, διαστάσεις και είδη τιμών των dataset.

Dataset	URL
MNIST [5]	<a href="http://yann.lecun.com/exdb/mnist">http://yann.lecun.com/exdb/mnist</a>
Binarized MNIST	<a href="https://github.com/yburda/iwae/tree/master/datasets/BinaryMNIST">https://github.com/yburda/iwae/tree/master/datasets/BinaryMNIST</a>
CIFAR-10 [6]	<a href="https://www.cs.toronto.edu/~kriz/cifar.html">https://www.cs.toronto.edu/~kriz/cifar.html</a>
OMNIGLOT	<a href="https://github.com/yburda/iwae/tree/master/datasets/OMNIGLOT">https://github.com/yburda/iwae/tree/master/datasets/OMNIGLOT</a>
Cropped YALE Faces [7]	<a href="http://vision.ucsd.edu/extyaleb/CroppedYaleBZip/CroppedYale.zip">http://vision.ucsd.edu/extyaleb/CroppedYaleBZip/CroppedYale.zip</a>
ORL Face Database	<a href="http://www.cl.cam.ac.uk/research/dtg/attarchive/facedatabase.html">http://www.cl.cam.ac.uk/research/dtg/attarchive/facedatabase.html</a>
MovieLens 100k	<a href="https://grouplens.org/datasets/movielens">https://grouplens.org/datasets/movielens</a>

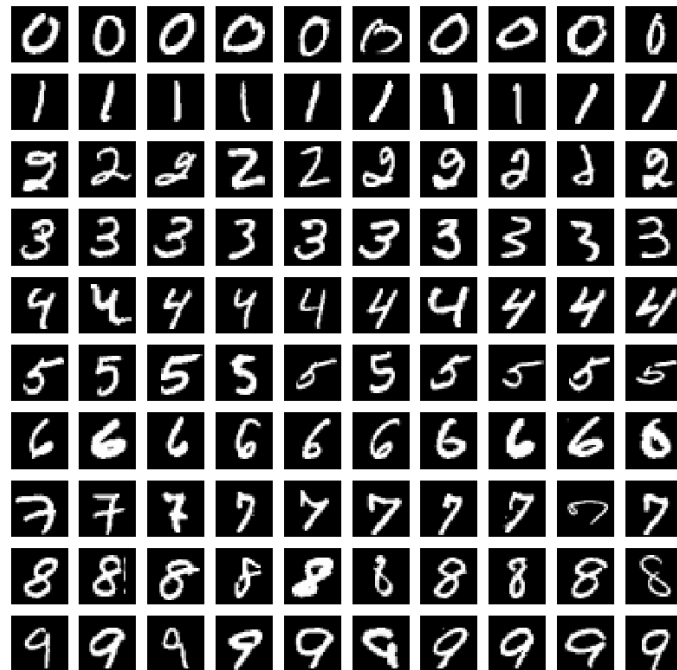
ΠΙΝΑΚΑΣ 4.3: Σύνδεσμοι για τα dataset.

#### ΣΗΜΕΙΩΣΕΙΣ:

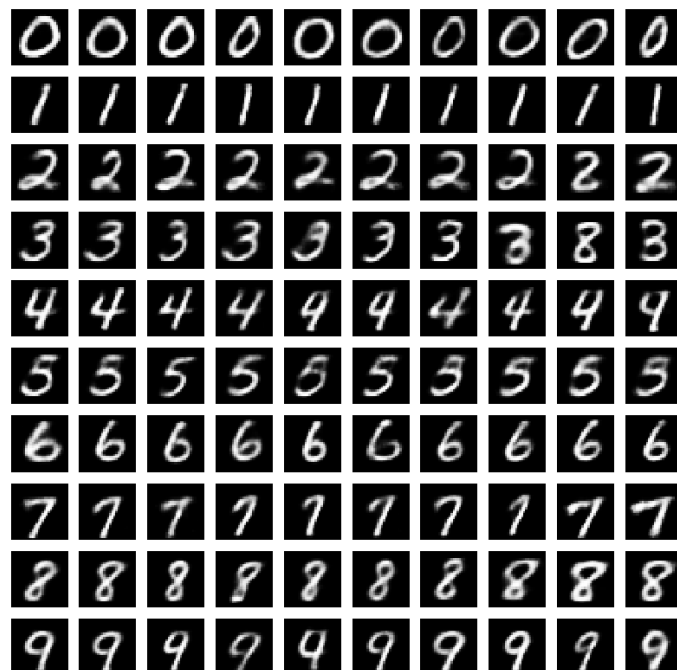
- Ο σύνδεσμος για το dataset YALE faces πιθανόν να μην είναι έγκυρος.
- Το dataset "ORL Face Database" περιλαμβάνει πολύ λίγα παραδείγματα και ως εκ τούτου τα αποτελέσματα είναι ανακριβή. Δεν συνιστούμε αυτό το dataset στα πειράματά μας. Ωστόσο, το dataset υπόσχεται καλά αποτελέσματα μόνο στον [K-NN αλγόριθμο συμπλήρωσης ελλিপών τιμών](#).
- Όλα τα dataset που παίρνουν τιμές στο διάστημα  $[0, 255]$  κανονικοποιούνται στο διάστημα  $[0, 1]$ , για ευκολότερους υπολογισμούς στη διαδικασία εκπαίδευσης.

## 4.4 Αποτελέσματα Πειραμάτων

Ακολουθούν τα αποτελέσματα του VAE αλγορίθμου, με χρήση TensorFlow, στο MNIST dataset:



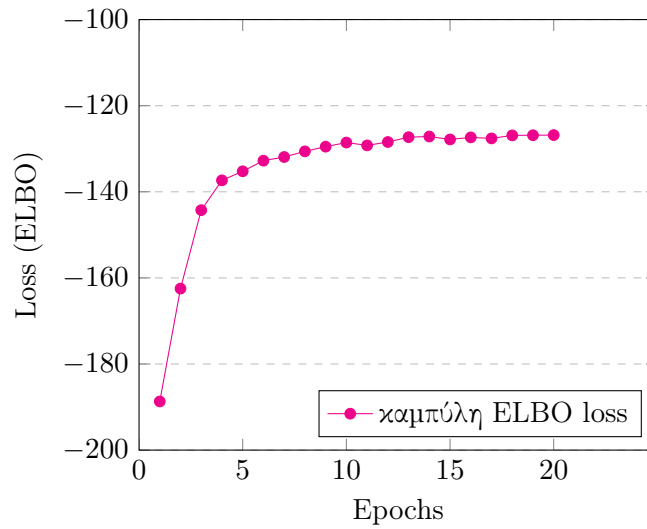
(a) Αρχικά Δεδομένα



(b) VAE Εποχή 20

ΕΙΚΟΝΑ 4.4: MNIST VAE σε TensorFlow.  
root mean squared error: 0.142598

ELBO σε κάθε εποχή (epoch), VAE στο MNIST dataset σε TensorFlow

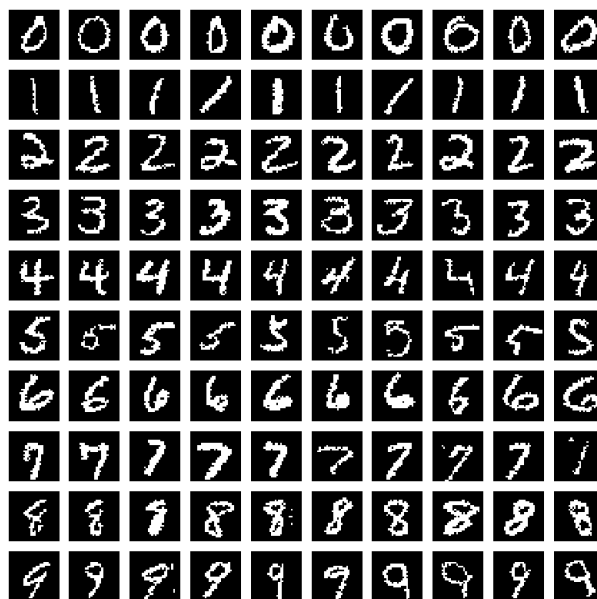


Metric	Value
last epoch loss (ELBO)	-126.83
root mean squared error (RMSE)	0.17121448655010216
mean absolute error (MAE)	0.07299246278044173

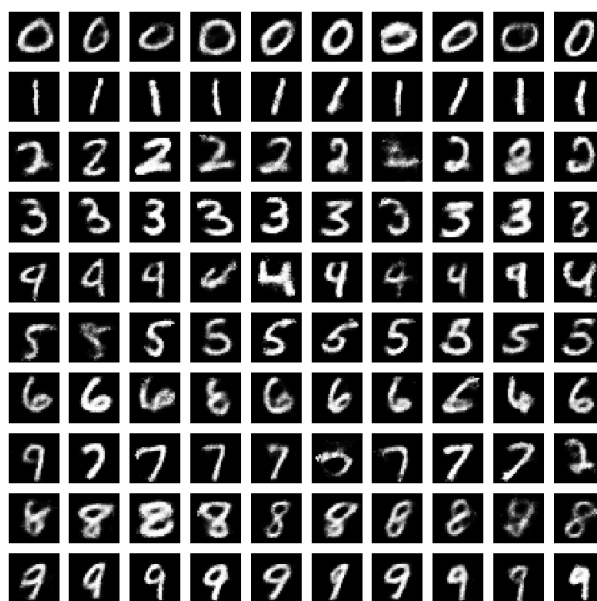
ΠΙΝΑΚΑΣ 4.4: VAE στο MNIST dataset σε TensorFlow.

Όπως εξηγήσαμε στο Κεφάλαιο 3, [Ενότητα 2](#), το ELBO ελαχιστοποιείται σε κάθε εποχή αντί να μεγιστοποιείται, επειδή χρησιμοποιήσαμε αλγόριθμο καθοδικής κλίσης (gradient descent), αντί για αλγόριθμο ανοδικής κλίσης (gradient ascent). Επίσης, μπορούμε να συμπεράνουμε ότι το MAE είναι καλύτερος δείκτης από το RMSE, επειδή το RMSE είναι μικρότερο του 1. Αυτό σημαίνει ότι το τετράγωνό του (δηλαδή το MSE) είναι μικρότερο από το RMSE και επίσης μικρότερο από 1.

Ακολουθούν αποτελέσματα του VAE αλγορίθμου, χρησιμοποιώντας **TensorFlow**, στο **Binarized MNIST** dataset:



(a) Αρχικά Δεδομένα



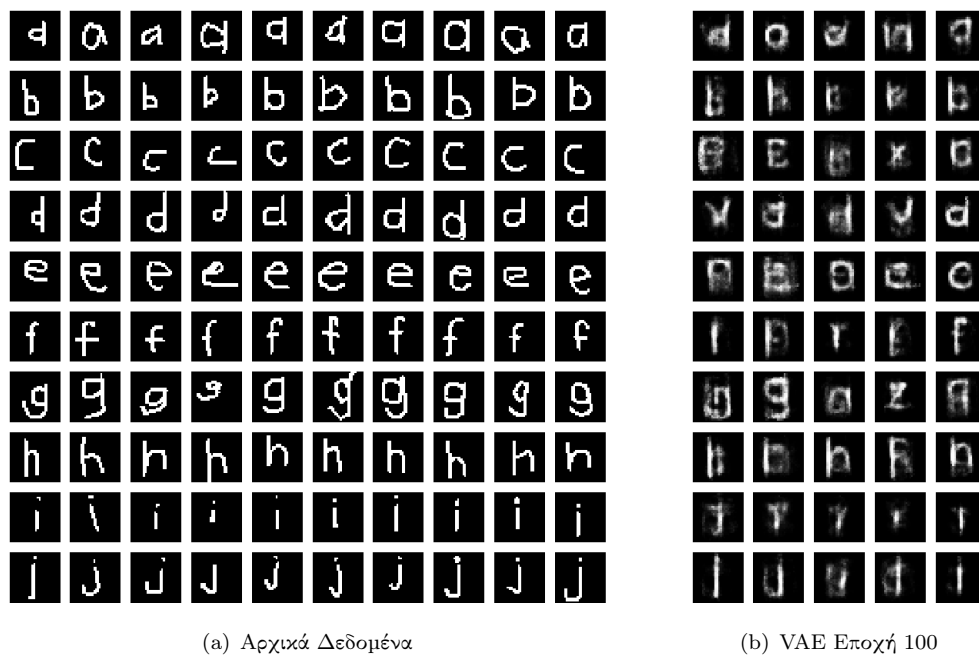
(b) VAE Εποχή 50

ΕΙΚΟΝΑ 4.5: MNIST VAE σε TensorFlow.

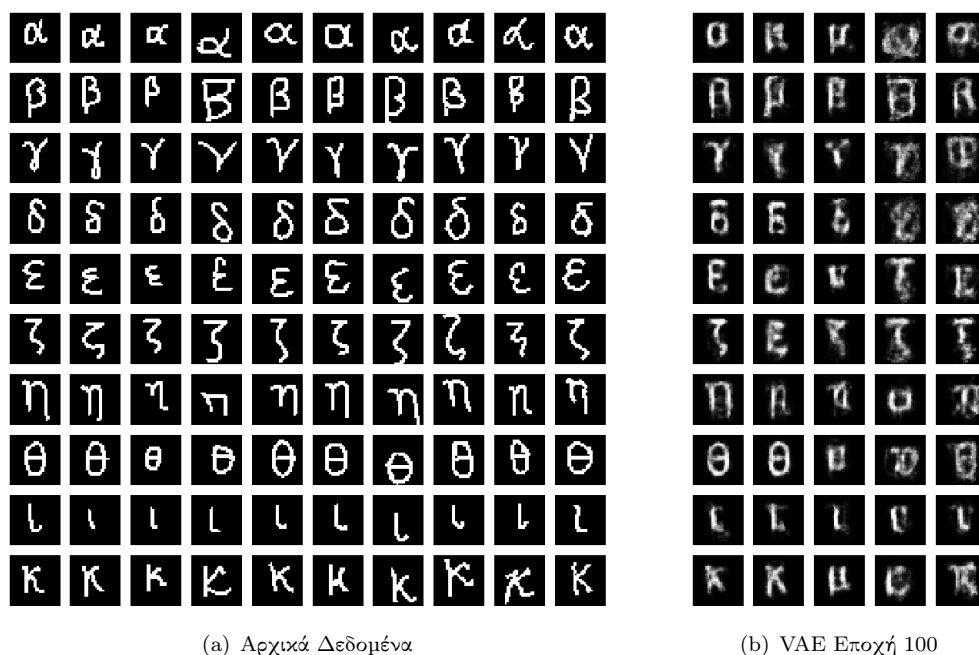
Metric	Value
last epoch loss (ELBO)	-109.82
root mean squared error (RMSE)	0.186880795395
mean absolute error (MAE)	0.0685835682327

ΠΙΝΑΚΑΣ 4.5: VAE στο Binarized MNIST dataset σε TensorFlow.

Ακολουθούν αποτελέσματα του VAE αλγορίθμου, με χρήση **PyTorch**, στο **OMNIGLOT** dataset:

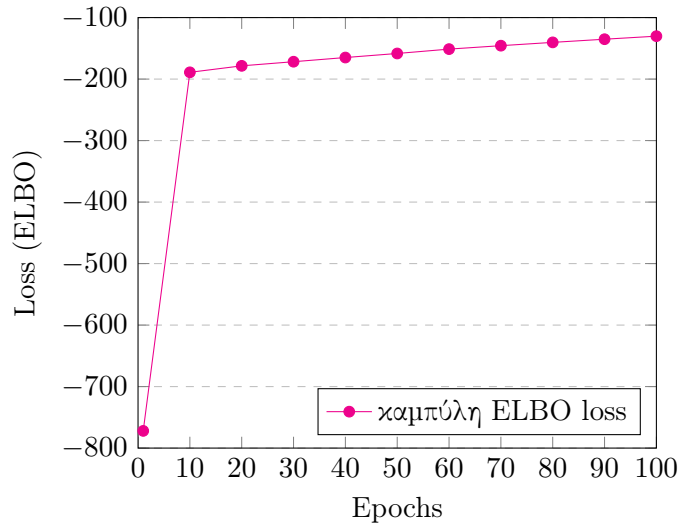


ΕΙΚΟΝΑ 4.6: OMNIGLOT Αγγλικό αλφάβητο, χαρακτήρες 1-10, αρχικοί και ανακατασκευασμένοι.



ΕΙΚΟΝΑ 4.7: OMNIGLOT Ελληνικό αλφάβητο, χαρακτήρες 1-10, αρχικοί και ανακατασκευασμένοι.

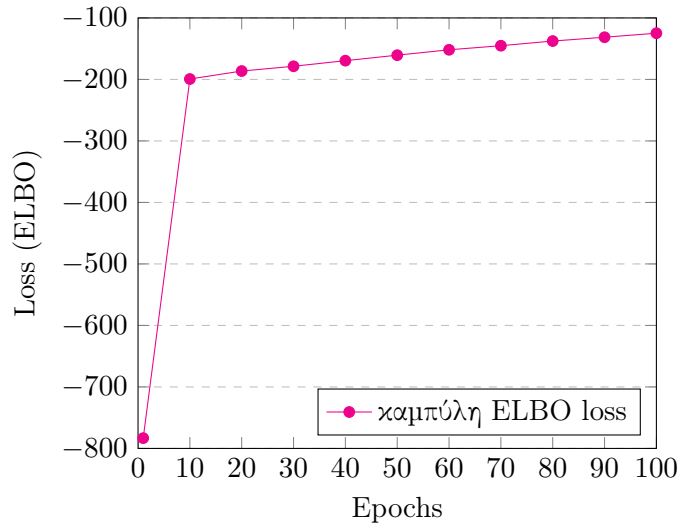
ELBO σε κάθε εποχή (epoch), VAE στο OMNIGLOT Αγγλικό dataset, σε PyTorch



Metric	Value
last epoch loss (ELBO)	-130.22
root mean squared error (RMSE)	0.2134686176531402
mean absolute error (MAE)	0.09105986614619989

ΠΙΝΑΚΑΣ 4.6: VAE στο OMNIGLOT Αγγλικό dataset σε PyTorch.

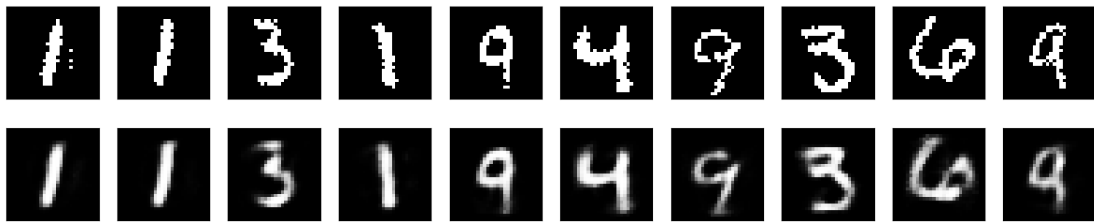
ELBO σε κάθε εποχή (epoch), VAE στο OMNIGLOT Ελληνικό dataset, σε PyTorch



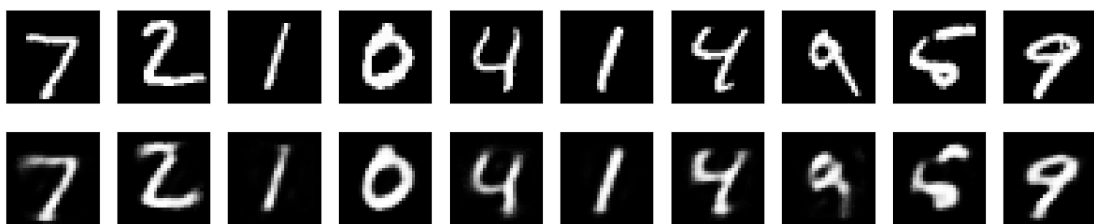
Metric	Value
last epoch loss (ELBO)	-124.85
root mean squared error (RMSE)	0.20157381435292054
mean absolute error (MAE)	0.08379960438030741

ΠΙΝΑΚΑΣ 4.7: VAE στο OMNIGLOT Ελληνικό dataset σε PyTorch.

Ακολουθούν αποτελέσματα του VAE αλγορίθμου, με χρήση του **Keras**, πάνω στα **Binarized MNIST**, **MNIST**, **CIFAR-10** και **OMNIGLOT** dataset:

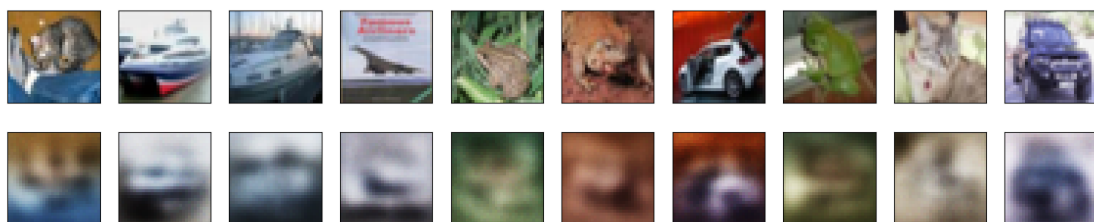


(a) Binarized MNIST

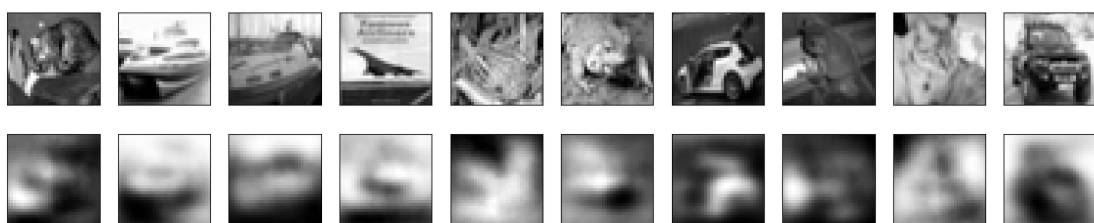


(b) MNIST

ΕΙΚΟΝΑ 4.8: Binarized MNIST και (Real-valued) MNIST ψηφία, αρχικά και ανακατασκευασμένα.

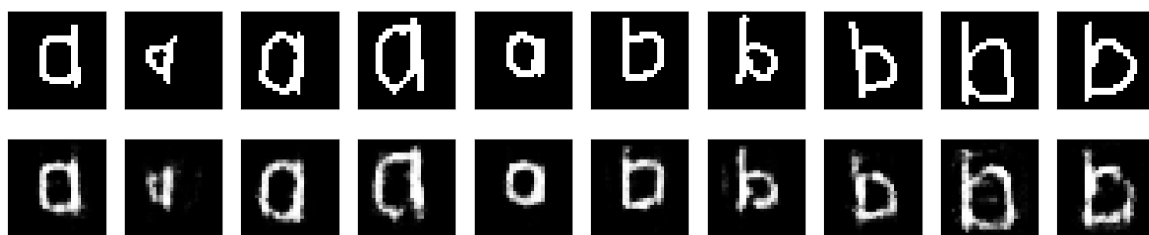


(a) CIFAR-10 RGB

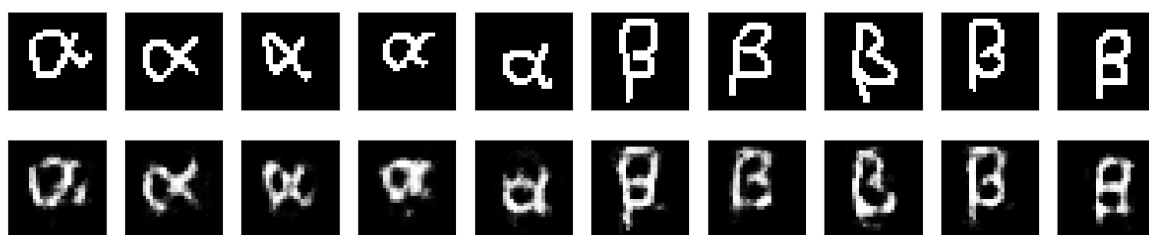


(b) CIFAR-10 Grayscale

ΕΙΚΟΝΑ 4.9: CIFAR-10 εικόνες, RGB και grayscale, αρχικές και ανακατασκευασμένες.



(a) OMNIGLOT Αγγλικό αλφάβητο



(b) OMNIGLOT Ελληνικό αλφάβητο

ΕΙΚΟΝΑ 4.10: OMNIGLOT Αγγλικοί και Ελληνικοί χαρακτήρες, αρχικοί και ανακατασκευασμένοι.

**ΣΗΜΕΙΩΣΗ:** Οι εικόνες στο CIFAR-10 RGB dataset περιλαμβάνουν pixel με τρία κανάλια (κόκκινο, πράσινο και μπλε) τα οποία παίρνουν πραγματικές τιμές, επομένως ο VAE δεν αναμένεται να έχει ικανοποιητικά αποτελέσματα σε αυτές.

Metric	Value
last epoch loss (ELBO)	0.0635
root mean squared error (RMSE)	0.1361456340713906
mean absolute error (MAE)	0.03993093576275391

ΠΙΝΑΚΑΣ 4.8: VAE στο Binarized MNIST dataset, σε Keras.

Metric	Value
last epoch loss (ELBO)	0.0040
root mean squared error (RMSE)	0.00102286
mean absolute error (MAE)	0.00063948194

ΠΙΝΑΚΑΣ 4.9: VAE στο MNIST dataset, σε Keras.

Metric	Value
last epoch loss (ELBO)	0.6198
root mean squared error (RMSE)	0.17352526
mean absolute error (MAE)	0.1368697

ΠΙΝΑΚΑΣ 4.10: VAE στο CIFAR-10 RGB dataset, σε Keras.



Metric	Value
last epoch loss (ELBO)	0.6166
root mean squared error (RMSE)	0.14963366
mean absolute error (MAE)	0.11517575

ΠΙΝΑΚΑΣ 4.11: VAE στο CIFAR-10 Grayscale dataset, σε Keras.

Metric	Value
last epoch loss (ELBO)	0.1764
root mean squared error (RMSE)	0.23234108227525616
mean absolute error (MAE)	0.11663868188603484

ΠΙΝΑΚΑΣ 4.12: VAE στο OMNIGLOT Αγγλικό dataset, σε Keras.

Metric	Value
last epoch loss (ELBO)	0.1840
root mean squared error (RMSE)	0.23350568189372273
mean absolute error (MAE)	0.1205574349168227

ΠΙΝΑΚΑΣ 4.13: VAE στο OMNIGLOT Ελληνικό dataset, σε Keras.

## ΠΑΡΑΤΗΡΗΣΕΙΣ:

- Τα ELBO loss που υπολογίστηκαν από το VAE σε Keras είναι σημαντικά χαμηλότερες από τις απώλειες που επέφεραν οι υλοποιήσεις σε TensorFlow και PyTorch.
- Οι εικόνες που ανακατασκευάστηκαν από το VAE με το MNIST dataset είναι πιο κοντά στα αρχικά δεδομένα από τις εικόνες που ανακατασκευάστηκαν με το Binarized MNIST dataset. Ο δείκτης σφάλματος (RMSE, MAE) όπως και το ELBO για το MNIST dataset είναι πολύ κοντά στο 0, πράγμα που σημαίνει ότι ο VAE έχει κάνει overfit στα δεδομένα εκπαίδευσης. Αυτό το αποτέλεσμα δεν ήταν αναμενόμενο καθώς ο VAE θα έπρεπε να συμπεριφέρεται καλύτερα στα δυαδικά δεδομένα. Η εξήγηση γι' αυτό πρέπει να αναζητηθεί στη λειτουργία της βιβλιοθήκης Keras και στους αλγόριθμους που χρησιμοποιεί για να υπολογίσει τη συνάρτηση απώλειας.
- Τα αποτελέσματα στο CIFAR-10 dataset είναι καλύτερα για τις Grayscaled εικόνες αντί των RGB εικόνων.
- Τα αποτελέσματα στο OMNIGLOT dataset είναι ελαφρώς καλύτερα για το dataset της Αγγλικής γλώσσας έναντι του dataset της Ελληνικής γλώσσας, κυρίως λόγω έλλειψης σε δεδομένα. Υπάρχουν 520 εικόνες Αγγλικών χαρακτήρων και 480 εικόνες Ελληνικών χαρακτήρων.

## Κεφάλαιο 5

# Αλγόριθμοι Συμπλήρωσης Ελλιπών Τιμών & Variational Autoencoders

### 5.1 Ένας Απλός Αλγόριθμος Συμπλήρωσης Ελλιπών Τιμών Χρησιμοποιώντας K-NN Συνεργατικό Φιλτράρισμα (CF)

Ο παραδοσιακός αλγόριθμος για ένα σύστημα συστάσεων (recommendation system) είναι το συνεργατικό φιλτράρισμα (collaborative filtering). Μια πολύ συνηθισμένη εφαρμογή για αυτή τη μέθοδο είναι η πρόβλεψη βαθμολογίας ενός χρήστη για ταινίες που δεν έχει ακόμα βαθμολογήσει, με βάση την ομοιότητα μεταξύ των βαθμολογιών του χρήστη για άλλες ταινίες και τις βαθμολογίες άλλων χρηστών. Η ομοιότητα συνημιτόνων (cosine similarity) είναι ένας τρόπος να μετρήσουμε την ομοιότητα μεταξύ χρηστών. Στον αλγόριθμο αυτό, θα χρησιμοποιήσουμε Ευκλίδειες αποστάσεις. Το **MovieLens** dataset είναι το πιο δημοφιλές σε σχετικές περιπτώσεις. Το **MNIST** dataset μπορεί επίσης να είναι μια κατάλληλη εφαρμογή του αλγορίθμου, με μικρές τροποποιήσεις.

Μπορούμε να εφαρμόσουμε συνεργατικό φιλτράρισμα στο **MNIST** dataset, χρησιμοποιώντας μια παραλλαγή του K-NN (K Nearest Neighbors) αλγορίθμου για regression, ως εξής:

1. Αποθηκεύουμε το MNIST dataset  $\mathbf{X\_train}$  με εικόνες από ψηφία στη μνήμη. Οι διαστάσεις των δεδομένων εκπαίδευσης,  $\mathbf{X\_train}$  είναι  $N \times D$ , όπου  $N$  είναι ο αριθμός των εικόνων εκπαίδευσης και  $D$  είναι ο αριθμός των pixel σε κάθε εικόνα. Για το MNIST dataset:  $D = 784$ . Επίσης, αποθηκεύουμε τα δεδομένα ελέγχου σε μια μεταβλητή,  $\mathbf{X\_test}$ .
2. Τροποποιούμε τα δεδομένα εκπαίδευσης (train) και ελέγχου (test), εισάγοντας ελλιπείς τιμές (missing values). Σκεφτήκαμε δύο τρόπους να επιλέξουμε και να κατασκευάσουμε τις ελλιπείς τιμές. Ένας τρόπος είναι να αντικαταστήσουμε όλα τα δεξιά, αριστερά, επάνω και κάτω pixel ή καθόλου pixel σε κάθε εικόνα με τιμές που δείχνουν ότι τα pixel είναι ελλιπή. Ένας δεύτερος τρόπος είναι να διαλέξουμε τυχαία pixel από κάθε εικόνα και να τα αντικαταστήσουμε με την προκαθορισμένη τιμή για ελλιπή pixel. Αν τα pixel παίρνουν τιμές στο διάστημα  $[0, 1]$  (1 για μαύρο pixel και 0 άσπρο pixel), μια καλή επιλογή για την τιμή που θα αναπαραστήσει τα ελλιπή δεδομένα είναι το 0.5, που αντιστοιχεί σε pixel με γκρι χρώμα.
3. Για τα δεδομένα ελέγχου, διαλέγουμε τα κοντινότερα  $\mathbf{K}$  δεδομένα (εικόνες) τα οποία ο αλγόριθμος θα λάβει υπόψη. Ο αλγόριθμος ενός κοντινότερου γείτονα αντιστοιχεί σε  $\mathbf{K}=1$ .

4. Για κάθε δεδομένο ελέγχου υπολογίζουμε τις Ευκλείδειες αποστάσεις προς κάθε δεδομένο εκπαιδευσης, βρίσκουμε τα **K** κοντινότερα και αποθηκεύουμε τα ευρετήρια (indices) των K κοντινότερων δεδομένων εκπαιδευσης σε μια μεταβλητή. Πρέπει όμως να προσέξουμε να μην προτιμήσουμε κοντινότερους γείτονες με ελλιπείς τιμές. Για να προσπεράσουμε αυτό το πρόβλημα, προτείνουμε την παρακάτω διαδικασία:

Κάνουμε τα pixel των δεδομένων εκπαιδευσης ίσα με την ελλιπή τιμή (missing value), στα σημεία όπου το τρέχον δεδομένο ελέγχου (test instance) έχει ελλιπείς τιμές. Αποθηκεύουμε το αποτέλεσμα στη μεταβλητή με όνομα **X\_train\_common**. Έτσι, η διαφορά **X\_train\_common - X\_test\_common**, στα σημεία όπου το τρέχον δεδομένο ελέγχου έχει ελλιπείς τιμές, θα είναι 0. Η μεταβλητή **X\_test\_common** θα κατασκευαστεί αμέσως μετά.

```
X_test_i # the current test instance
s = np.where(X_test_i == missing_value)
X_train_common = np.array(X_train)
X_train_common[:, s] = missing_value
```

ΚΩΔΙΚΑΣ 5.1: X\_train\_common

Επαναλαμβάνουμε το τρέχον δεδομένο ελέγχου Ntrain φορές και κάνουμε τις τιμές των δεδομένων ελέγχου ίσες με τον μέσο όρο όλων των δεδομένων εκπαιδευσης, εκεί όπου το τρέχον δεδομένο εκπαιδευσης κάθε φορά έχει ελλιπείς τιμές.

Αποθηκεύουμε το αποτέλεσμα σε μια μεταβλητή με όνομα **X\_test\_common**.

```
Ntrain # number of train examples
D # number of pixels (aka dimensions)
X_test_i # the current test instance
X_test_common = np.zeros((Ntrain, D))
mean_values = np.mean(X_train, axis=0) # 1 x D array
for k in range(Ntrain):
    X_test_common[k, :] = X_test_i
    s = np.where(X_train[k, :] == missing_value)
    if len(s[0]) != 0:
        X_test_common[k, s] = mean_values[s]
```

ΚΩΔΙΚΑΣ 5.2: X\_test\_common

5. Εξάγουμε τα δεδομένα (εικόνες) των οποίων τα ευρετήρια (indices) αντιστοιχούν στα σημεία των K κοντινότερων δεδομένων εκπαιδευσης, από τα δεδομένα εκπαιδευσης **X\_train** και τα αποθηκεύουμε σε μια μεταβλητή.
6. Προτείνουμε δύο λύσεις για να προχωρήσουμε περαιτέρω.
- Μπορούμε να χρησιμοποιήσουμε συντελεστές βαρών, ανάλογα με τις αποστάσεις μεταξύ του τρέχοντος παραδείγματος ελέγχου και τα K κοντινότερα δεδομένα εκπαιδευσης. Το κοντινότερο δεδομένο παίρνει το βάρος με τη μεγαλύτερη τιμή. Κάθε pixel του  $k$ -οστού δεδομένου εκπαιδευσης πολλαπλασιάζεται με το αντίστοιχο βάρος του  $w_k$ . Οι τιμές των

βαρών ανατίθενται ως εξής:

$$\mathbf{w}_k = \text{softmax}(-\mathbf{d}_k) = \frac{e^{-\mathbf{d}_k}}{\sum_{i=1}^K e^{-\mathbf{d}_i}},$$

$$\text{όπου: } w_1 \geq w_2 \geq w_3 \geq \dots \geq w_K,$$

$$\text{και } w_1 + w_2 + w_3 + \dots + w_K = \sum_{j=1}^K w_j = \sum_{j=1}^K \frac{e^{-d_j}}{\sum_{i=1}^K e^{-d_i}} = 1$$

όπου  $d_k$  συμβολίζει την απόσταση από το  $k$ -οστό κοντινότερο δεδομένο εκπαίδευσης

- Εναλλακτικά, μπορούμε να πάρουμε το άθροισμα όλων των σειρών από την εξαγόμενη μεταβλητή (κάθε σειρά αναπαριστά ένα δεδομένο εκπαίδευσης). Έχοντας υπολογίσει το άθροισμα των  $K$  κοντινότερων δεδομένων εκπαίδευσης σε κάθε pixel, μπορούμε να διαιρέσουμε το διάνυσμα των αθροισμάτων με το  $\mathbf{K}$ . Το αποτέλεσμα θα είναι μια μέση πρόβλεψη στη τιμή του κάθε pixel του τρέχοντος δεδομένου ελέγχου. Αυτή η μέθοδος θα ήταν ισοδύναμη με την προηγούμενη, αν είχαμε αναθέσει σε όλα τα βάρη ίσες τιμές:  $w_k = \frac{1}{K}$ , για κάθε  $k$ .

7. Αναθέτουμε τις τιμές του διανύσματος με τις προβλεπόμενες τιμές, μόνο στα αντίστοιχα σημεία (indices) με ελλειπείς τιμές για το τρέχον δεδομένο.
8. Επαναλαμβάνουμε τα βήματα 3-7 για όλα τα δεδομένα ελέγχου (test data).
9. Υπολογίζουμε τη ρίζα της μέσης τιμής του τετραγώνου του σφάλματος (RMSE) μεταξύ των προβλεπόμενων δεδομένων ελέγχου και των πραγματικών δεδομένων ελέγχου, έτσι ώστε να εκτιμήσουμε την απόδοση του αλγορίθμου, με τον ακόλουθο τρόπο:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^N \sum_{j=1}^D (X_{i,j} - \tilde{X}_{i,j})^2}$$

όπου  $X$  είναι τα αρχικά δεδομένα και  $\tilde{X}$  είναι τα δεδομένα εκπαίδευσης με τις προβλεπόμενες ελλειπείς τιμές. Όσο πιο κοντά είναι αυτός ο δείκτης στο 0, τόσο πιο αποδοτικός είναι ο αλγόριθμος.

Η συνολική διαδικασία της πρόβλεψης των pixel μπορεί να αναπαρασταθεί με χρήση πινάκων ως εξής:

$$\text{closest\_data}_{ij} = \begin{bmatrix} p_{11} & p_{12} & p_{13} & \dots & p_{1D} \\ p_{21} & p_{22} & p_{23} & \dots & p_{2D} \\ & & & & \\ p_{k1} & p_{k2} & p_{k3} & \dots & p_{kD} \end{bmatrix} \times \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_K \end{bmatrix} = \begin{bmatrix} p_{11} \cdot w_1 & p_{12} \cdot w_1 & \dots & p_{1D} \cdot w_1 \\ p_{21} \cdot w_2 & p_{22} \cdot w_2 & \dots & p_{2D} \cdot w_2 \\ & & & & \\ p_{K1} \cdot w_K & p_{K2} \cdot w_K & \dots & p_{KD} \cdot w_K \end{bmatrix}$$

όπου  $\mathbf{p}$  σημαίνει pixel,  $\mathbf{K}$  είναι ο αριθμός των κοντινότερων δεδομένων εκπαίδευσης που θα λάβουμε υπόψη και  $D$  είναι ο αριθμός των pixel. Επομένως, το διάνυσμα των προβλεπόμενων pixel για το

τρέχον δεδομένο ελέγχου θα είναι το ακόλουθο:

$$predicted\_pixels = \sum_{i=1}^K closest\_data_{ij}$$

Από αυτό το διάνυσμα θα αναθέσουμε τιμές μόνο στα ελλιπή pixel του τρέχοντος δεδομένου ελέγχου. Συνοψίζοντας, ο K-NN αλγόριθμος χρησιμοποιείται κυρίως για κατηγοριοποίηση. Ωστόσο, εδώ δείξαμε πως μπορούμε να τροποποιήσουμε τον αλγόριθμο για να εφαρμόσουμε regression, όπως για παράδειγμα αν θέλουμε να κάνουμε πρόβλεψη πραγματικών τιμών για pixel.

## 5.2 Ένας Προτεινόμενος Αλγόριθμος Συμπλήρωσης Ελλιπών Τιμών Χρησιμοποιώντας Variational Autoencoder

Μπορούμε να χρησιμοποιήσουμε variational autoencoders για να προβλέψουμε ελλιπείς τιμές σε εικόνες εκπαίδευσης. Η κύρια ιδέα είναι να τροποποιήσουμε τον αλγόριθμο VAE έτσι ώστε να κρατάει ακέραια τα δεδομένα εκεί όπου οι τιμές δεν είναι ελλιπείς και να αλλάζει μόνο τα σημεία με ελλιπείς τιμές, σε κάθε επανάληψη.

Πρώτον, κατασκευάζουμε το dataset με ελλιπείς τιμές,  $X\_train\_missing$ , από το αρχικό dataset  $X\_train$ . Αφήνουμε το  $\frac{1}{5}$  του dataset όπως ήταν. Από τα υπόλοιπα  $\frac{4}{5}$  των δεδομένων, αντικαθιστούμε το μισό κομμάτι κάθε εικόνας, με την "ελλιπή τιμή" (missing value). Η "ελλιπής τιμή" θα μπορούσε να οριστεί ίση με το μέσο του διαστήματος  $[lowest\_value, highest\_value]$ , όπου  $lowest\_value$  και  $highest\_value$  είναι η χαμηλότερη και η υψηλότερη τιμή που εμφανίζεται στο dataset, αντίστοιχα. Για παράδειγμα, αυτό το διάστημα, στο MNIST dataset είναι  $[0, 1]$ , άρα η "ελλιπής τιμή" είναι ίση με 0.5. Το κομμάτι που επιλέγουμε να διαγράψουμε και να αντικαταστήσουμε με ελλιπείς τιμές μπορεί να είναι το επάνω, κάτω, αριστερά ή δεξιά μισό της εικόνας. Μια εναλλακτική επιλογή είναι να διαλέξουμε τυχαία pixel, τα μισά από κάθε εικόνα και να τα αντικαταστήσουμε με την ελλιπή τιμή. Αφού κατασκευάσουμε το dataset με τις ελλιπείς τιμές,  $X\_train\_missing$ , χρειαζόμαστε να κατασκευάσουμε ένα πίνακα με δυαδικές τιμές (0 ή 1), ο οποίος θα αποθηκεύει την πληροφορία σχετικά με το ποια σημεία των εικόνων αντικαταστήσαμε με ελλιπείς τιμές. Ονομάζουμε αυτόν τον πίνακα  $X\_train\_masked$ . Αν κάποιο pixel δεν αντικαταστάθηκε, το αντίστοιχο pixel στον πίνακα  $X\_train\_masked$  θα είναι 1. Αντίθετα, αν ένα pixel σε μια εικόνα αντικαταστάθηκε από την ελλιπή τιμή, το αντίστοιχο pixel στον πίνακα  $X\_train\_masked$  θα είναι 0. Επιπλέον, πρέπει να καθορίσουμε τον πίνακα που θα περιλαμβάνει τις προβλεπόμενες τιμές των ελλιπών pixel. Ονομάζουμε αυτόν τον πίνακα  $X\_filled$ . Αρχικοποιούμε τον πίνακα  $X\_filled$  να είναι ίδιος με τον πίνακα  $X\_train\_missing$ .

Τελικά, τρέχουμε μια τροποποιημένη έκδοση του VAE αλγορίθμου. Όπως αναφέραμε νωρίτερα, σε κάθε επανάληψη, πρέπει να αντικαθιστούμε τα pixel μόνο εκεί όπου υπήρχαν ελλιπείς τιμές, με τις

προβλεπόμενες. Την ίδια στιγμή, πρέπει να διατηρούμε ακέραια τα pixel με τις μη ελλιπείς τιμές. Για αυτό το σκοπό, θα φανεί χρήσιμος ο πίνακας  $X\_train\_masked$ .

Ακολουθεί ψευδοκώδικας για όλη τη διαδικασία που περιγράψαμε παραπάνω:

```

for epoch = 0 to epochs - 1
    iterations = N / batch_size

    for i = 0 to iterations - 1
        start_index = i * batch_size
        end_index = (i + 1) * batch_size

        // fetch the batch data, labels and masked batch data
        batch_data = X_filled(start_index:end_index, :)
        batch_labels = y_train(start_index:end_index)
        masked_batch_data = X_train_masked(start_index:end_index, :)

        // train the batch data using the VAE process
        cur_samples = train(batch_data, params)

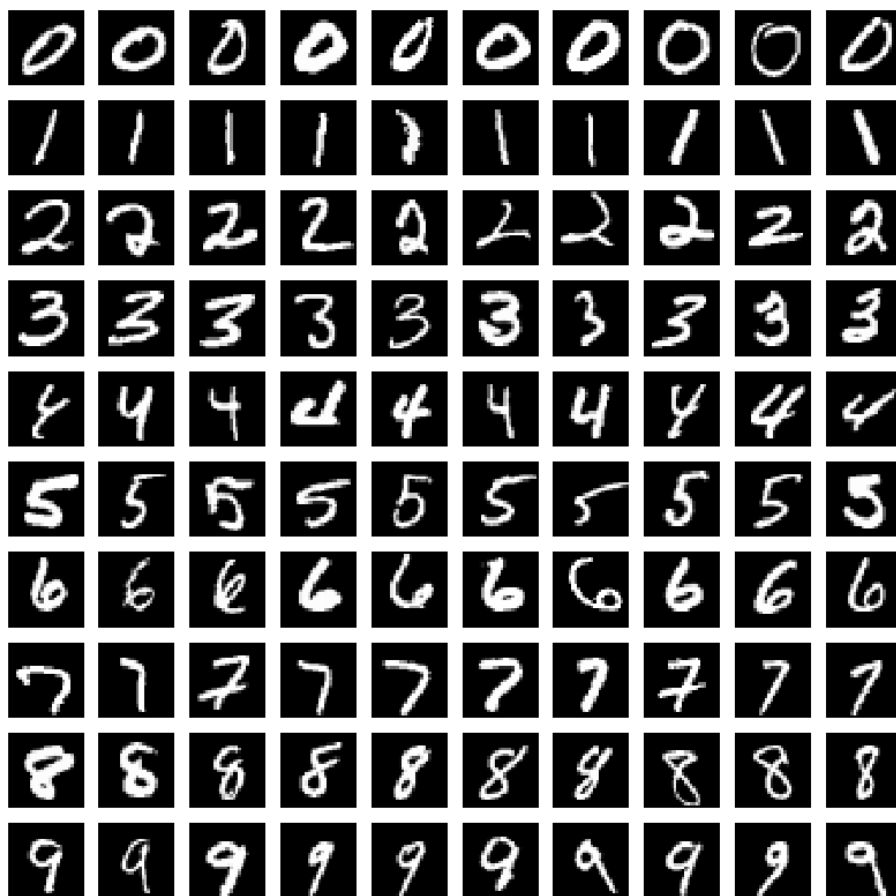
        // "." denotes element-wise multiplication
        // The "cur_samples" will take values from the "batch_data"
        // where the pixels are observed (with masked values=1)
        // and will keep intact its values from the VAE training
        // where the pixels are missing (with masked values=0).
        cur_samples = masked_batch_data .* batch_data +
                      (1 - masked_batch_data) .* cur_samples
        X_filled(start_index:end_index, :) = cur_samples

```

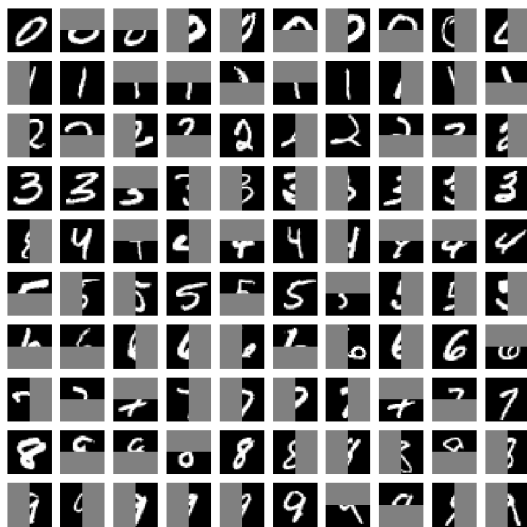
ΚΩΔΙΚΑΣ 5.3: ψευδοκώδικας VAE αλγορίθμου συμπλήρωσης ελλιπών τιμών

### 5.3 Αποτελέσματα Πειραμάτων του K-NN Αλγορίθμου Συμπλήρωσης Ελλιπών Τιμών και Αλγορίθμων Συμπλήρωσης Ελλιπών Τιμών με Χρήση VAE

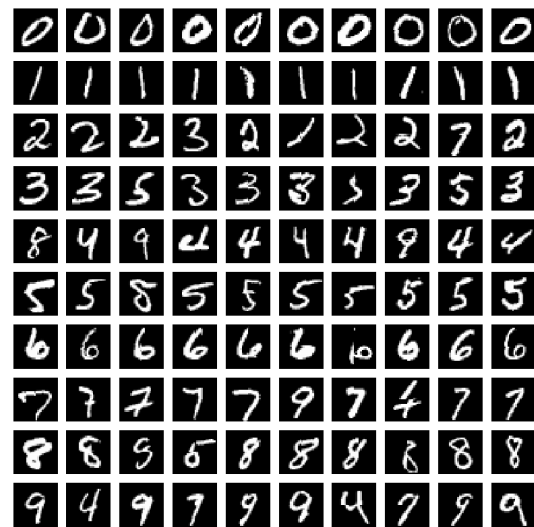
Μετά από εκτελέσεις του αλγορίθμου συνεργατικού φιλτραρίσματος στο MNIST dataset με ελλιπείς τιμές, για διάφορες τιμές του  $K$ , και εκτελέσεις του VAE αλγορίθμου αλλιπών τιμών σε PyTorch, καταλήξαμε στα παρακάτω αποτελέσματα. Όπως πρόκειται να δούμε, κάποιες εικόνες ελέγχου (test) των ψηφίων 3, 5, 8 ανακατασκευάστηκαν με το άλλο μισό των ψηφίων 3 ή 5 ή 8. Αυτό το ατυχές συμβάν είναι λογικό, γιατί τα ψηφία 3, 5 και 8 μοιάζουν μεταξύ τους, εν αντιθέσει με τα άλλα ψηφία.



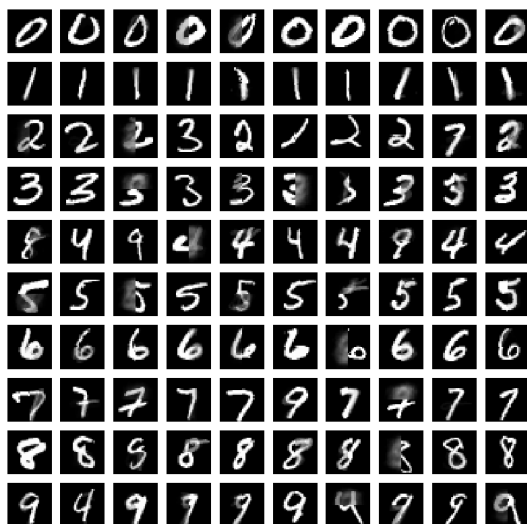
ΕΙΚΟΝΑ 5.1: Αρχικά MNIST Δεδομένα Ελέγχου (Test)



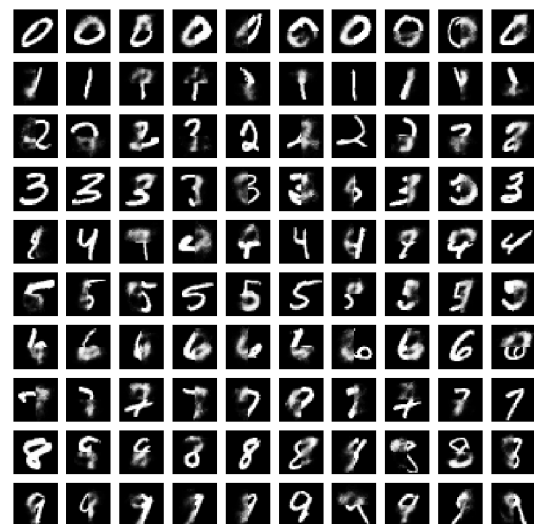
(a) Δεδομένα Εκπαίδευσης με Δομημένες (Structured) Ελλιπίες Τιμές



(b) 1-NN



(c) 100-NN



(d) VAE σε PyTorch Εποχή 200

ΕΙΚΟΝΑ 5.2: MNIST 1-NN, 100-NN & VAE αλγόριθμοι ελλιπών τιμών.

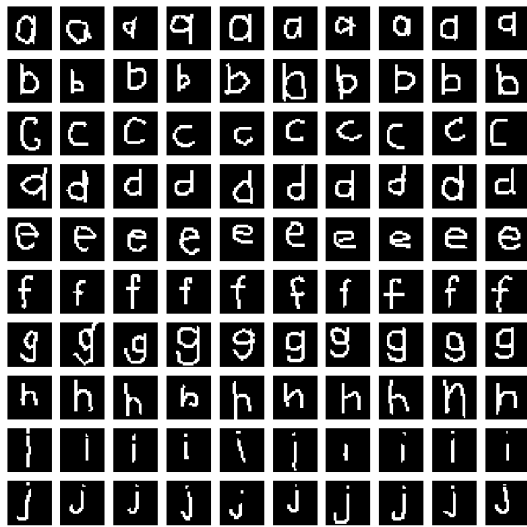
Το ελάχιστο root mean squared error (RMSE) και το ελάχιστο mean absolute error (MAE) που επιτεύχθηκαν ήταν για  $K=100$ , αλλά ο αλγόριθμος ήταν ο πιο χρονοβόρος συγκριτικά με τους άλλους.

Method	RMSE	MAE	Time
1-NN	0.171569	0.0406822	190.26 sec
100-NN	0.148223	0.0396628	233.19 sec
VAE	0.168031	0.0499518	178.09 sec

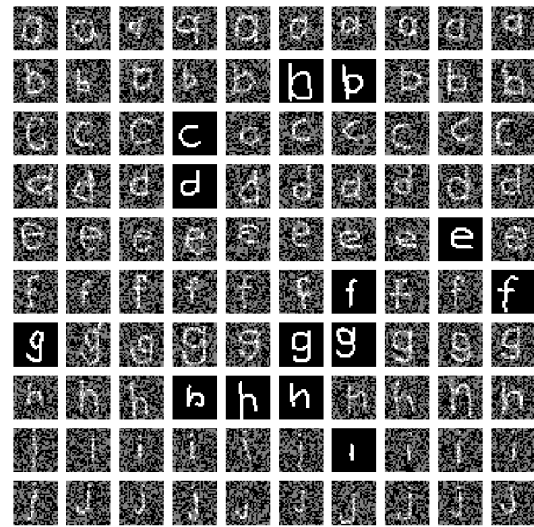
ΠΙΝΑΚΑΣ 5.1: Αλγόριθμοι συμπλήρωσης ελλιπών τιμών στο MNIST dataset.



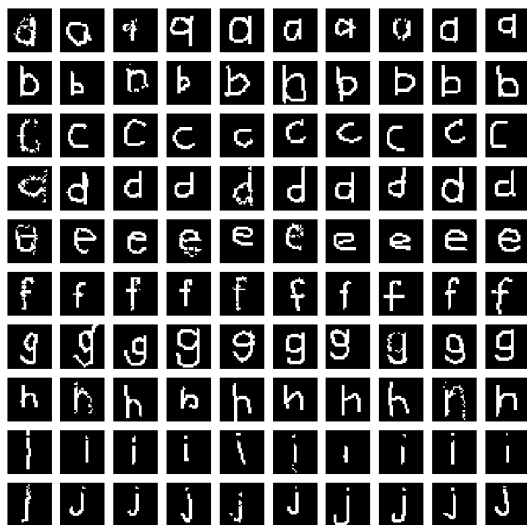
Ακολουθούν τα αποτελέσματα των αλγορίθμων συμπλήρωσης ελλιπών τιμών, στο OMNIGLOT Αγγλικό dataset: OMNIGLOT English dataset.



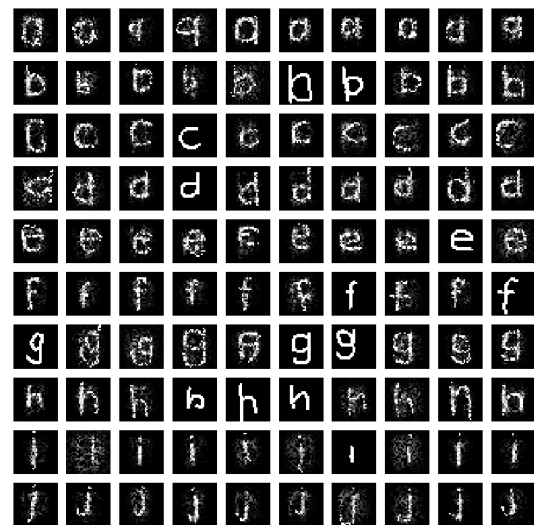
(a) Αρχικά Δεδομένα



(b) Ελλιπή Δεδομένα



(c) 1-NN



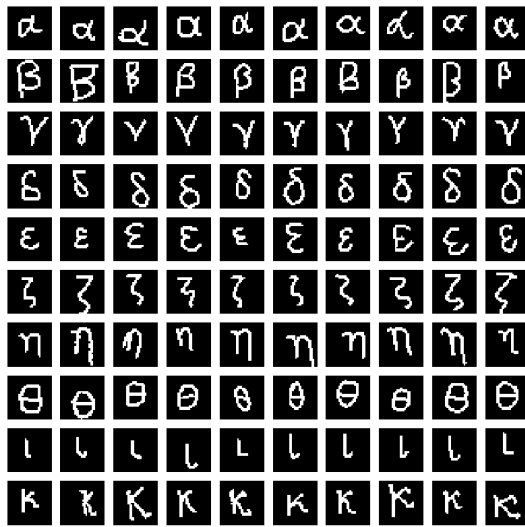
(d) VAE σε PyTorch Εποχή 100

ΕΙΚΟΝΑ 5.3: 1-NN & VAE αλγόριθμοι ελλιπών τιμών στο OMNIGLOT Αγγλικό αλφάβητο, χαρακτήρες 1-10, αρχικές, τυχαίες ελλιπείς τιμές και ανακατασκευασμένες.

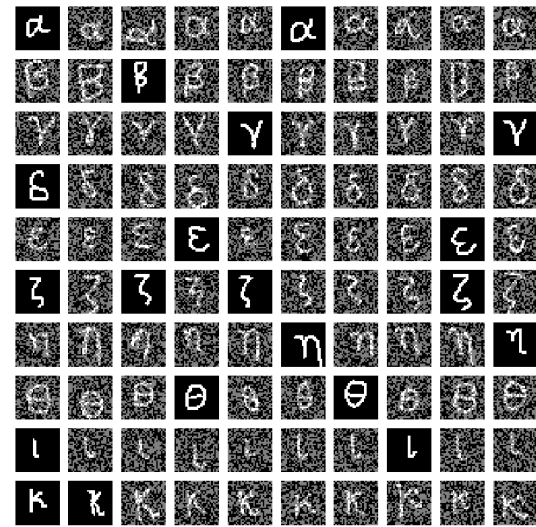
Method	RMSE	MAE	Time
1-NN	0.0890025270986	0.00792144982993	8.64 sec
VAE	0.168031	0.0499518	178.09 sec

ΠΙΝΑΚΑΣ 5.2: Αλγόριθμος συμπλήρωσης ελλιπών τιμών στο OMNIGLOT English dataset.

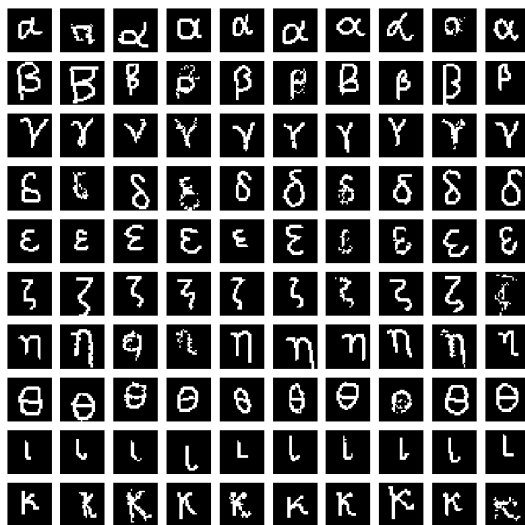
Ακολουθούν τα αποτελέσματα των αλγορίθμων συμπλήρωσης ελλιπών τιμών, με χρήση **PyTorch**, στο **OMNIGLOT** Ελληνικό dataset.



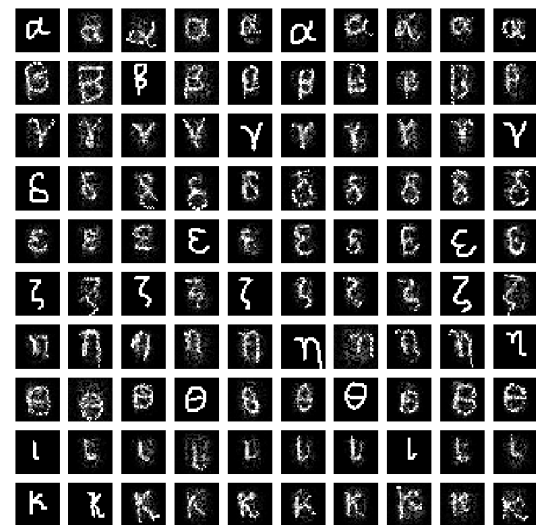
(a) Αρχικά Δεδομένα



(b) Ελλιπή Δεδομένα



(c) 1-NN



(d) VAE σε PyTorch Εποχή 100

ΕΙΚΟΝΑ 5.4: 1-NN & VAE αλγόριθμοι ελλιπών τιμών στο OMNIGLOT Ελληνικό αλφάβητο, χαρακτήρες 1-10, αρχικές, τυχαίες ελλιπείς τιμές και ανακατασκευασμένες.

Method	RMSE	MAE	Time
1-NN	0.0846893853663	0.00717229199372	9.37 sec
VAE	0.162737553068	0.0524327812139	245.66 sec

ΠΙΝΑΚΑΣ 5.3: Αλγόριθμος συμπλήρωσης ελλιπών τιμών στο OMNIGLOT Ελληνικό dataset.

**ΣΗΜΕΙΩΣΗ:** Ο δείκτης του χρόνου (time) για τους K-NN αλγόριθμους δείχνει τη διάρκεια των υπολογισμών των αποστάσεων και των προβλέψεων, ενώ ο δείκτης του χρόνου για τους VAE αλγόριθμους δείχνει τις διάρκειες των βρόχων εκπαίδευσης των VAE.

## 5.4 K-NN vs VAE Αλγόριθμος Συμπλήρωσης Ελλιπών Τιμών στο MovieLens Dataset

Τρέξαμε τους K-NN και VAE σε TensorFlow αλγόριθμους συμπλήρωσης ελλιπών τιμών στο MovieLens 100k dataset, `ua set` και συγκρίναμε τις προβλεπόμενες βαθμολογίες. Τα αποτελέσματα φαίνονται παρακάτω:

Metric	Value
root mean squared error (RMSE)	0.0803184360998
mean absolute error (MAE)	0.0209178842034
match percentage	91.2514516501 %
Method	Mean Rating
VAE in TensorFlow	1.60717332671
K-NN Missing Values algorithm	1.61095209334

ΠΙΝΑΚΑΣ 5.4: Σύγκριση μεταξύ των 10-NN και VAE σε TensorFlow αλγόριθμων συμπλήρωσης ελλιπών τιμών.

Από τον πίνακα αυτόν και κρίνοντας από τα σφάλματα, συμπεραίνουμε ότι οι προβλέψεις των δύο αλγόριθμων για τις βαθμολογίες των χρηστών είναι πολύ κοντινές. Μάλιστα, όταν κάναμε στρογγυλοποίηση εξαλείφοντας τα δεκαδικά ψηφία των βαθμολογιών, καταλήξαμε ότι τα αποτελέσματα ταυτίζονται κατά περίπου 91%. Τέλος, η μέση βαθμολογία για τον VAE αλγόριθμο σε TensorFlow είναι περίπου 1.61, ενώ η μέση βαθμολογία για τον K-NN αλγόριθμο συμπλήρωσης ελλιπών τιμών είναι περίπου 1.6. Οι βαθμολογίες είναι στρογγυλοποιημένες ώστε να παίρνουν ακέραιες τιμές στο διάστημα  $[1, 5]$  και η τιμή που δείχνει ότι ένας χρήστης δεν έχει βαθμολογήσει μια ταινία είναι 0.

## Κεφάλαιο 6

# Γραφική Διεπαφή (GUI) για Variational Autoencoder & Αλγορίθμους Συμπλήρωσης Ελλιπών Τιμών

Μαζί με το project αυτής της διατριβής έχει υλοποιηθεί μια εφαρμογή γραφικής διεπαφής (GUI), σε προγραμματιστική γλώσσα Python 3, κάνοντας χρήση της βιβλιοθήκης **Tkinter**.

Πρώτα πηγαίνετε στο directory "**vaes\_gui**" της εργασίας και εγκαταστήστε τις απαιτούμενες Python βιβλιοθήκες, πληκτρολογώντας:

```
pip install -r dependencies.txt
```

ΚΩΔΙΚΑΣ 6.1: εντολή για εγκατάσταση Python εξαρτήσεων

Για να εκτελέσετε το GUI από το τερματικό, πληκτρολογείτε την εντολή:

```
python vaes_gui.py
```

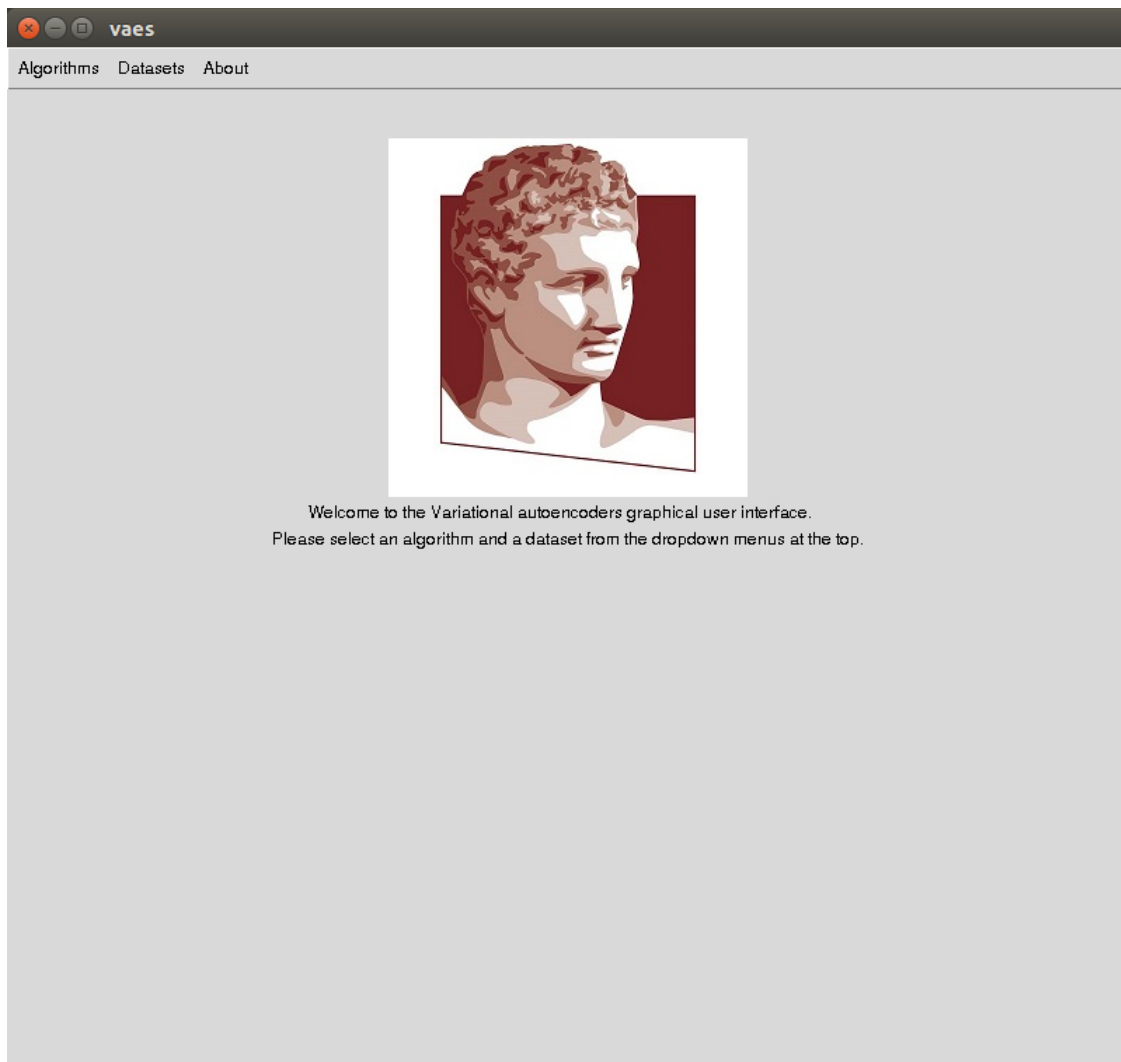
ΚΩΔΙΚΑΣ 6.2: εντολή για εκτέλεση του GUI

Για να δημιουργήσετε ένα εκτελέσιμο αρχείο για το GUI (".exe"), το οποίο θα μπορείτε να ανοίξετε ανά πάσα ώρα σε Windows περιβάλλον, πληκτρολογείτε:

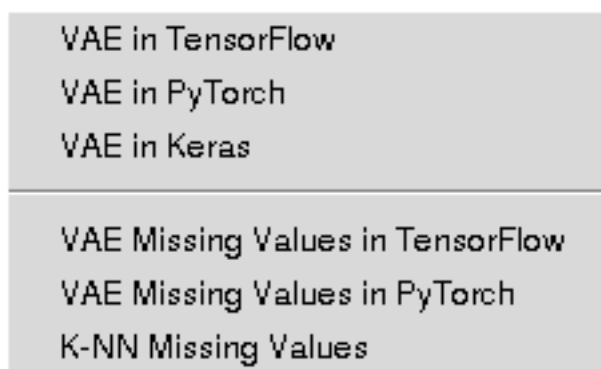
```
pip install pyinstaller  
pyinstaller vaes_gui.spec
```

ΚΩΔΙΚΑΣ 6.3: εντολή για δημιουργία εκτελέσιμου αρχείου για το GUI

Στη συνέχεια, κατεβάστε όλα τα dataset από τα URL στο αρχείο "**datasets\_urls.txt**" και μετακινήστε τα στον νέο φάκελο "**dist**". Εντός του φακέλου, υπάρχει ένας άλλος φάκελος με όνομα "**vaes\_gui**", ο οποίος περιλαμβάνει το εκτελέσιμο αρχείο "**vaes\_gui.exe**".



ΕΙΚΟΝΑ 6.1:  
Σελίδα καλωσορίσματος του GUI.

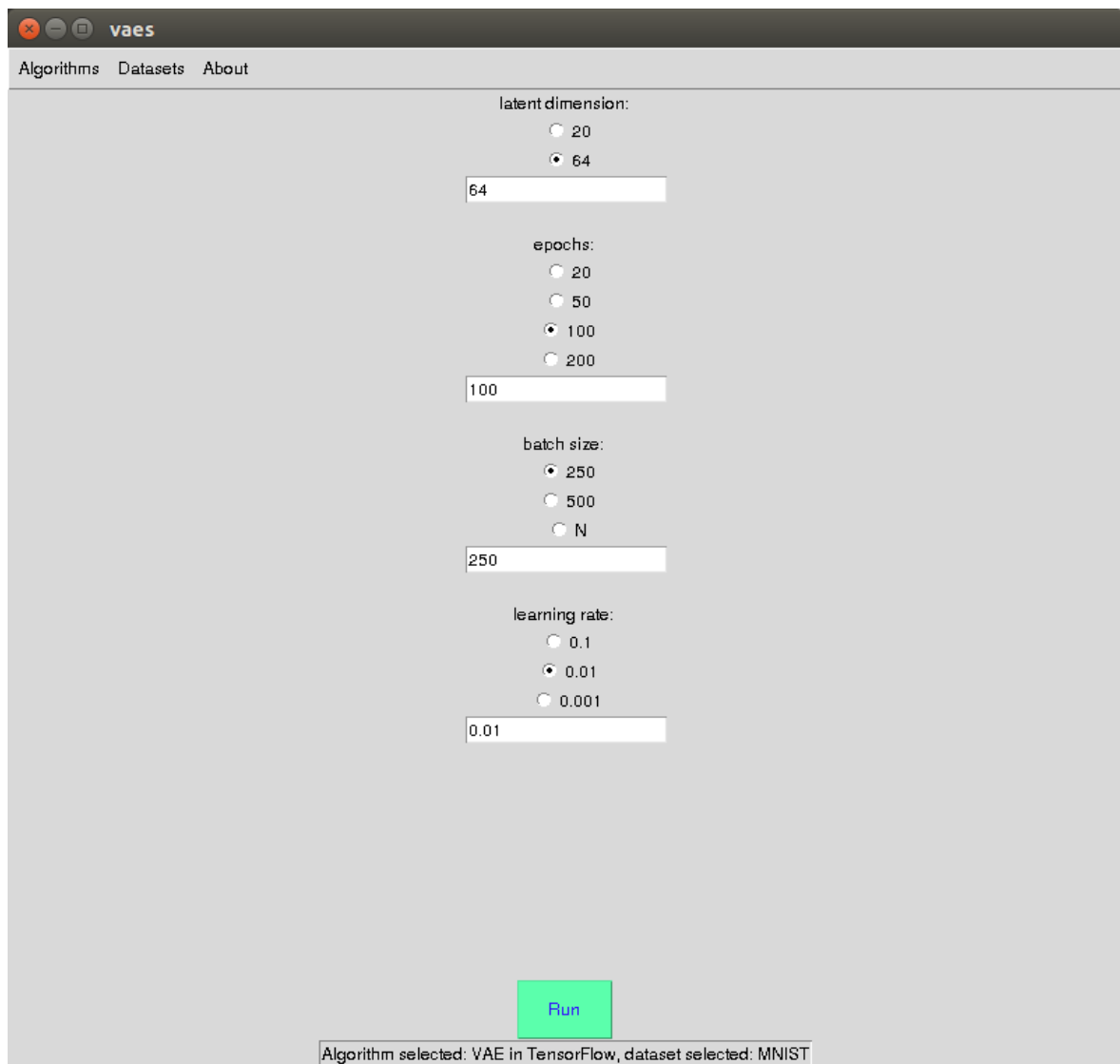


(a) GUI dropdown menu Αλγορίθμων.

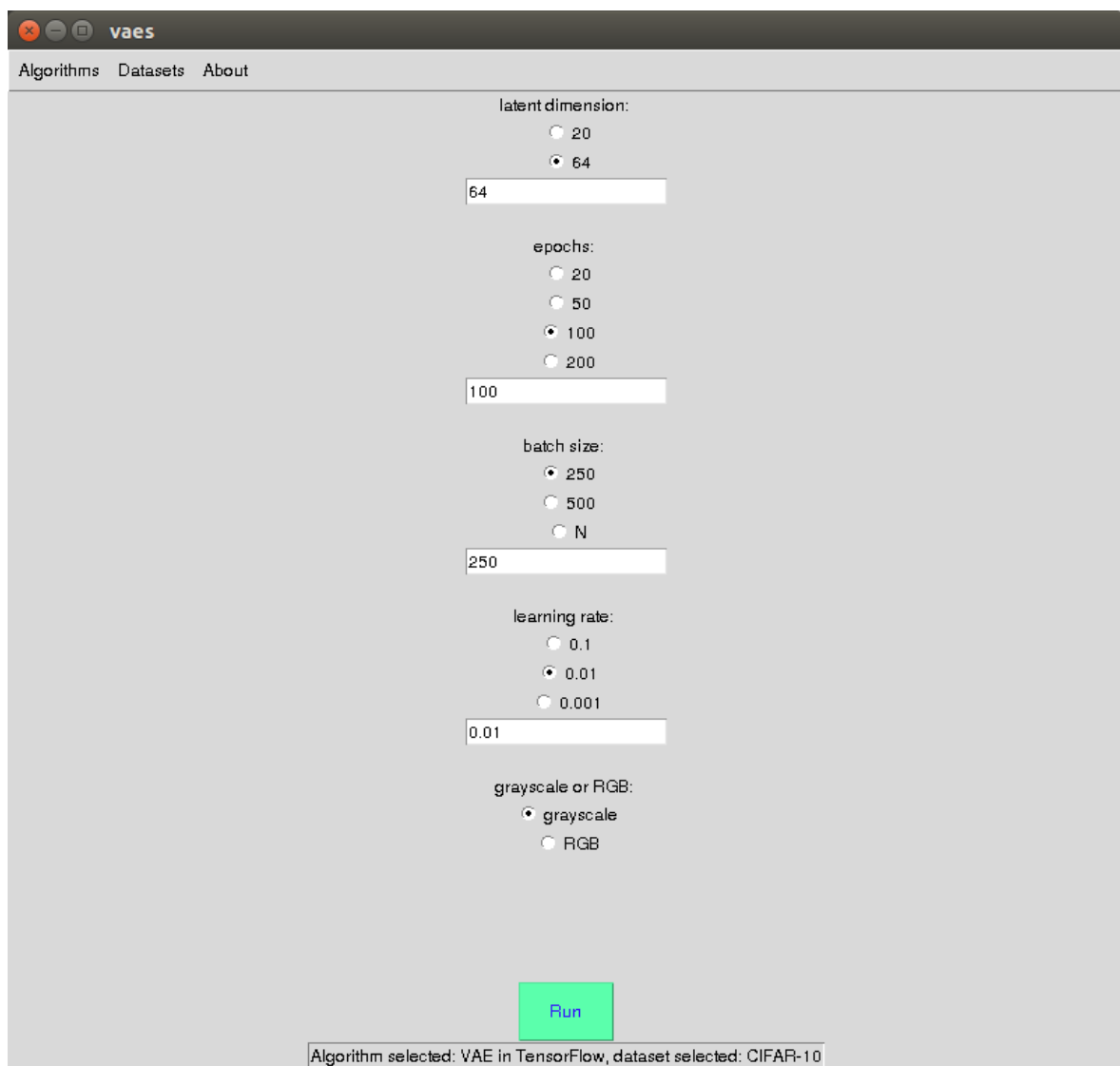


(b) GUI Datasets dropdown menu.

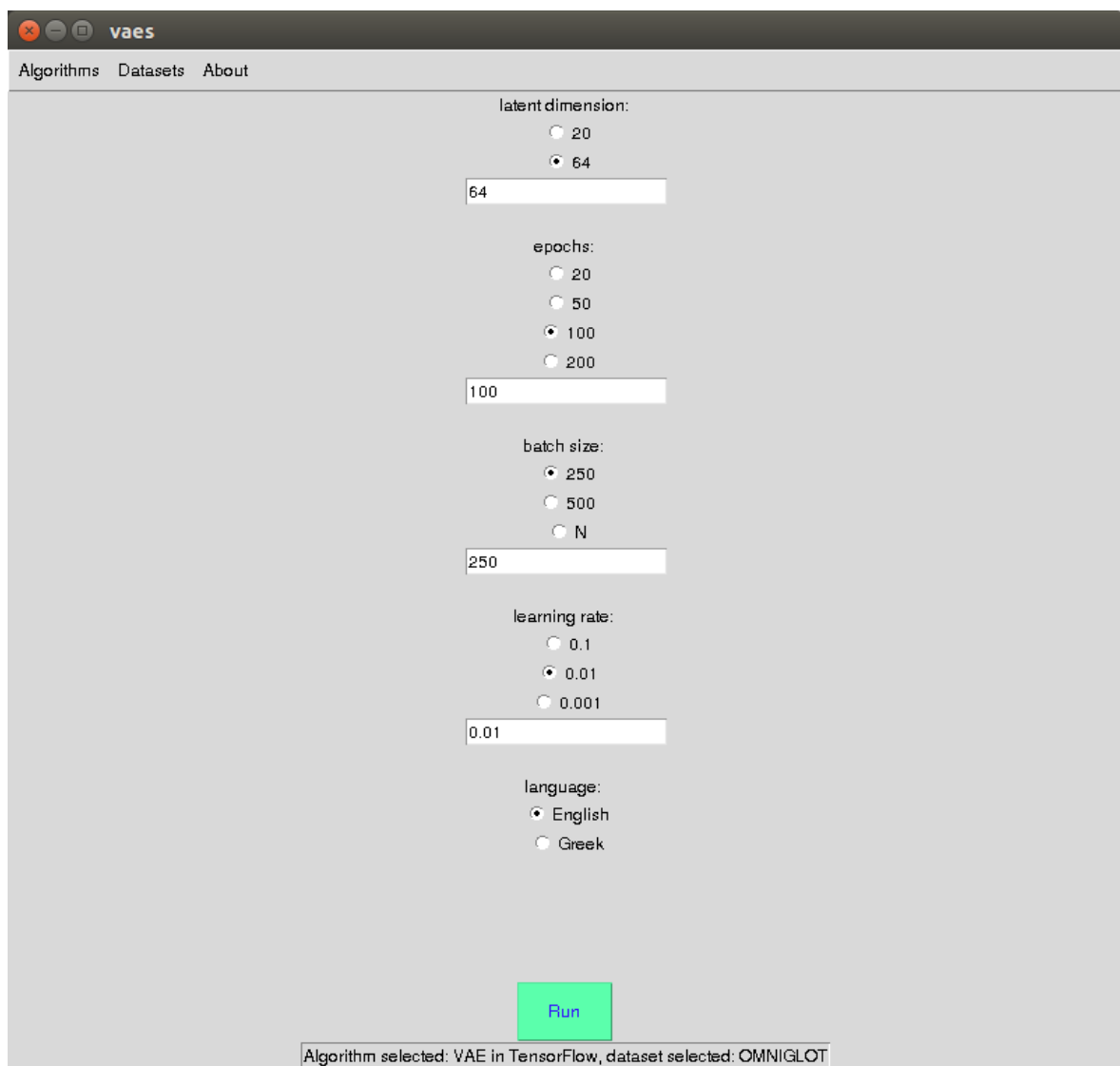
ΕΙΚΟΝΑ 6.2: Dropdown menus.



ΕΙΚΟΝΑ 6.3: GUI VAE σε TensorFlow, MNIST dataset.

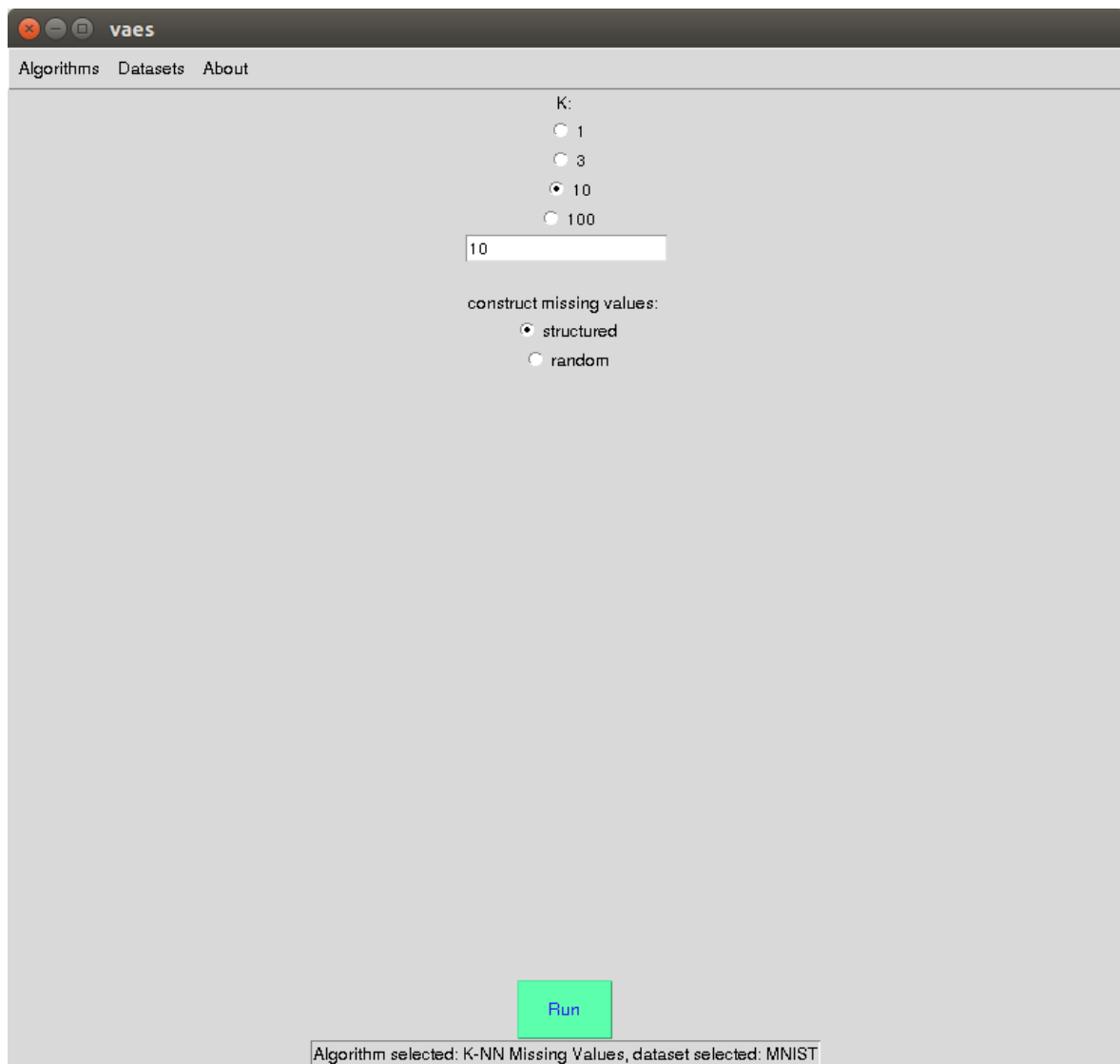


ΕΙΚΟΝΑ 6.4: GUI VAE σε TensorFlow, CIFAR-10 dataset.

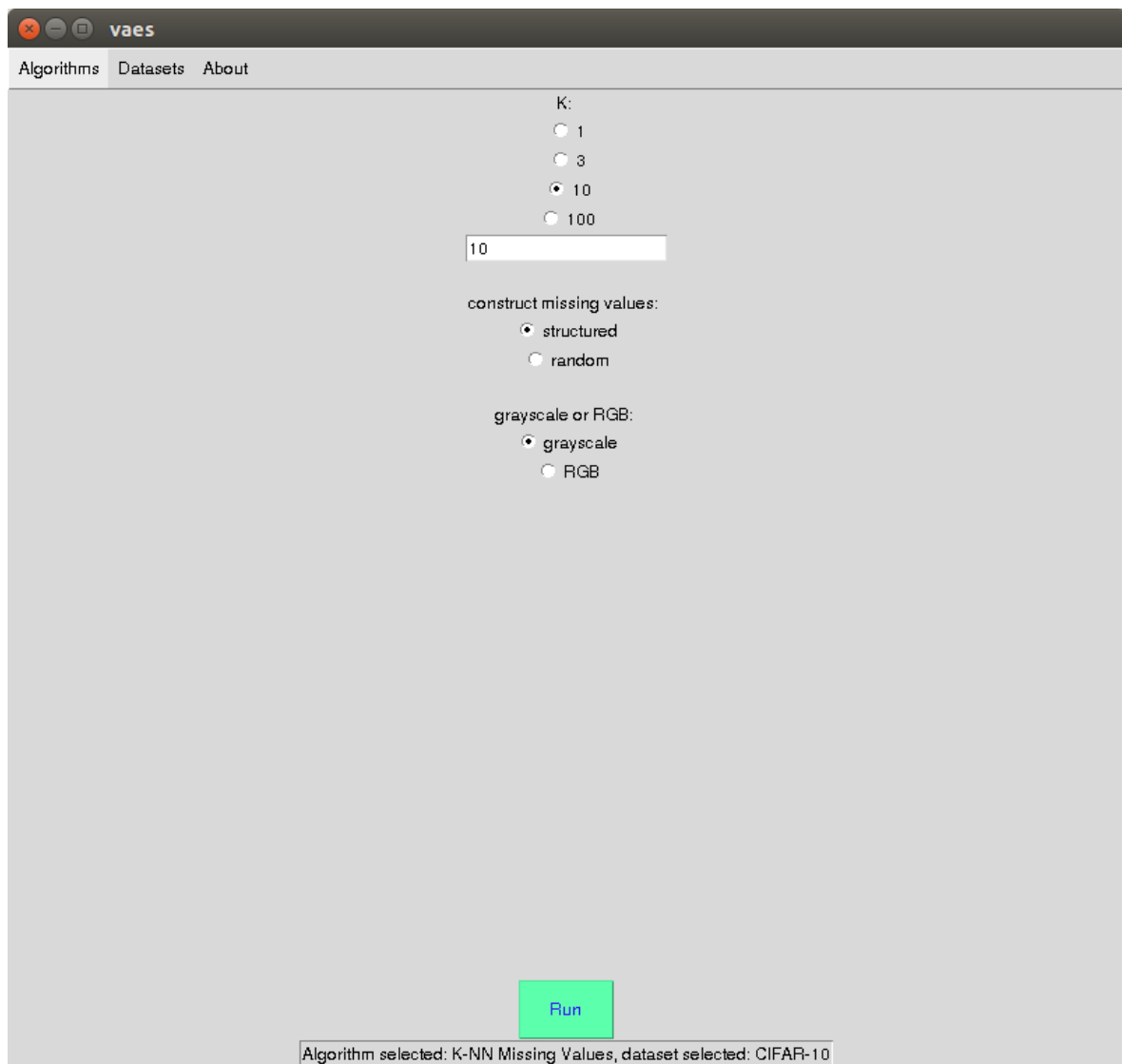


ΕΙΚΟΝΑ 6.5: GUI VAE σε TensorFlow, OMNIGLOT dataset.

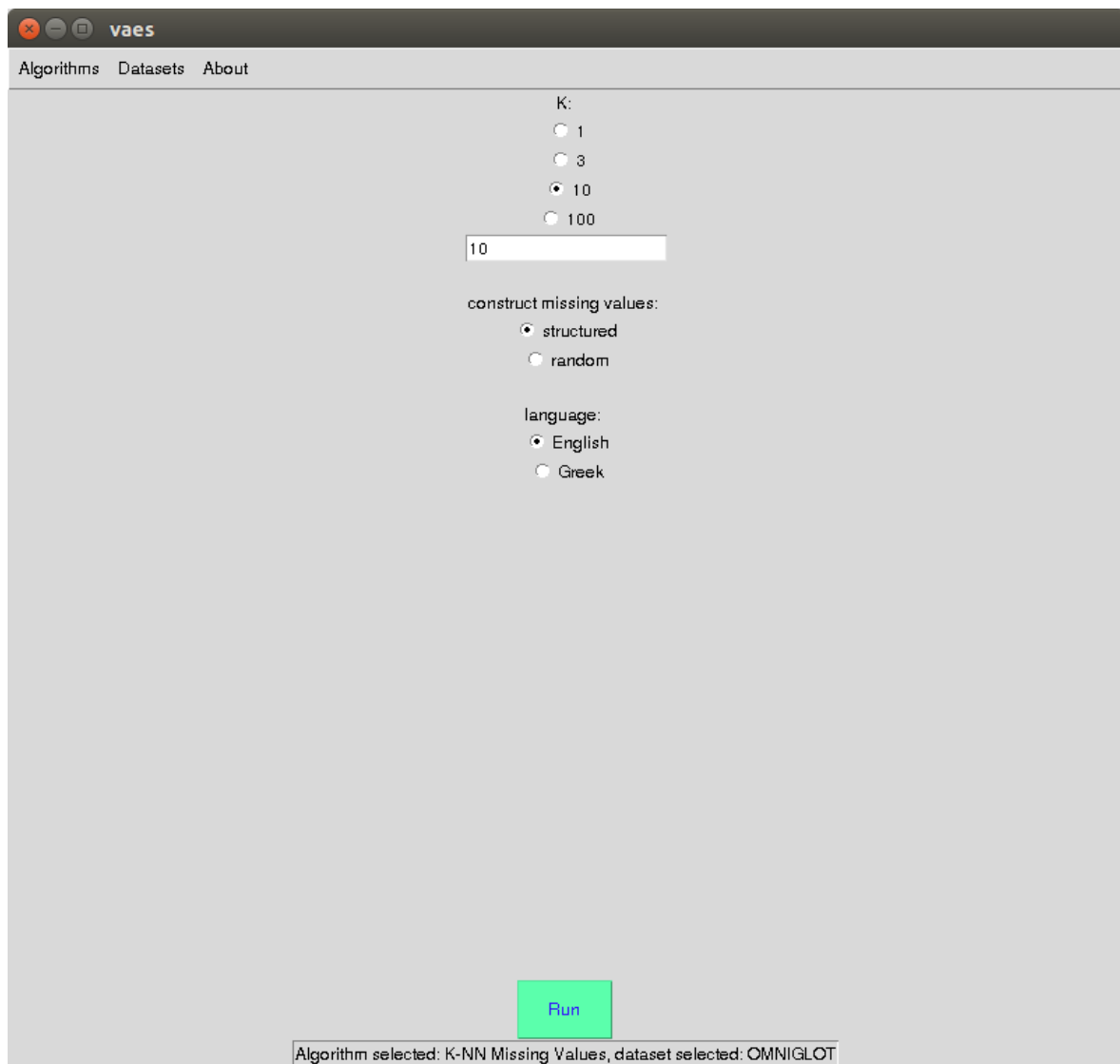




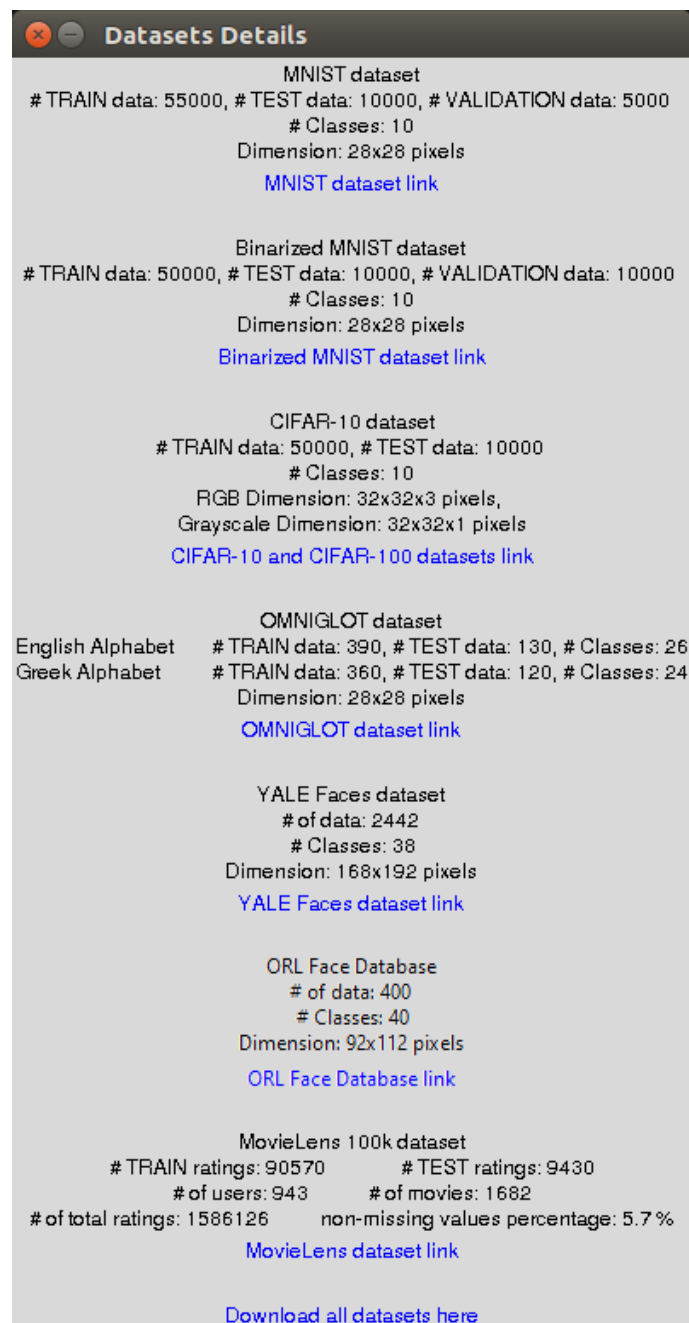
ΕΙΚΟΝΑ 6.6: GUI K-NN αλγόριθμος Ελλιπών Τιμών, MNIST dataset.



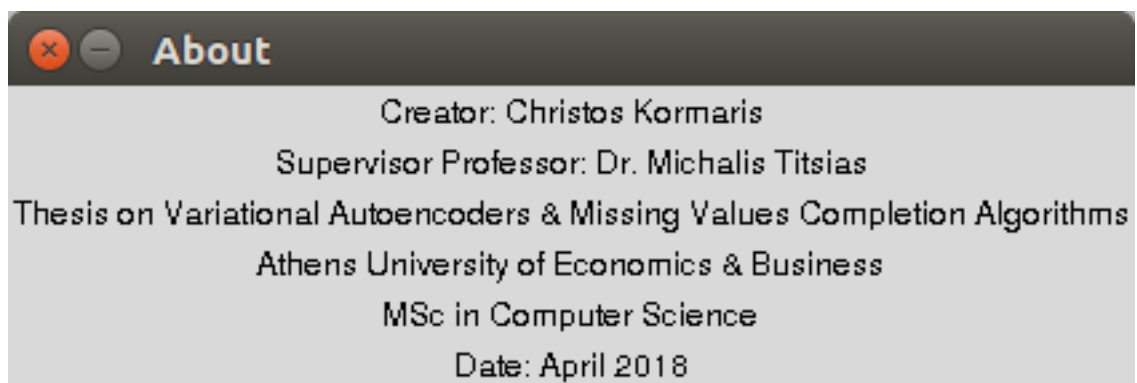
ΕΙΚΟΝΑ 6.7: GUI K-NN αλγόριθμος Ελλιπών Τιμών, CIFAR-10 dataset.



ΕΙΚΟΝΑ 6.8: GUI K-NN αλγόριθμος Ελλιπών Τιμών, OMNIGLOT dataset.



ΕΙΚΟΝΑ 6.9: GUI λεπτομέρειες των dataset.



ΕΙΚΟΝΑ 6.10: GUI σχετικά.

## Κεφάλαιο 7

# Περαιτέρω Εφαρμογές των Variational Autoencoder

(Aaron Courville, Ian Goodfellow, Yoshua Bengio 2016. [4])

Οι autoencoder έχουν εφαρμοστεί επιτυχώς σε διεργασίες όπως: **1) ελάττωση διάστασης** και **2) ανάκτηση πληροφοριών**. Η ελάττωση διάστασης ήταν από μια από τις πρώτες εφαρμογές της αντιπροσωπευτικής μάθησης (representation learning) και της βαθιάς μάθησης (deep learning).

**Αναπαραστάσεις σε χαμηλότερη διάσταση** μπορούν να βελτιώσουν την απόδοση σε πολλές διεργασίες, όπως η κατηγοριοποίηση. Μοντέλα μικρότερου όγκου καταναλώνουν λιγότερη μνήμη και χρόνο εκτέλεσης. Τα στοιχεία που παρέχονται από τη χαρτογράφηση στο χώρο χαμηλότερης διάστασης βοηθούν στη γενίκευση (generalization).

Μια διαδικασία που ωφελεί παραπάνω από το συνηθισμένο στην ελάττωση διάστασης είναι η **ανάκτηση πληροφοριών**. Ανάκτηση πληροφοριών είναι η διαδικασία αναζήτησης καταχωρήσεων, από μια βάση δεδομένων που μοιάζουν σε μια καταχώρηση-ερώτημα (query entry). Αυτή η διαδικασία αντλεί τα ίδια πλεονεκτήματα με άλλες διεργασίες, από την ελάττωση διάστασης, και έχει επιπλέον το πλεονέκτημα ότι η αναζήτηση μπορεί να γίνει υπερβολικά επαρκής σε συγκεκριμένα είδη χώρων χαμηλής διάστασης. Ειδικά, αν εκπαιδεύσουμε τον αλγόριθμο ελάττωσης διάστασης ώστε να παράγει έναν κωδικό που να είναι μικρής διάστασης και δυαδικός, τότε μπορούμε να αποθηκεύσουμε όλες τις καταχωρήσεις της βάσης δεδομένων σε ένα πίνακα κατακερματισμού (hash table), αντιστοιχίζοντας διανύσματα δυαδικού κωδικού σε καταχωρίσεις (entries). Αυτός ο πίνακας κατακερματισμού μας επιτρέπει να εκτελέσουμε ανάκτηση πληροφοριών επιστρέφοντας όλες τις καταχωρήσεις της βάσης δεδομένων που έχουν τον ίδιο δυαδικό κωδικό με το ερώτημα (query). Μπορούμε επίσης να ψάξουμε σε λίγο λιγότερο παρόμοιες καταχωρήσεις, πολύ αποτελεσματικά, αναστρέφοντας απλώς μεμονωμένα bits από την κωδικοποίηση του ερωτήματος. Αυτή η προσέγγιση στην ανάκτηση πληροφοριών μέσω της ελάττωσης διάστασης και της κωδικοποίησης σε δυαδικό σύστημα ονομάζεται **σημασιολογικός κατακερματισμός** (semantic hashing).

# Παράρτημα Α

## Θεωρητικό Υπόβαθρο

### A.1 Κανόνας του Bayes

Κανόνας του Bayes για υπό συνθήκη (conditional) πιθανότητες:

$$P(A|B) \cdot P(B) = P(B|A) \cdot P(A) \Rightarrow$$

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Επίσης, ισχύει η παρακάτω εξίσωση:

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

### A.2 Συνάρτηση Softmax

Έστω  $X$  ένα διάνυσμα ή ένας πίνακας, το δεδομένο εισόδου. Η συνάρτηση softmax κατασκευάζει βάρη για κάθε στοιχείο της εισόδου  $X$ . Το στοιχείο με τη μεγαλύτερη τιμή θα πάρει το βάρος με τη μεγαλύτερη τιμή. Επίσης, το άθροισμα όλων των βαρών πρέπει να είναι ίσο με 1. Επομένως, η συνάρτηση softmax δίνεται από τον παρακάτω τύπο:

$$w_i = \text{softmax}(X) = \frac{e^{X_i}}{\sum_{j=1}^N e^{X_j}},$$

$$\text{and } w_1 + w_2 + w_3 + \dots + w_N = \sum_{i=1}^N w_i = \sum_{i=1}^N \frac{e^{X_i}}{\sum_{j=1}^N e^{X_j}} = \frac{\sum_{i=1}^N e^{X_i}}{\sum_{j=1}^N e^{X_j}} = 1$$

### A.3 Εντροπία Θεωρίας Πληροφοριών (Entropy)

Η εντροπία πληροφορίας ορίζεται ως η μέση ποσότητα πληροφορίας που παράγεται από μια στοχαστική πηγή δεδομένων.

Ο τύπος της εντροπίας πληροφορίας μιας κατανομής  $P$  είναι ο ακόλουθος:

$$H(P) = - \sum_x P(x) \cdot \log_b P(x)$$

όπου  $b$  είναι η βάση του λογαρίθμου

Η εντροπία πληροφορία υπολογίζεται τυπικά σε bits (εναλλακτικά επωνομαζόμενα "**Shannons**") για  $b=2$  ή ορισμένες φορές σε "φυσικές μονάδες" (natural units) (**nats**) για  $b=e$  (σταθερά του Euler), ή σε δεκαδικά ψηφία (επωνομαζόμενα "dits", "bans", ή "hartleys") για  $b=10$ . Η μονάδα μέτρησης εξαρτάται από τη βάση του λογαρίθμου που χρησιμοποιείται για να καθορίσει την εντροπία.

### A.4 Cross-Entropy

Έστω  $Q$  μία "αφύσικη" πιθανοτική κατανομή και έστω  $P$  η "πραγματική" κατανομή για κάποιο σύνολο συμβάντων. Το cross entropy μεταξύ δύο πιθανοτικών κατανομών  $P$  και  $Q$  υπολογίζει τον μέσο αριθμό από bit που χρειάζονται για να αναγνωριστεί συμβάν επιλεγμένο από το σύνολο.

Ο τύπος του cross entropy δύο διακριτών κατανομών  $P$  και  $Q$  είναι ο ακόλουθος:

$$H(P, Q) = - \sum_x P(x) \cdot \log_b Q(x)$$

Παρομοίως, ο τύπος του cross entropy δύο συνεχών κατανομών  $P$  και  $Q$  είναι ο ακόλουθος:

$$H(P, Q) = - \int_{-\infty}^{+\infty} P(x) \cdot \log_b Q(x) dx \Rightarrow$$

$$H(P, Q) = -E_P[-\log Q(x)]$$

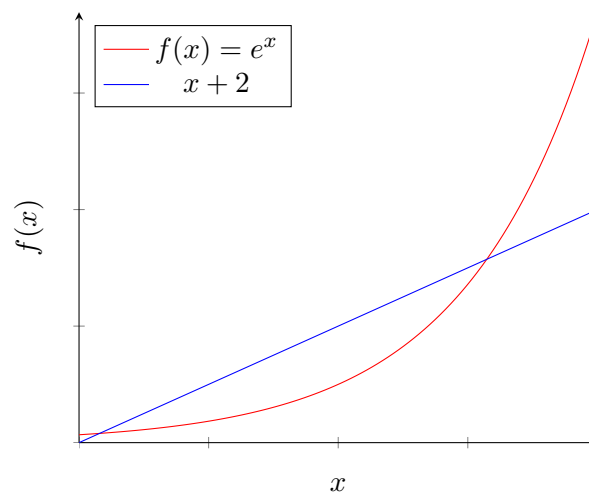
όπου  $b$  είναι η βάση του λογαρίθμου



## A.5 Ανισότητα του Jensen

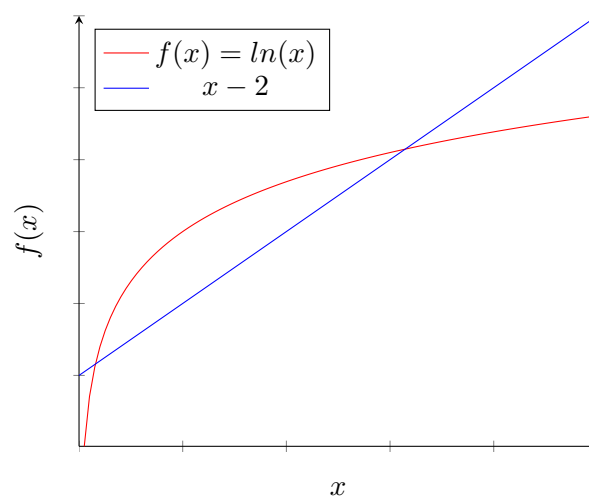
Η ανισότητα του Jensen, η οποία πήρε το όνομά της από το Δανό μαθηματικό Johan Jensen, συσχετίζει την τιμή μιας κυρτής (ή κοίλης) συνάρτησης ενός ολοκληρώματος με το ολοκλήρωμα μιας κυρτής συνάρτησης. Γενικεύει τη δήλωση ότι μια τέμνουσα γραμμή σε μια κυρτή συνάρτηση βρίσκεται πάνω από τη γραφική παράσταση της συνάρτησης.

$$f(E[X]) \leq E[f(X)], \text{ όπου } f \text{ είναι μια κυρτή συνάρτηση}$$



Με την ίδια λογική, μια τέμνουσα γραμμή σε μια κοίλη συνάρτηση βρίσκεται κάτω από τη γραφική παράσταση της συνάρτησης.

$$f(E[X]) \geq E[f(X)], \text{ όπου } f \text{ είναι μια κοίλη συνάρτηση}$$



## A.6 Kullback-Leibler (KL) Απόκλιση

Αν έχουμε δύο ξεχωριστές πιθανοτικές κατανομές  $P(x)$  και  $Q(x)$  για την ίδια Τ.Μ.  $x$ , μπορούμε να μετρήσουμε πόσο διαφορετικές είναι αυτές οι δύο κατανομές χρησιμοποιώντας την Kullback-Leibler (KL) απόκλιση (divergence):

$$D_{KL}[P \parallel Q] = E_{X \sim P}[\log \frac{P(X)}{Q(X)}] = E_{X \sim P}[\log P(X) - \log Q(X)]$$

Στην περίπτωση των διακριτών τυχαίων μεταβλητών, αυτή η απόκλιση είναι η επιπλέον ποσότητα πληροφορίας (μετρημένη σε bit αν χρησιμοποιούμε λογάριθμο βάσης 2, αν και στη Μηχανική Μάθηση συνήθως χρησιμοποιούμε nat και το φυσικό λογάριθμο) που απαιτείται για να στείλουμε ένα μήνυμα που περιέχει σύμβολα σχεδιασμένα από μια πιθανοτική κατανομή  $P$ , όταν χρησιμοποιούμε μια κωδικοποίηση που έχει σχεδιαστεί για να ελαχιστοποιεί το μέγεθος των μηνυμάτων που

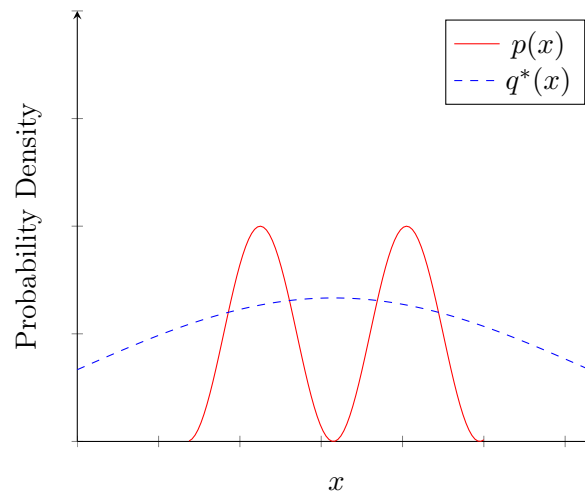
εξάγονται από την κατανομή  $Q$ . Η KL κατανομή έχει πολλές χρήσιμες ιδιότητες, κυρίως το ότι είναι μη αρνητική. Η KL απόκλιση είναι 0 αν και μόνο αν η  $P$  και η  $Q$  είναι η ίδια κατανομή, στην περίπτωση των διακριτών μεταβλητών, ή είναι ίσες "σχεδόν παντού", στην περίπτωση των συνεχών μεταβλητών. Επειδή η KL απόκλιση είναι μη αρνητική και μετράει τη διαφορά μεταξύ δύο κατανομών, συχνά έχει την έννοια ότι μετράει κάποιου είδους απόσταση μεταξύ αυτών των κατανομών. Ωστόσο, δεν είναι πραγματικό μέτρο απόστασης, επειδή δεν είναι συμμετρικό:  $D_{KL}(P \parallel Q) \neq D_{KL}(Q \parallel P)$  για κάποιο  $P$  και  $Q$ . Αυτή η ασυμμετρία σημαίνει ότι υπάρχουν σημαντικές συνέπειες στον αν θα επιλέξουμε να χρησιμοποιήσουμε  $D_{KL}(P \parallel Q)$  ή  $D_{KL}(Q \parallel P)$ . Μια ποσότητα που είναι στενά συνδεδεμένη με την KL απόκλιση είναι το cross-entropy,  $H(P, Q) = H(P) + D_{KL}(P \parallel Q)$ , η οποία είναι παρόμοια με την KL απόκλιση, αλλά χωρίς τον όρο στα αριστερά:

$$H(P, Q) = -E_{X \sim P} \log(Q(X))$$

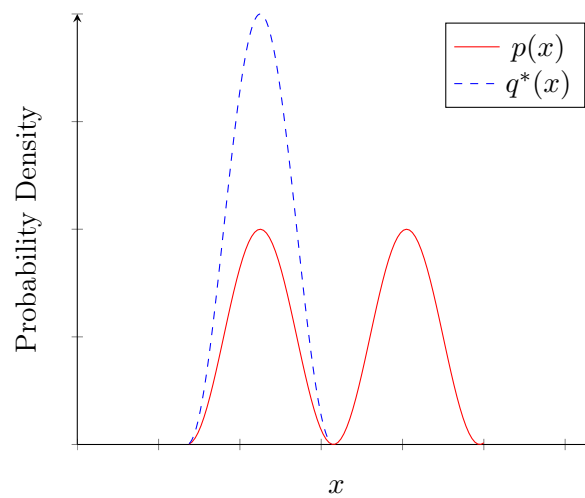
Η ελαχιστοποίηση του cross-entropy ως προς  $Q$  είναι ισοδύναμη με την ελαχιστοποίηση της KL απόκλισης, επειδή το  $Q$  δεν συμμετέχει στον παραλειπόμενο όρο. Υπολογίζοντας πολλές από αυτές τις ποσότητες, πολύ συχνά μπορεί να συναντήσουμε εκφράσεις της μορφής  $0 \cdot \log 0$ . Κατά κανόνα, στα πλαίσια της Θεωρίας Πληροφοριών, χειριζόμαστε αυτές τις εκφράσεις ως  $\lim_{x \rightarrow 0} x \cdot \log x = 0$ .

Η KL απόκλιση είναι ασύμμετρη:

$$q^* = \operatorname{argmin}_q D_{KL}(p || q)$$



$$q^* = \operatorname{argmin}_q D_{KL}(q || p)$$



## A.7 Mean Absolute Error (MAE)

Μέσο Απόλυτο Σφάλμα

$$\frac{\sum_{i=1}^N |X_i - X_{recon\_i}|}{N}$$

## A.8 Mean Squared Error (MSE)

Μέσο Τετραγωνικό Σφάλμα

$$\frac{\sum_{i=1}^N (X_i - X_{recon\_i})^2}{N}$$

## A.9 Το Ξυράφι του Occam

Το ξυράφι του Occam είναι επιστημονική αρχή που δηλώνει ότι μεταξύ συγχρουόμενων υποθέσεων, αυτή με τους λιγότερες εικασίες πρέπει να υπερισχύει. Αποδίδεται στον Άγγλο φιλόσοφο Λογικής και φραγκισκανό μοναχό του 14ου αιώνα, Γουλιέλμο του Όκαμ. Αυτή η αρχή στα Λατινικά εκφράζεται ως εξής: **"Pluralitas non est ponenda sine neccesitate"**. Το γνωμικό αυτό πρωτοδιατυπώνεται από τους Πυθαγόρειους για περισσότερο από δυο χιλιετίες, όπως μας πληροφορεί ο φιλόσοφος Πρόκλος, του οποίου το Πανεπιστήμιο ήταν στην οδό Ηρώδου Αττικού, θαμμένο κάτω από το πεζοδρόμιο, μπροστά στο θέατρο του Διονύσου. Στα Αρχαία Ελληνικά διατυπώνεται ως εξής: **"ἐξ ἐλαχίστων καὶ ἀπλουστάτων ὑποθέσεων δεικνύναι τὰ ζητούμενα"**. Στην απλούστερη διατύπωσή του, το ξυράφι του Occam εκφράζεται ως εξής: **"Κανείς δεν θα πρέπει να προβαίνει σε περισσότερες εικασίες από όσες είναι απαραίτητες"**.

## A.10 Root Mean Squared Error (RMSE)

Ρίζα Μέσου Τετραγωνικού Σφάλματος

$$\sqrt{\frac{\sum_{i=1}^N (X_i - X_{recon\_i})^2}{N}}$$

## A.11 Κανόνες Ενημέρωσης Βαρών των Νευρωνικών Δικτύων

Αυτοί είναι οι γενικευμένοι κανόνες ενημέρωσης βαρών που χρησιμοποιούνται στους αλγόριθμους προς τα πίσω διάδοσης (back-propagation).

- **Batch gradient descent**

Κανόνες ενημέρωσης για τις παραμέτρους του κωδικοποιητή (encoder):

$$W = W - \eta \cdot \frac{\partial E}{\partial W}$$

- **Stochastic gradient descent**

Κανόνες ενημέρωσης για τις παραμέτρους του κωδικοποιητή (encoder):

$$\text{for each } x_i : W = W - \eta \cdot \frac{\partial E_i}{\partial W}$$

- **Batch gradient ascent**

Κανόνες ενημέρωσης για τις παραμέτρους του κωδικοποιητή (encoder):

$$W = W + \eta \cdot \frac{\partial E}{\partial W}$$

- **Stochastic gradient ascent**

Κανόνες ενημέρωσης για τις παραμέτρους του κωδικοποιητή (encoder):

$$\text{for each } x_i : W = W + \eta \cdot \frac{\partial E_i}{\partial W}$$

όπου  $W$  είναι τα βάρη του Νευρωνικού Δικτύου και  $E$  είναι η συνάρτηση κόστους που θα μεγιστοποιηθεί (gradient ascent) ή ελαχιστοποιηθεί (gradient descent).

# Παράρτημα Β

## Πιθανοτικές Κατανομές

### Β.1 Κατανομή Bernoulli

(Christopher M. Bishop 2006. [8]) Αυτή είναι η κατανομή για μια απλή δυαδική μεταβλητή  $x \in 0, 1$  που αναπαριστά, για παράδειγμα, το αποτέλεσμα της ρίψης ενός νομίσματος. Κυριαρχείται από μια και μόνο συνεχή παράμετρο  $\mu \in [0, 1]$  που αναπαριστά την πιθανότητα του  $x = 1$ .

$$Bern(x|\mu) = \mu^x \cdot (1 - \mu)^{1-x}$$

$$E[x] = \mu$$

$$VAR[x] = \mu \cdot (1 - \mu)$$

$$mode[x] = \begin{cases} 1, & \text{if } \mu \geq 0.5 \\ 0, & \text{διαφορετικά} \end{cases}$$

$$H[x] = -\mu \cdot \ln \mu - (1 - \mu) \cdot \ln(1 - \mu)$$

Η Bernoulli είναι μια ειδική περίπτωση της διωνυμικής κατανομής στην περίπτωση μιας μεμονωμένης παρατήρησης. Η συζυγής εκ των προτέρων (conjugate prior) κατανομή της για  $\mu$  είναι η κατανομή beta.

pmf	Mean	Variance
$\begin{cases} q = (1 - p), & \text{for } k = 0 \\ p, & \text{for } k = 1 \end{cases}$	p	$p \cdot (1 - p) = p \cdot q$

ΠΙΝΑΚΑΣ Β.1: Κατανομή Bernoulli, συνάρτηση μάζας πιθανότητας (pmf), μέση τιμή και διασπορά.

## B.2 Διωνυμική Κατανομή

(Christopher M. Bishop 2006. [8]) Η διωνυμική κατανομή μας δίνει την πιθανότητα της παρατήρησης  $m$  περιστατικών της μεταβλητής  $x = 1$  σε ένα σύνολο από  $N$  δείγματα από μια Bernoulli κατανομή, όπου η πιθανότητα της μεταβλητής  $x = 1$  είναι  $\mu \in [0, 1]$ .

$$\text{Bin}(k|N, \mu) = \binom{N}{k} \cdot \mu^k \cdot (1 - \mu)^{N-k}$$

$$E[k] = N \cdot \mu$$

$$\text{VAR}[k] = N \cdot \mu \cdot (1 - \mu)$$

$$\text{mode}[k] = \lfloor (N + 1) \cdot \mu \rfloor$$

όπου  $\lfloor (N + 1) \cdot \mu \rfloor$  συμβολίζει το μεγαλύτερο ολοκλήρωμα που είναι μικρότερο ή ίσο με  $(N + 1) \cdot \mu$ , και η ποσότητα:

$$\binom{N}{k} = \frac{N!}{k!(N - k)!}$$

συμβολίζει τον αριθμό των τρόπων που μπορούμε να επιλέξουμε  $m$  αντικείμενα από ένα σύνολο  $N$  πανομοιότυπων αντικειμένων.

Εδώ το  $m!$ , το οποίο προφέρεται ως 'παραγοντικό  $m$ ', συμβολίζει το γινόμενο  $m \times (m - 1) \times \dots \times 2 \times 1$ . Η ιδιαίτερη περίπτωση όπου το  $N = 1$  είναι γνωστή ως η κατανομή Bernoulli, και για μεγάλα  $N$  η διωνυμική κατανομή είναι προσεγγιστικά Gaussian. Η συζυγής εκ των προτέρων (conjugate prior) κατανομή της για  $\mu$  είναι η beta κατανομή.

pmf	Mean	Variance
$\binom{n}{k} = \frac{n!}{k!(n-k)!}$	$n \cdot p$	$n \cdot p \cdot (1 - p)$

ΠΙΝΑΚΑΣ Β.2: Διωνυμική κατανομή, συνάρτηση μάζας πιθανότητας (pmf), μέση τιμή και διασπορά.

### B.3 Κατανομή Gauss (ή Κανονική Κατανομή)

(Christopher M. Bishop 2006. [8]) Η Gaussian είναι η πιο ευρέως χρησιμοποιημένη κατανομή για συνεχείς μεταβλητές. Είναι επίσης γνωστή ως κανονική κατανομή. Στην περίπτωση μιας μεμονωμένης μεταβλητής  $x \in (-\infty, \infty)$  κυριαρχείται από δύο παραμέτρους, τη μέση τιμή  $\mu \in (-\infty, \infty)$  και τη διασπορά  $\sigma^2 > 0$ .

$$N(x|\mu, \sigma^2) = \frac{1}{\sqrt{2 \cdot \pi \sigma^2}} \cdot e^{-\frac{1}{2 \cdot \sigma^2} \cdot (x-\mu)^2}$$

$$E[x] = \mu$$

$$VAR[x] = \sigma^2$$

$$mode[x] = \mu$$

$$H[x] = \frac{1}{2} \cdot \ln \sigma^2 + \frac{1}{2} \cdot (1 + \ln(2 \cdot \pi))$$

Το αντίστροφο της διασποράς  $\tau = \frac{1}{\sigma^2}$  ονομάζεται ακρίβεια, και η τετραγωνική ρίζα της διασποράς  $\Sigma$  ονομάζεται τυπική απόκλιση. Η συζυγής εκ των προτέρων (conjugate prior) κατανομή της για  $\mu$  είναι πάλι η κατανομή Gauss, ενώ η συζυγής εκ των προτέρων κατανομή της για  $\tau$  είναι η κατανομή Gaussian-gamma. Για το διάνυσμα  $D$  διαστάσεων  $x$ , η κατανομή Gauss κυριαρχείται από το  $D$ -διάστατο διάνυσμα μέσων τιμών  $\mu$  και ο  $D \times D$  πίνακας συνδιακύμανσης  $\Sigma$  πρέπει να είναι συμμετρικός και θετικά ορισμένος.

$$N(x|\mu, \Sigma) = \frac{1}{(2 \cdot \pi)^{D/2}} \cdot \frac{1}{\sqrt{|\Sigma|}} \cdot e^{-\frac{1}{2} \cdot (x-\mu)^T \cdot \Sigma^{-1} \cdot (x-\mu)}$$

$$E[x] = \mu$$

$$Cov[x] = \Sigma$$

$$mode[x] = \mu$$

$$H[x] = \frac{1}{2} \cdot \ln |\Sigma| + \frac{D}{2} \cdot (1 + \ln(2 \cdot \pi))$$

Ο αντίστροφος του πίνακα συνδιακύμανσης  $\Lambda = \Sigma^{-1}$  είναι ο πίνακας ακρίβειας, ο οποίος είναι επίσης θετικός και θετικά ορισμένος. Σύμφωνα με το κεντρικό οριακό θεώρημα, οι μέσες τιμές τυχαίων μεταβλητών τείνουν σε κατανομή Gauss και το άθροισμα δύο Gaussian μεταβλητών είναι πάλι Gaussian. Η Gaussian είναι κατανομή που μεγιστοποιεί την εντροπία για μια δεδομένη διασπορά (ή συνδιακύμανση). Κάθε γραμμικός μετασχηματισμός μιας Gaussian τυχαίας μεταβλητής είναι πάλι Gaussian. Η περιθώρια (marginal) κατανομή μιας πολυποίκιλης (multivariate) Gaussian ως προς ένα υποσύνολο των μεταβλητών είναι και αυτή Gaussian, και παρομοίως η υπό συνθήκη (conditional)



κατανομή είναι επίσης Gaussian. Αν έχουμε περιθώρια κατανομή Gauss για  $x$  και υπό συνθήκη κατανομή Gauss για  $y$  δεδομένου του  $x$  υπό τη μορφή:

$$P(x) = N(x|\mu, \Lambda^{-1})$$

$$P(y|x) = N(y|A \cdot x + b, L^{-1})$$

τότε η περιθώρια κατανομή του  $y$ , και η υπό συνθήκη κατανομή του  $x$  δεδομένου του  $y$ , ορίζονται ως εξής:

$$P(y) = N(y|A \cdot \mu + b, L^{-1} + A \cdot \Lambda^{-1} \cdot A^T)$$

$$P(x|y) = N(x|\Sigma A^T \cdot L(y - b) + \Lambda \cdot \mu, \Sigma)$$

όπου:

$$\Sigma = (\Lambda + A^T \cdot L \cdot A)^{-1}.$$

Αν έχουμε μια από κοινού κατανομή Gauss  $N(x|\mu, \Sigma)$  με  $\Lambda = \Sigma^{-1}$  και ορίζουμε τους παρακάτω διαχωρισμούς:

$$x = \begin{bmatrix} x_a \\ x_b \end{bmatrix}, \mu = \begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix}, \Lambda = \begin{bmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{bmatrix}$$

τότε η υπό συνθήκη κατανομή  $P(x_a|x_b)$  δίνεται από τον τύπο:

$$P(x_a|x_b) = N(x|\mu_{a|b}, \Lambda_{aa}^{-1})$$

$$\mu_{a|b} = \mu_a - \Lambda_{aa}^{-1} \cdot \Lambda_{ab} \cdot (x_b - \mu_b)$$

και η περιθώρια κατανομή  $P(x_a)$  δίνεται από τον τύπο:

$$P(x_a) = N(x_a|\mu_a, \Sigma_a).$$

pmf	Mean	Variance
$\frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{1}{2\sigma^2} \cdot (x-\mu)^2}$	$\mu$	$\sigma^2$

ΠΙΝΑΚΑΣ Β.3: Κατανομή Gauss, συνάρτηση μάζας πιθανότητας (pmf), μέση τιμή και διασπορά.

## Β.4 Νόμος Μεγάλος Αριθμών (N.M.A.)

(Σταύρος Τουμπής, Γιάννης Κοντογιάννης 2011. [9]) Το 1<sup>st</sup> θεμελιώδες θεώρημα της θεωρίας πιθανοτήτων είναι ο Νόμος των Μεγάλων Αριθμών (N.M.A.).

### Ορισμός

Έστω  $X_i, i = 1, 2, \dots$  μια ακολουθία από ανεξάρτητες τυχαίες μεταβλητές (T.M.) με την ίδια κατανομή, μέση τιμή  $E(X) = \mu$  και διασπορά  $VAR(X) = \sigma^2$ . Το 1<sup>st</sup> θεμελιώδες αποτέλεσμα των πιθανοτήτων, ο N.M.A., μας εξασφαλίζει ότι η πιθανότητα ο εμπειρικός μέσος όρος:

$$\bar{X}_N \approx \frac{1}{N} \cdot \sum_{i=1}^N X_i$$

να παρουσιάζει απόκλιση από το  $\mu$  μεγαλύτερη από  $\epsilon$  τείνει στο 0, δηλαδή:

$$\lim_{N \rightarrow \infty} P(|\bar{X}_N - \mu| > \epsilon) = 0$$

Επιπλέον, το ακριβές ποσοτικό άνω φράγμα για την πιθανότητα ο εμπειρικός μέσος όρος  $\bar{X}_N$ : να απέχει από τη μέση τιμή  $\mu$  κατά τουλάχιστον  $\epsilon$  έχει αποδειχθεί ότι είναι ίσο με:

$$P(|\bar{X}_N - \mu| \geq \epsilon) \leq \frac{\sigma^2}{N \cdot \epsilon^2}$$

όπου:  $E(\bar{X}_N) = \mu$  και  $Var(\bar{X}_N) = \frac{\sigma^2}{N}$ , για κάθε  $N$

## Β.5 Κεντρικό Οριακό Θεώρημα (Κ.Ο.Θ.)

(Σταύρος Τουμπής, Γιάννης Κοντογιάννης 2011. [9]) Το 2<sup>nd</sup> θεμελιώδες θεώρημα της θεωρίας πιθανοτήτων είναι το Κεντρικό Οριακό Θεώρημα (Κ.Ο.Θ.). Μέσω αυτού του θεωρήματος μπορούμε να εκτιμήσουμε με υψηλή ακρίβεια το όριο  $\lim_{N \rightarrow \infty} P(|\bar{X}_N - \mu| > \epsilon)$ . Συγκεκριμένα, το Κ.Ο.Θ. μας λέει ότι υπό συγκεκριμένες περιστάσεις, η κατανομή του εμπειρικού μέσου  $\bar{X}_N$  μπορεί να υπολογιστεί με υψηλή ακρίβεια από την Gaussian (γνωστή και ως κανονική) κατανομή με παραμέτρους  $\mu$  και  $\frac{\sigma^2}{N}$ , όπου  $\sigma^2$  η διασπορά των τυχαίων μεταβλητών (Τ.Μ.)  $X_i$ .

### Ορισμός

Έστω μια ακολουθία από ανεξάρτητες τυχαίες μεταβλητές (Τ.Μ.)  $X_1, X_2, \dots$  που έχουν όλες την ίδια κατανομή και κατά συνέπεια την ίδια μέση τιμή  $\mu = E(X_i)$  και την ίδια διασπορά  $\sigma^2 = VAR(X_i)$ .

Έστω το κανονικοποιημένο άθροισμα:

$$\bar{S}_N = \frac{(\frac{1}{N} \cdot \sum_{i=1}^N X_i - \mu)}{\sigma \cdot \sqrt{\frac{1}{N}}} = \frac{1}{\sigma \cdot \sqrt{N}} \cdot \sum_{i=1}^N (X_i - \mu) = \frac{(\sum_{i=1}^N X_i) - N \cdot \mu}{\sigma \cdot \sqrt{N}}$$

Έστω επίσης  $a, b \in \mathbb{R}, a < b$ . Ισχύουν τα ακόλουθα, καθώς  $N \rightarrow \infty$ :

- $P(a \leq \bar{S}_N \leq b) \rightarrow \Phi(b) - \Phi(a)$
- $P(\bar{S}_N \leq b) \rightarrow \Phi(b)$
- $P(\bar{S}_N \geq a) \rightarrow 1 - \Phi(a)$

όπου  $\Phi(x)$  είναι η Gaussian (γνωστή και ως κανονική) κατανομή

## Παράρτημα C

# Προγραμματιστικές Υλοποιήσεις για Variational Autoencoder σε Python

### C.1 Υλοποίηση VAE με χρήση βιβλιοθήκης TensorFlow

Ακολουθεί μια υλοποίηση VAE αλγορίθμου, με βιβλιοθήκη TensorFlow:

```
1 import numpy as np
2
3 X_train # matrix with the original data
4 y_train # labels of the matrix "X_train"
5 N # the number of training examples
6 input_dim # the input dimensionality D
7 hidden_encoder_dim # number of neurons in the encoder
8 hidden_decoder_dim # number of neurons in the decoder
9 latent_dim # Z_dim
10 epochs # the number of total epochs for the training
11 batch_size # the number of data trained in each iteration
12 learning_rate # the rate of learning "eta" for the optimization algorithm
13 log_dir # the location where the tensorboard graph will be saved
14
15 x, loss_summ, apply_updates, summary_op, saver, elbo, x_recon_samples =
16     vae(batch_size, input_dim, hidden_encoder_dim, hidden_decoder_dim, latent_dim,
17         lr)
18 X_recon = np.zeros((N, input_dim))
19
20 with tf.Session() as sess:
21     summary_writer = tf.summary.FileWriter(log_dir, graph=sess.graph)
22
23     for epoch in range(epochs):
24         iterations = int(N / batch_size)
25
26         for i in range(iterations):
27             start_index = i * batch_size
28             end_index = (i + 1) * batch_size
29
30             batch_data = X_train[start_index:end_index, :]
```

```
31     batch_labels = y_train[start_index:end_index]
32
33     feed_dict = {x: batch_data}
34     loss_str, _, summary_str, cur_elbo, cur_samples =
35         sess.run([loss_summ, apply_updates, summary_op,
36                 elbo, x_recon_samples], feed_dict=feed_dict)
37
38     X_recon[start_index:end_index, :] = cur_samples
39
40     summary_writer.add_summary(loss_str, epoch)
41     summary_writer.add_summary(summary_str, epoch)
```

ΚΩΔΙΚΑΣ C.1: Κύριος Βρόχος του VAE σε TensorFlow

Ακολουθεί η υλοποίηση της συνάρτησης "vae", της υλοποίησης με TensorFlow:

```

1 import tensorflow as tf
2
3 # HELPER FUNCTIONS #
4
5 def initialize_weight_variable(shape, name=''):
6     initial = tf.truncated_normal(shape, stddev=0.001,
7                                   name='truncated_normal_' + name)
8     return tf.Variable(initial, name=name)
9
10
11 def initialize_bias_variable(shape, name=''):
12     initial = tf.constant(0., shape=shape)
13     return tf.Variable(initial, name=name)
14
15 #####
16
17 def vae(batch_size, input_dim, hidden_encoder_dim, hidden_decoder_dim,
18         latent_dim, lr=0.01):
19     # Input placeholder
20     with tf.name_scope('input_data'):
21         x = tf.placeholder('float', [batch_size, input_dim], name='X')
22
23     # ===== Q(Z|X) = Q(Z) - Encoder NN ===== #
24
25     # The encoder is a neural network with 2 hidden layers.
26     with tf.name_scope('encoder'):
27         with tf.name_scope('Phis'):
28             # phi1: M1 x D
29             phi1 = initialize_weight_variable([hidden_encoder_dim,
30                                               input_dim], name='phi1')
31             # bias_phi1: 1 x M1
32             bias_phi1 = initialize_bias_variable([hidden_encoder_dim],
33                                                  name='bias_phi1')
34
35             # phi_mu: Z_dim x M1
36             phi_mu = initialize_weight_variable([latent_dim,
37                                                hidden_encoder_dim], name='phi_mu')
38             # bias_phi_mu: 1 x Z_dim
39             bias_phi_mu = initialize_bias_variable([latent_dim],
40                                                    name='bias_phi_mu')
41
42             # phi_logvar: Z_dim x M1
43             phi_logvar = initialize_weight_variable([latent_dim,
44                                                    hidden_encoder_dim], name='phi_logvar')
45             # bias_phi_logvar: 1 X Z_dim
46             bias_phi_logvar = initialize_bias_variable([latent_dim],

```

```

47         name='bias_phi_logvar')
48
49     with tf.name_scope('hidden_layer1'):
50         # Hidden layer 1 activation function of the encoder
51         # hidden_layer_encoder: N x M1
52         # RELU
53         hidden_layer_encoder = tf.nn.relu(tf.nn.bias_add(tf.matmul(
54             x, tf.transpose(phi1)),
55             bias_phi1))
56
57     with tf.name_scope('hidden_layer2_mu'):
58         # mu_encoder: N x Z_dim
59         mu_encoder = tf.nn.bias_add(tf.matmul(hidden_layer_encoder,
60             tf.transpose(phi_mu)), bias_phi_mu)
61
62     with tf.name_scope('hidden_layer2_logvar'):
63         # the log sigma^2 of the encoder
64         # logvar_encoder: N x Z_dim
65         logvar_encoder = tf.nn.bias_add(tf.matmul(hidden_layer_encoder,
66             tf.transpose(phi_logvar)),
67             bias_phi_logvar)
68
69     with tf.name_scope('sample_E'):
70         # Sample epsilon
71         # epsilon: N x Z_dim
72         epsilon = tf.random_normal((batch_size, latent_dim),
73             mean=0.0, stddev=1.0, name='epsilon')
74
75     with tf.name_scope('construct_Z'):
76         # Sample epsilon from the Gaussian distribution. #
77         # std_encoder: N x Z_dim
78         std_encoder = tf.exp(logvar_encoder / 2)
79         # z: N x Z_dim
80         z = tf.add(mu_encoder, tf.multiply(std_encoder, epsilon),
81             name='Z')
82
83     # ===== P(X|Z) - Decoder NN ===== #
84
85     # The encoder is a neural network with 2 hidden layers.
86     with tf.name_scope('decoder'):
87         with tf.name_scope('Thetas'):
88             # theta1: M2 x Z_dim
89             theta1 = initialize_weight_variable([hidden_decoder_dim,
90                 latent_dim], name='theta1')
91             # bias_theta1: 1 x M2
92             bias_theta1 = initialize_bias_variable([hidden_decoder_dim],
93                 name='bias_theta1')
94

```

```

95     # theta2: D x M2
96     theta2 = initialize_weight_variable([input_dim,
97                                         hidden_decoder_dim], name='theta2')
98     # bias_theta2: 1 x D
99     bias_theta2 = initialize_bias_variable([input_dim],
100                                           name='bias_theta2')
101
102     with tf.name_scope('hidden_layer1'):
103         # Hidden layer 1 activation function of the decoder
104         # hidden_layer_decoder: N x M2
105         # RELU
106         hidden_layer_decoder = tf.nn.relu(tf.nn.bias_add(tf.matmul(z,
107                                                             tf.transpose(theta1)), bias_theta1))
108
109     with tf.name_scope('output_layer'):
110         # x_hat: N x D
111         x_hat = tf.nn.bias_add(tf.matmul(hidden_layer_decoder,
112                                          tf.transpose(theta2)), bias_theta2)
113
114     with tf.name_scope('reconstructed_data'):
115         # X_recon_samples: N x D, reconstructed data
116         x_recon_samples = tf.nn.sigmoid(x_hat, name='X_recon_samples')
117
118     # ===== TRAINING ===== #
119
120     with tf.name_scope('ELBO'):
121         with tf.name_scope('reconstruction_cost'):
122             # log P(X)
123             # reconstruction_cost: N x 1
124             reconstruction_cost = tf.reduce_sum(tf.nn.
125                                                  sigmoid_cross_entropy_with_logits(logits=x_hat,
126                                                                 labels=x), reduction_indices=1)
127
128         with tf.name_scope('KL_divergence'):
129             # KLD: N x 1
130             KLD = 0.5 * tf.reduce_sum(1 + logvar_encoder -
131                                       tf.square(mu_encoder) - tf.exp(logvar_encoder),
132                                       reduction_indices=1)
133
134         elbo = tf.reduce_mean(reconstruction_cost - KLD,
135                               name='lower_bound')
136
137     loss_summ = tf.summary.scalar('ELBO', elbo)
138
139     phis = [phi1, bias_phi1,
140            phi_mu, bias_phi_mu,
141            phi_logvar, bias_phi_logvar]
142     thetas = [theta1, bias_theta1,

```



```

143         theta2, bias_theta2]
144     var_list = phis + thetas
145
146     # Adam Optimizer (WORKS BEST!) #
147     grads_and_vars = tf.train.AdamOptimizer(learning_rate=lr). \
148         compute_gradients(loss=elbo, var_list=var_list)
149     apply_updates = tf.train.AdamOptimizer(learning_rate=lr). \
150         apply_gradients(grads_and_vars=grads_and_vars)
151
152     # add op for merging summary
153     summary_op = tf.summary.merge_all()
154
155     # add Saver ops
156     saver = tf.train.Saver()
157
158     return x, loss_summ, apply_updates, summary_op, saver,
159         elbo, x_recon_samples

```

#### ΚΩΔΙΚΑΣ C.2: Συνάρτηση VAE σε TensorFlow

**ΣΗΜΕΙΩΣΗ:** Για το βελτιστοποιητή (optimizer) των παραγώγων (gradient), χρησιμοποιήσαμε τον Adam optimizer (tf.train.AdamOptimizer), επειδή επιφέρει τα καλύτερα αποτελέσματα. Όλοι οι optimizer που παρέχει η βιβλιοθήκη TensorFlow είναι οι εξής:

- tf.train.GradientDescentOptimizer
- tf.train.AdadeltaOptimizer
- tf.train.AdagradOptimizer
- tf.train.AdagradDAOptimizer
- tf.train.MomentumOptimizer
- **tf.train.AdamOptimizer**
- tf.train.FtrlOptimizer
- tf.train.ProximalGradientDescentOptimizer
- tf.train.ProximalAdagradOptimizer
- tf.train.RMSPropOptimizer

## C.2 Υλοποίηση VAE με χρήση βιβλιοθήκης PyTorch

Ακολουθεί μια υλοποίηση VAE αλγορίθμου, με βιβλιοθήκη PyTorch:

```
1 import numpy as np
2
3 X_train # matrix with the original data
4 y_train # labels of the matrix "X_train"
5 N # the number of training examples
6 input_dim # the input dimensionality D
7 hidden_encoder_dim # number of neurons in the encoder
8 hidden_decoder_dim # number of neurons in the decoder
9 latent_dim # Z_dim
10 epochs # the number of total epochs for the training
11 batch_size # the number of data trained in each iteration
12 learning_rate # the rate of learning "eta" for the optimization algorithm
13
14 # initialize_weights:
15 # A function that initializes the weights of the encoder and
16 # the decoder neural networks.
17 params # the weights of the encoder and the decoder neural networks,
18         # that are updated on each iteration
19 solver # the optimizer used for training (e.g SGD)
20 params, solver = initialize_weights(input_dim, hidden_encoder_dim, \
                                     hidden_decoder_dim, latent_dim, lr)
21
22 for epoch in range(epochs):
23     iterations = int(N / batch_size)
24
25     for i in range(iterations):
26         start_index = i * batch_size
27         end_index = (i + 1) * batch_size
28
29         batch_data = X_train[start_index:end_index, :]
30         batch_labels = y_train[start_index:end_index]
31
32         cur_samples, cur_elbo = train(batch_data, batch_size, \
                                     latent_dim, params, solver)
33
34         X_train[start_index:end_index, :] = cur_samples
```

ΚΩΔΙΚΑΣ C.3: Κύριος Βρόχος του VAE σε PyTorch

Ακολουθεί η υλοποίηση της συνάρτησης "initialize\_weights" της υλοποίησης με PyTorch:

```

1 import numpy as np
2 import torch
3 import torch.optim as optim
4 from torch.autograd import Variable
5
6
7 def xavier_init(size):
8     in_dim = size[0]
9     xavier_stddev = 1. / np.sqrt(in_dim / 2.)
10    return Variable(torch.randn(*size) * xavier_stddev,
11                    requires_grad=True)
12
13
14 def initialize_weights(X_dim, hidden_encoder_dim, hidden_decoder_dim,
15                       Z_dim, lr=0.01):
16
17     # ===== Encoder Parameters: Phis ===== #
18
19     # M1 x D
20     phi1 = xavier_init(size=[hidden_encoder_dim, X_dim])
21     # 1 x M1
22     bias_phi1 = Variable(torch.zeros(hidden_encoder_dim),
23                          requires_grad=True)
24
25     # Z_dim x M1
26     phi_mu = xavier_init(size=[Z_dim, hidden_encoder_dim])
27     # 1 x Z_dim
28     bias_phi_mu = Variable(torch.zeros(Z_dim),
29                            requires_grad=True)
30
31     # Z_dim x M1
32     phi_logvar = xavier_init(size=[Z_dim, hidden_encoder_dim])
33     # 1 x Z_dim
34     bias_phi_logvar = Variable(torch.zeros(Z_dim),
35                                requires_grad=True)
36
37     # ===== Decoder Parameters: Thetas ===== #
38
39     # M2 x Z_dim
40     theta1 = xavier_init(size=[hidden_decoder_dim, Z_dim])
41     # 1 x M2
42     bias_theta1 = Variable(torch.zeros(hidden_decoder_dim),
43                            requires_grad=True)
44
45     # D x M2
46     theta2 = xavier_init(size=[X_dim, hidden_decoder_dim])

```

```

47  # 1 x D
48  bias_theta2 = Variable(torch.zeros(X_dim),
49                          requires_grad=True)
50
51  # ===== TRAINING =====
52
53  phis = [phi1, bias_phi1,
54          phi_mu, bias_phi_mu,
55          phi_logvar, bias_phi_logvar]
56  thetas = [theta1, bias_theta1,
57            theta2, bias_theta2]
58  params = phis + thetas
59
60  solver = optim.Adam(params, lr=lr)
61
62  return params, solver

```

ΚΩΔΙΚΑΣ C.4: Αρχικοποίηση των Βαρών του VAE σε PyTorch

**ΣΗΜΕΙΩΣΗ:** Για το βελτιστοποιητή (optimizer) των παραγώγων (gradient), χρησιμοποιήσαμε τον Adam optimizer (`torch.optim.Adam`), επειδή επιφέρει τα καλύτερα αποτελέσματα. Όλοι οι optimizer που παρέχει η βιβλιοθήκη PyTorch είναι οι εξής:

- `torch.optim.Adadelta`
- `torch.optim.Adagrad`
- **`torch.optim.Adam`**
- `torch.optim.SparseAdam`
- `torch.optim.AdaMax`
- `torch.optim.ASGD`
- `torch.optim.LBFGS`
- `torch.optim.RMSprop`
- `torch.optim.Rprop`
- `torch.optim.SGD`

Ακολουθεί η υλοποίηση της συνάρτησης "train" της υλοποίησης με PyTorch:

```

1 import torch
2 import torch.nn.functional as nn
3 from torch.autograd import Variable
4
5
6 # @: denotes matrix multiplication
7 def train(x, mb_size, Z_dim, params, solver):
8     x = Variable(torch.from_numpy(x).float())
9
10    # This
11    phi1 = params[0] # M1 x D
12    bias_phi1 = params[1] # 1 x M1
13    phi_mu = params[2] # Z_dim x M1
14    bias_phi_mu = params[3] # 1 x Z_dim
15    phi_logvar = params[4] # Z_dim x M1
16    bias_phi_logvar = params[5] # 1 x Z_dim
17
18    # Thetas
19    theta1 = params[6] # M2 x Z_dim
20    bias_theta1 = params[7] # 1 x M2
21    theta2 = params[8] # D x M2
22    bias_theta2 = params[9] # 1 x D
23
24    # ===== Q(Z|X) = Q(Z) - Encoder NN ===== #
25
26    hidden_layer_encoder = nn.relu(x @ torch.transpose(phi1, 0, 1) +
27                                   bias_phi1.repeat(mb_size, 1))
28    mu_encoder = hidden_layer_encoder @ torch.transpose(phi_mu, 0, 1) +
29                bias_phi_mu.repeat(hidden_layer_encoder.size(0), 1)
30    logvar_encoder = hidden_layer_encoder @ torch.transpose(phi_logvar, 0, 1)
31                    + bias_phi_logvar.repeat(hidden_layer_encoder.size(0), 1)
32
33    # Sample epsilon from the Gaussian distribution. #
34    epsilon = Variable(torch.randn(mb_size, Z_dim))
35    # std_encoder: N x Z_dim
36    std_encoder = torch.exp(logvar_encoder / 2)
37    # Sample the latent variables Z. #
38    z = mu_encoder + std_encoder * epsilon
39
40    # ===== P(X|Z) - Decoder NN ===== #
41
42    hidden_layer_decoder = nn.relu(z @ torch.transpose(theta1, 0, 1) +
43                                   bias_theta1.repeat(z.size(0), 1))
44    x_hat = hidden_layer_decoder @ torch.transpose(theta2, 0, 1) +
45            bias_theta2.repeat(hidden_layer_decoder.size(0), 1)
46    x_recon_samples = nn.sigmoid(x_hat)

```

```
47
48 # Loss #
49 recon_loss = nn.binary_cross_entropy(x_recon_samples,
50                                     x, size_average=False) / mb_size
51 kl_loss = torch.mean(0.5 * torch.sum(1 + logvar_encoder
52                                     - mu_encoder ** 2 - torch.exp(logvar_encoder), 1))
53 elbo_loss = recon_loss - kl_loss
54
55 # Backward #
56 elbo_loss.backward()
57
58 # Update #
59 solver.step()
60
61 # Housekeeping #
62 for p in params:
63     p.grad.data.zero_()
64
65 # convert to numpy data type and return
66 x_recon_samples = x_recon_samples.data.numpy()
67 elbo_loss = elbo_loss.data.numpy()[0]
68
69 return x_recon_samples, elbo_loss
```

ΚΩΔΙΚΑΣ C.5: Συνάρτηση Εκπαίδευσης του VAE σε PyTorch

### C.3 Υλοποίηση VAE με χρήση βιβλιοθήκης Keras

Ακολουθεί υλοποίηση VAE αλγορίθμου, με βιβλιοθήκη Keras (χρησιμοποιώντας το TensorFlow ως backend):

```

1 import numpy as np
2 from keras.callbacks import TensorBoard
3 from keras.layers import Input, Dense
4 from keras.models import Model
5
6 X_train # matrix with the original train data
7 X_test  # matrix with the test data
8 X_valid # matrix with the cross-validation data
9 input_dim # the input dimensionality D
10 latent_dim # Z_dim
11 epochs # the number of total epochs for the training
12 batch_size # the number of data trained in each iteration
13 log_dir # the location where the tensorboard graph will be saved
14
15 # this is our input placeholder
16 input_img = Input(shape=(input_dim,))
17
18 # 'encoded' is the encoded representation of the input
19 encoded = Dense(latent_dim, activation='relu')(input_img)
20 # 'decoded' is the lossy reconstruction of the input
21 decoded = Dense(input_dim, activation='sigmoid')(encoded)
22
23 # this model maps an input to its reconstruction
24 autoencoder = Model(input_img, decoded)
25
26 # ===== Q(Z|X) = Q(Z) - Encoder NN ===== #
27
28 # this model maps an input to its encoded representation
29 encoder = Model(input_img, encoded)
30
31 # ===== P(X|Z) - Decoder NN ===== #
32
33 # create a placeholder for an encoded latent_dim input
34 encoded_input = Input(shape=(latent_dim,))
35 # retrieve the last layer of the autoencoder model
36 decoder_layer = autoencoder.layers[-1]
37 # create the decoder model
38 decoder = Model(encoded_input, decoder_layer(encoded_input))
39
40 # compile the model
41 # optimizer: Adadelta
42 # loss function: binary_crossentropy loss

```

```
43 autoencoder.compile(optimizer='adadelta', loss='binary_crossentropy')
44
45 # fit the model
46 autoencoder.fit(X_train, X_train,
47                 epochs=epochs,
48                 batch_size=batch_size,
49                 shuffle=True,
50                 validation_data=(X_valid, X_valid),
51                 callbacks=[TensorBoard(log_dir=log_dir,
52                                       histogram_freq=0, write_graph=False)])
53
54 # encode and decode some digits
55 # note that we take them from the *test* set
56 encoded_imgs = encoder.predict(X_test)
57 decoded_imgs = decoder.predict(encoded_imgs)
```

ΚΩΔΙΚΑΣ C.6: Κύριος Βρόχος του VAE σε Keras



## C.4 K-NN Αλγόριθμος Συμπλήρωσης Ελλιπών Τιμών σε Python

Ακολουθεί ο πλήρης κώδικας για τον K-NN αλγόριθμο συμπλήρωσης ελλιπών τιμών σε Python:

```

1 X_train_missing # TRAIN data with missing values
2 X_test_missing # TEST data with missing values
3 K = 1
4 missing_value = 0.5
5
6 X_test_predicted = kNNMatrixCompletion(X_train_missing,
7                                         X_test_missing, K, missing_value)
8
9 # Customized Euclidean distance function.
10 def sqdist(X_train, X_test_instance, missing_value):
11     Ntrain = X_train.shape[0]
12     D = X_test.shape[1]
13
14     X_train_common = np.array(X_train)
15     s = np.where(X_test_instance == missing_value)
16     X_train_common[:, s] = missing_value
17
18     X_test_common = np.zeros((Ntrain, D))
19     mean_values = np.mean(X_train, axis=0) # 1 x D array
20     for k in range(Ntrain):
21         X_test_common[k, :] = X_test_instance
22         s = np.where(X_train[k, :] == missing_value)
23         if len(s[0]) != 0:
24             X_test_common[k, s] = mean_values[s]
25
26     dist = np.matrix(np.sqrt(np.sum(np.square(X_train_common - X_test_common),
27                                         axis=1))).T
28     return dist
29
30 def kNNMatrixCompletion(X_train, X_test, K, missing_value,
31                         use_softmax_weights=True, binarize=False):
32     Ntest = X_test.shape[0]
33
34     X_test_predicted = np.array(X_test)
35     for i in range(Ntest):
36         print('data%i' % i)
37         X_test_i = np.squeeze(np.array(X_test[i, :]))
38
39         distances = sqdist(X_train, X_test_i, missing_value)
40
41         # sort the distances in ascending order
42         # and store the indices in a variable

```

```

43     closest_data_indices = \
44         (np.argsort(distances, axis=0)).tolist()
45
46     # select the top k indices
47     closest_k_data_indices = closest_data_indices[:K]
48     closest_k_data = \
49         np.squeeze(X_train[closest_k_data_indices, :])
50
51     closest_k_data_distances = \
52         np.squeeze(distances[closest_k_data_indices, :]).T
53
54     if use_softmax_weights:
55         # distribute weights, in the following manner:
56         # wk = e-k / sum{i=1}^K e-i
57         # such that: w1 + w2 + ... wK = 1
58         a = 1
59         weights = softmax(-a * closest_k_data_distances)
60     else:
61         # ALTERNATIVE: all weights equal to 1 / K
62         weights = np.ones((K, 1)) / float(K)
63
64     closest_k_data = np.multiply(closest_k_data, weights)
65
66     # sum across all rows for each column, returns a column vector
67     predicted_pixels = np.sum(closest_k_data, axis=0).T
68
69     X_test_predicted[i, np.where(X_test_i == missing_value)] = \
70         predicted_pixels[np.where(X_test_i == missing_value)].T
71
72     if binarize:
73         X_test_predicted = np.round(X_test_predicted + 0.3)
74
75     return X_test_predicted

```

ΚΩΔΙΚΑΣ C.7: K-NN αλγόριθμος συμπλήρωσης ελλειπόν τιμών σε Python

## C.5 VAE Αλγόριθμος Συμπλήρωσης Ελλιπών Τιμών σε PyTorch

Ακολουθεί ο πλήρης κώδικας για το VAE αλγόριθμο συμπλήρωσης ελλιπών τιμών σε Python, με χρήση **PyTorch**.

```

1 import numpy as np
2
3 X_train # matrix with the original data
4 y_train # labels of the matrix "X_train"
5 X_train_missing # matrix with missing values
6 X_train_masked # matrix with 0s where the pixel are missing and
7                 # 1s where the pixel are not missing
8 X_filled = np.array(X_train_missing) # the final matrix of predicted pixels
9 N # the number of training examples
10 input_dim # the input dimensionality D
11 hidden_encoder_dim # number of neurons in the encoder
12 hidden_decoder_dim # number of neurons in the decoder
13 latent_dim # Z_dim
14 epochs # the number of total epochs for the training
15 batch_size # the number of data trained in each iteration
16 learning_rate # the rate of learning "eta" for the optimization algorithm
17
18 # initialize_weights:
19 # initializes the weights of the encoder and the decoder.
20 # params: the weights of the encoder and the decoder
21 # solver: the optimizer used for training (e.g SGD)
22 params, solver = initialize_weights(input_dim, hidden_encoder_dim,
23                                     hidden_decoder_dim, latent_dim, lr)
24
25 for epoch in range(epochs):
26     iterations = int(N / batch_size)
27
28     for i in range(iterations):
29         start_index = i * batch_size
30         end_index = (i + 1) * batch_size
31
32         batch_data = X_filled[start_index:end_index, :]
33         batch_labels = y_train[start_index:end_index]
34         masked_batch_data = X_train_masked[start_index:end_index, :]
35
36         cur_samples, cur_elbo = train(batch_data, batch_size,
37                                     latent_dim, params, solver)
38
39         cur_samples = np.multiply(masked_batch_data, batch_data) +
40                         np.multiply(1 - masked_batch_data, cur_samples)
41         X_filled[start_index:end_index, :] = cur_samples

```

ΚΩΔΙΚΑΣ C.8: VAE αλγόριθμος συμπλήρωσης ελλιπών τιμών σε PyTorch

**ΣΗΜΕΙΩΣΕΙΣ:**

- Οι υλοποιήσεις των συναρτήσεων `"initialize_weights"` και `"train"` βρίσκονται στο Παράρτημα.
- Η υλοποίηση για το VAE αλγόριθμο συμπλήρωσης ελλιπών τιμών σε Python, με χρήση **TensorFlow**, δε συμπεριλαμβάνεται επειδή είναι παρόμοια με την υλοποίηση σε PyTorch.

# Βιβλιογραφία

- [1] Yoshua Bengio. Learning deep architectures for ai. *Foundations and Trends® in Machine Learning*, 2(1):1–127, 2009. ISSN 1935-8237. doi: 10.1561/22000000006. <http://dx.doi.org/10.1561/22000000006>.
- [2] Carl Doersch. Tutorial on variational autoencoders. *Variational Autoencoders*, June 2016. <https://arxiv.org/abs/1606.05908>.
- [3] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *Variational Autoencoders*, December 2013. <https://arxiv.org/abs/1312.6114>.
- [4] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. Deep learning. *Deep Learning*, 2016. <http://www.deeplearningbook.org>.
- [5] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, November 1998.
- [6] Alex Krizhevsky. Learning multiple layers of features from tiny images. *Machine Learning*, pages 34–37, April 2009.
- [7] A.S. Georgiades, P.N. Belhumeur, and D.J. Kriegman. From few to many: Illumination cone models for face recognition under variable lighting and pose. *IEEE Trans. Pattern Anal. Mach. Intelligence*, 23(6):643–660, 2001.
- [8] Christopher M. Bishop. Pattern recognition and machine learning. *Machine Learning*, 2006.
- [9] Γιάννης Κοντογιάννης και Σταύρος Τουμπής. Σημειώσεις μαθήματος "Πιθανοτήτων". *Οικονομικό Πανεπιστήμιο Αθηνών*, 2011.

# Σύνδεσμοι

- Εικόνες Google <https://images.google.com/>
- Analytics Vidhya,  
An Introduction to Implementing Neural Networks using TensorFlow  
<https://www.analyticsvidhya.com/blog/2016/10/an-introduction-to-implementing-neural-networks-using-tensorflow/>
- Arxiv Insights YouTube κανάλι,  
Variational Autoencoders  
<https://www.youtube.com/watch?v=9zKuYvjFFS8>
- Augustinus Kristiadi's Blog, Many flavors of Autoencoder  
<https://wiseodd.github.io/techblog/2016/12/03/autoencoders/>
- Christopher Bourez's blog,  
Symbolic computing and deep learning tutorial with Tensorflow/Theano: learn basic commands of 2 libraries for the price of 1  
<http://christopher5106.github.io/big/data/2016/03/06/symbolic-computing-and-deep-learning-tutorial-on-theano-and-google-tensorflow.html>
- Edureka YouTube κανάλι,  
Introduction To TensorFlow | Deep Learning Using TensorFlow | TensorFlow Tutorial | Edureka  
<https://www.youtube.com/watch?v=uh2Fh6df7Lg&t=1562s>
- Fast Forward Labs Blog,
  1. Introducing Variational Autoencoders (in Prose and Code)  
<http://blog.fastforwardlabs.com/2016/08/12/introducing-variational-autoencoders-in-prose-and.html>
  2. Under the Hood of the Variational Autoencoder (in Prose and Code)  
<http://blog.fastforwardlabs.com/2016/08/22/under-the-hood-of-the-variational-autoencoder-in.html>
- Jaan Altosaar, Tutorial - What is a Variational Autoencoder?  
<https://jaan.io/what-is-variational-autoencoder-vae-tutorial/>
- kevin frans, Variational Autoencoders Explained  
<http://kvfrans.com/variational-autoencoders-explained/>

- Nando de Freitas YouTube κανάλι,  
Deep Learning Lecture 14: Karol Gregor on Variational Autoencoders and Image Generation  
<https://www.youtube.com/watch?v=P78QYjWh5sM>
- Siraj Raval YouTube κανάλι,  
How to Generate Images - Intro to Deep Learning #14  
<https://www.youtube.com/watch?v=3-UDwk1U77s>
- studio otoro, Generating Large Images from Latent Vectors  
<http://blog.otoro.net/2016/04/01/generating-large-images-from-latent-vectors/>
- The Keras Blog, Building Autoencoders in Keras  
<https://blog.keras.io/building-autoencoders-in-keras.html>
- videolectures.net,  
Building Machines that Imagine and Reason: Principles and Applications of Deep Generative Models  
[http://videolectures.net/deeplearning2016\\_mohamed\\_generative\\_models/](http://videolectures.net/deeplearning2016_mohamed_generative_models/)
- videolectures.net,  
Variational Autoencoder and Extensions  
[http://videolectures.net/deeplearning2015\\_courville\\_autoencoder\\_extension/](http://videolectures.net/deeplearning2015_courville_autoencoder_extension/)
- Xitong Yang's Blog, Understanding the Variational Lower Bound  
<https://xyang35.github.io/2017/04/14/variational-lower-bound/>
- wiseodd GitHub, Generative Models  
<https://github.com/wiseodd/generative-models>

- Βικιπαίδεια,

1. Autoencoder  
<https://en.wikipedia.org/wiki/Autoencoder>
2. Bayes' theorem  
[https://en.wikipedia.org/wiki/Bayes'\\_theorem](https://en.wikipedia.org/wiki/Bayes'_theorem)
3. Softmax function  
[https://en.wikipedia.org/wiki/Softmax\\_function](https://en.wikipedia.org/wiki/Softmax_function)
4. Cross entropy  
[https://en.wikipedia.org/wiki/Cross\\_entropy](https://en.wikipedia.org/wiki/Cross_entropy)
5. Entropy (information theory)  
[https://en.wikipedia.org/wiki/Entropy\\_\(information\\_theory\)](https://en.wikipedia.org/wiki/Entropy_(information_theory))
6. Jensen's inequality  
[https://en.wikipedia.org/wiki/Jensen's\\_inequality](https://en.wikipedia.org/wiki/Jensen's_inequality)
7. Mean absolute error  
[https://en.wikipedia.org/wiki/Mean\\_absolute\\_error](https://en.wikipedia.org/wiki/Mean_absolute_error)
8. Mean squared error  
[https://en.wikipedia.org/wiki/Mean\\_squared\\_error](https://en.wikipedia.org/wiki/Mean_squared_error)
9. Kullback–Leibler divergence  
[https://en.wikipedia.org/wiki/Kullback-Leibler\\_divergence](https://en.wikipedia.org/wiki/Kullback-Leibler_divergence)
10. Ξυράφι του Όχαμ  
[https://el.wikipedia.org/wiki/Ξυράφι\\_του\\_Όχαμ](https://el.wikipedia.org/wiki/Ξυράφι_του_Όχαμ)
11. Root-mean-square error  
[https://en.wikipedia.org/wiki/Root-mean-square\\_deviation](https://en.wikipedia.org/wiki/Root-mean-square_deviation)
12. TensorFlow  
<https://en.wikipedia.org/wiki/TensorFlow>
13. Variational Bayesian methods  
[https://en.wikipedia.org/wiki/Variational\\_Bayesian\\_methods](https://en.wikipedia.org/wiki/Variational_Bayesian_methods)