



Outlines

- Overview
- Introduction
- Linear Algebra
- Probability
- Linear Regression
- Logistic Regression
- Linear Classification 1
- Linear Classification 2
- Kernel Methods
- Sparse Kernel Methods
- ixture Models and EM 1
- ixture Models and EM 2
- ural Networks 1
- ural Networks 2
- skip 1 Component Analysis
- encoders
- aphical Models 1
- Graphical Models 2
- Graphical Models 3
- Sampling
- Sequential Data 1
- Sequential Data 2

Statistical Machine Learning

Assignment Project Exam Help

Christian Walder

<https://eduassistpro.github.io>

College of Engineering and Computer Science
The Australian National University

Add WeChat edu_assist_pro
Canberra

Semester One, 2020.

(Many figures from C. M. Bishop, "Pattern Recognition and Machine Learning")



Motivation

CSIRO

Autoencoder

Assignment Project Exam Help

Part X

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



Motivation

Autoencoder

Number of layers

- expression of a function is **compact** when it has few computational elements, i.e. few degrees of freedom that need to be tuned by learning

- for a function to be better, it needs to be tuned by learning

- affine operations, sigmoid \Rightarrow logistic regression
 - fixed number of units (a.k.a. neurons)
- fixed kernel, affine operations \Rightarrow learnable
 - has two levels, with as many units as data points
- stacked neural network of multiple “linear transformation followed by a non-linearity” \Rightarrow deep neural network has arbitrary depth with arbitrary number of units per layer

Add WeChat `edu_assist_pro`



Motivation

Autoencoder

Easier to represent with more layers

An old result:

- functions that can be compactly represented by a depth k architecture might require an exponential number of computational elements to be represented by a depth $k - 1$ architecture
- Exa

<https://eduassistpro.github.io/>

A diagram of a parity circuit. It shows a sequence of inputs: p_1, p_2, \dots, p_d followed by t_1, t_2, \dots, t_o . The inputs are processed through a series of layers. The first layer has nodes labeled $1, d$. The second layer has nodes labeled $+$. The third layer has nodes labeled $0, o$.

- Theorem: NOT parity circuits of depth 2 have a size

Add WeChat `edu_assist_pro`

Analogous in modern deep learning:

- “Shallow networks require exponentially more parameters for the same number of modes” — Canadian deep learning mafia.



Motivation

Autoencoder

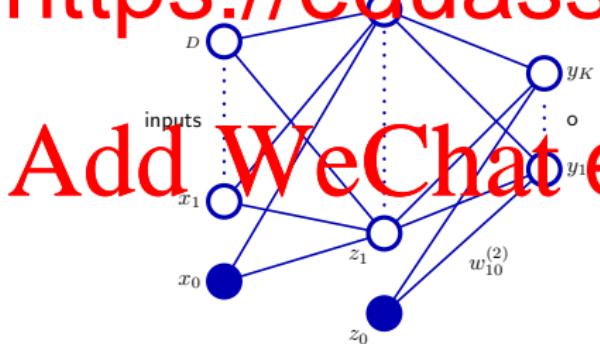
Recall: Multi-layer Neural Network Architecture

$$y_k(x, w) = \sigma \left(\sum_{j=0}^M w_{kj}^{(2)} h \left(\sum_{i=0}^D v_i^{(1)} x_i \right) \right)$$

Assignment Project Exam Help

where w

<https://eduassistpro.github.io>



Add WeChat edu_assist_pro

We could add more hidden layers



Motivation

Autoencoder

Assignment Project Exam Help

- Deep architectures get stuck in local minima or plateaus
- As a
gen
- Har
- 1 or 2 hidden layers seem to perform better
- 2006: Unsupervised pre-training, find distrib
representation

Add WeChat edu_assist_pro



Motivation

Autoencoder

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Bengio, "Learning Deep Architectures for AI", 2009



Motivation

Autoencoder

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

AlexNet / VGG-F network visualized by mNeuron.



Motivation

Autoencoder

Recall: PCA

- Idea: Linearly project the data points onto a lower dimensional subspace such that

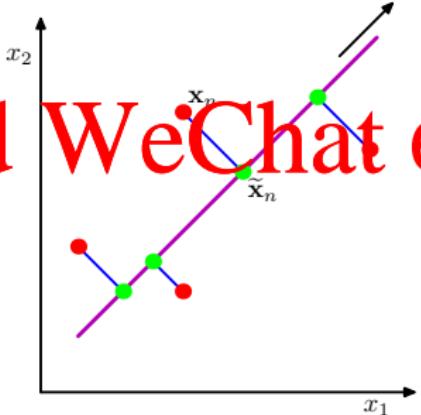
- the variance of the projected data is maximised; or
- the distortion error from the projection is minimised

- Both formulation lead to the same result.
- Ne

prin

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro





Motivation

Autoencoder

Assignment Project Exam Help

Principal Component Analysis is a linear transformation
(because it is a projection)

- The
- Lin
- Let

<https://eduassistpro.github.io>

$$ST(X + X') = ST(X) + S$$

- Similarly for multiplication with a scalar

⇒ multiple PCA layers pointless



Motivation



Autoencoder

Multiple PCA layers? - projection

- Let $X^T X = U \Lambda U^T$ be the eigenvalue decomposition of the covariance matrix (what is assumed about the mean?). Define U_k to be the matrix of the first k columns of U , for the k largest eigenvalues. Define Λ_k similarly
- Co
prin

<https://eduassistpro.github.io/>

- We perform PCA a second time, $Z^T Z$
- By the definition of Z and $X^T X$ and the ort

$$Z^T Z = (XU_k)^T (XU_k) = U_k^T X^T X U_k = U_k U \Lambda U^T U_k = \Lambda_k$$

- Hence $\Lambda_Z = \Lambda_k$ and V is the identity, therefore the second PCA has no effect

⇒ again, multiple PCA layers pointless



Motivation

Autoencoder

Autoencoder

- An autoencoder is trained to **encode** the input x into some representation $c(x)$ so that the input can be reconstructed from that representation
- the target output of the autoencoder is the autoencoder input itself

• Wit

crit

spa

<https://eduassistpro.github.io/>

- If the hidden layer is nonlinear, the autoencoder encodes differently from PCA, with the ability to capture more aspects of the input distribution
- Let f be the **decoder**. We want to minimise the reconstruction error

$$\sum_{n=1}^N \ell(x_n, f(c(x_n)))$$

Add WeChat edu_assist_pro



Motivation



Autoencoder

Cost function

- Recall: $f(c(x))$ is the reconstruction produced by the network
- Minimisation of the negative log likelihood of the reconstruction, given the encoding $c(x)$

Assignment Project Exam Help

<https://eduassistpro.github.io/>

- If x is binary
- If the inputs x_i are either binary or consider binomial probabilities, then the loss function is cross entropy

Add WeChat edu_assist_pro

$$-\log P(x|c(x)) = -x_i \log f_i(c(x)) + (1 - x_i) \log(1 - f_i(c(x)))$$

where $f_i(\cdot)$ is the i^{th} component of the decoder



Motivation

Autoencoder

Undercomplete Autoencoder

Assignment Project Exam Help

- Co
- $c(x)$
- Ca
- exa
- Hope code $c(x)$ is a distributed representa
- captures the main factors of variation in the data

Add WeChat edu_assist_pro



Motivation

Autoencoder

Stacking autoencoders

- Let c_j and f_j be the encoder and corresponding decoder of the j^{th} layer
- Let z_j be the representation after the encoder c_j
- We can define multiple layers of autoencoders recursively.
- For the j^{th} layer, $z_j = f_j(c_j(z_{j-1}))$
- Because feature z_2 can capture more complex patterns

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



Motivation

Autoencoder

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



Motivation

Autoencoder

Assignment Project Exam Help

- Latent features z_j in layer j can capture high level patterns.

$$z_j = c_j(c_{j-1}(\dots c_2(c_1(x)) \dots))$$

- The task is to construct a model $f(x)$ that takes input x and outputs y .
• In contrast to the generative models, the encoder c_j is constructed in an unsupervised fashion.
- Discard the decoding layers, and directly use the output of the encoder c_j as input to a supervised training method, such as logistic regression or SVM.
- Various such pre-trained networks are available on-line, e.g [VGG-19](https://eduassistpro.github.io/).

Add WeChat `edu_assist_pro`



Motivation

Autoencoder

Xavier Initialisation / ReLU

- Layer-wise unsupervised pre-training helps by extracting useful features for subsequent supervised backprop.
- Pre-training also avoids **saturation** (large magnitude arguments to, say, sigmoidal units).
- Simpler **Xavier initialization** can also avoid saturation.
- Let t acti

<https://eduassistpro.github.io>

$$\text{VAR}[z] = \mathbb{E}[(z - \mathbb{E}[z])^2] = \mathbb{E}[z^2] - \mathbb{E}[z]^2$$

Add WeChat edu_assist_pro

$$= \sum_{i=1}^n \mathbb{E}[(x_i w_i)^2] = \sum_{i=1}^n \mathbb{E}[x_i^2]$$

- So we set $\sigma = 1/\sqrt{m}$ to have “nice” activations.
- **Glorot initialization** takes care to have nice back-propagated signals — see the auto-encoder lab.
- **ReLU** activations $h(x) = \max(x, 0)$ also help in practice.



Motivation

Autoencoder

Assignment Project Exam Help

- if there is no other constraint, then an autoencoder with
 d_{di}
 d_{c}
- Av <https://eduassistpro.github.io>
 - Early stopping of stochastic gradient descent
 - Add noise in the encoding
 - Sparsity constraint on code $a(.)$

Add WeChat edu_assist_pro



Motivation



Autoencoder

Denoising autoencoder

Assignment Project Exam Help

- Add noise to input, keeping perfect example as output
- Aut
- Re

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro
where x noise free, \hat{x} corrupted



Motivation

Autoencoder

Image denoising

- Images with Gaussian noise added.

Assignment Project Exam Help

<https://eduassistpro.github.io>

- Denoised using Stacked Sparse Denoising Autoencoder

Add WeChat [edu_assist_pro](https://eduassistpro.github.io)



Motivation

Autoencoder

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Image from <http://cimg.eu/greycstoration/demonstration.shtml>

<http://cimg.eu/greycstoration/demonstration.shtml>

Undo text over image



Motivation

Autoencoder

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



Motivation

Autoencoder

Recall: Basis functions

Assignment Project Exam Help

- For fixed basis functions $\phi(x)$, we use domain knowledge for encoding features
- Ne $\phi_i(x)$ learn
- The transformations $\phi_i(\cdot)$ for a particular d longer be orthogonal, and furthermore may be variations of each other.
- We collect all the transformed features into a ma

Add WeChat edu_assist_pro



Motivation



Autoencoder

Sparse representations

- Idea: Have many hidden nodes, but only a few active for a particular code $c(x)$
- Student t prior on codes
- ℓ_1 penalty on coefficients α

• <https://eduassistpro.github.io/>

$$\min \sum_{n=1}^N \frac{1}{2} \|x_n - \Phi \alpha_n\|_2^2$$

Add WeChat edu_assist_pro

- Φ is overcomplete, no longer orthogonal
- Sparse \Rightarrow small number of non-zero α_i .
- Exact recovery under certain conditions (coherence):
 $\ell_1 \rightarrow \ell_0$.
- ℓ_1 regulariser \sim Laplace prior $p(\alpha_i) = \frac{\lambda}{2} \exp(-\lambda|\alpha_i|)$.



Motivation



Autoencoder

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat `edu_assist_pro`

$$\underbrace{y}_{\text{measurements}} = \underbrace{x_{\text{orig}}}_{\text{original image}} + \underbrace{\epsilon}_{\text{noise}}$$



Motivation



Autoencoder

Sparsity assumption

- Only have noisy measurements

Assignment Project Exam Help

measurements original image noise

- Giv

<https://eduassistpro.github.io/>

where $\|\cdot\|_0$ is the number of non-zero elements

- Φ is not necessarily features constructed from data.
- Minimise reconstruction error

$$\min_{\alpha} \sum_{n=1}^N \frac{1}{2} \|x_n - \Phi \alpha_n\|_2^2 + \lambda \|\alpha\|_0$$



Motivation

Autoencoder

Convex relaxation

- Want to minimise number of components

Assignment Project Exam Help

$$\min_{\alpha} \sum_{n=1}^N \frac{1}{2} \|x_n - \Phi \alpha_n\|_2^2 + \lambda \|\alpha\|_0$$

- but
- Rel <https://eduassistpro.github.io/>

$$\min_{\alpha} \sum_{n=1}^N \frac{1}{2} \|x_n - \Phi \alpha_n\|_2^2 +$$

Add WeChat edu_assist_pro

where $\|\alpha\|_1 = \sum_n |\alpha_n|$.

- In some settings does minimisation with ℓ_1 regularisation give the same solution as minimisation with ℓ_0 regularisation (exact recovery)?



Motivation



Autoencoder

Mutual coherence

Assignment Project Exam Help

- Expect to be ok when columns of Φ “not too parallel”
- Assume columns of Φ are normalised to unit norm
- Let $K(i, j) = \langle \phi_i, \phi_j \rangle$ be the inner product of the i -th and j -th columns of Φ .
- Define $M = \max_{i \neq j} |K(i, j)|$

$$M = M(\Phi) = \max_{i \neq j} |K(i, j)|$$

- If we have an orthogonal basis, then Φ is an orthogonal matrix, hence $K(i, j) = 0$ when $i \neq j$.
- However, if we have very similar columns, then $M \approx 1$.



Motivation

Autoencoder

Exact recovery conditions

Assignment Project Exam Help

- If a minimiser of the true ℓ_0 problem, α^* satisfies

$$\frac{1}{\| \alpha^* \|_0}$$

- the <https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

then the minimiser of the ℓ_1 relaxation has the same sparsity pattern as α^* .



Motivation

Autoencoder

References

- Yoshua Bengio, "Learning Deep Architectures for AI", Foundations and Trends in Machine Learning, 2009

• ht

• ht
au

<https://eduassistpro.github.io>

- Fuchs, "On Sparse Representations in Arbitrary Redundant Bases", IEEE Trans. Info. Theory

- Xavier Glorot and Yoshua Bengio, "Understanding the difficulty of training deep feedforward neural networks", 2010.

Add WeChat edu_assist_pro