



Outlines

Overview
Introduction
Linear Algebra
Probability
Linear Regression 1
Linear Regression 2
Linear Classification 1
Linear Classification 2
Kernel Methods
Sparse Kernel Methods
Mixture Models and EM 1
Mixture Models and EM 2
Neural Networks 1
Neural Networks 2
Principal Component Analysis
Autoencoders
Graphical Models 1
Graphical Models 2
Graphical Models 3
Sampling
Sequential Data 1
Sequential Data 2

Statistical Machine Learning

Assignment Project Exam Help

Christian Walder

Machine Learning Research Group

CSIRO Data61

and

College of Engineering and Computer Science

The Australian National University

Canberra

Semester One, 2020.

<https://powcoder.com>

Add WeChat powcoder

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



Assignment Project Exam Help

Part X

<https://powcoder.com>

Neural Network 3

Add WeChat powcoder



Motivation

Autoencoder

- 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 fixed number of training examples, expect that compact representations of the target function would yield better generalization

- Example representations

- affine operations, sigmoid \Rightarrow logistic regression has depth 1, fixed number of units (a.k.a. neurons)
- fixed kernel, affine operations \Rightarrow kernel machine (e.g. SVM 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



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

- Example, the d bit parity function

$$\text{parity} : (b_1, \dots, b_d) \in \{0, 1\}^d \mapsto \begin{cases} 1 & \text{if } \sum_{i=1}^d b_i \text{ is even} \\ 0 & \text{otherwise} \end{cases}$$

- **Theorem:** d -bit parity circuits of depth 2 have exponential size

Analogous in modern deep learning:

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

Recall: Multi-layer Neural Network Architecture

Statistical Machine
Learning

© 2020
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University



Motivation

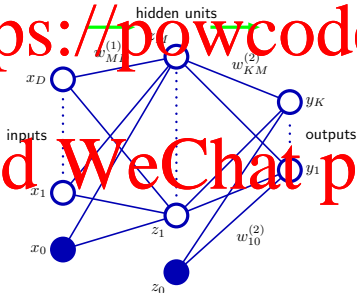
Autoencoder

Assignment Project Exam Help

where \mathbf{w} now contains all weight and bias parameters.

<https://powcoder.com>

Add WeChat powcoder



We could add more hidden layers



Motivation

Autoencoder

Assignment Project Exam Help

- Deep architectures get stuck in local minima or plateaus
- As architecture gets deeper, more difficult to obtain good generalisation
- Hard to initialise random weights well
- 1 or 2 hidden layers seem to perform better
- 2006: Unsupervised pre-training, find distributed representation

<https://powcoder.com>

Add WeChat powcoder



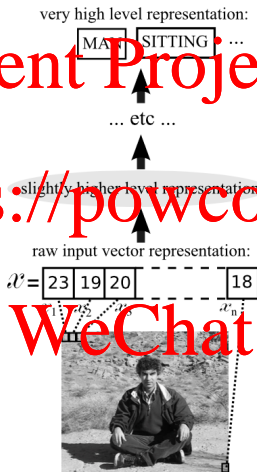
Motivation

Autoencoder

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

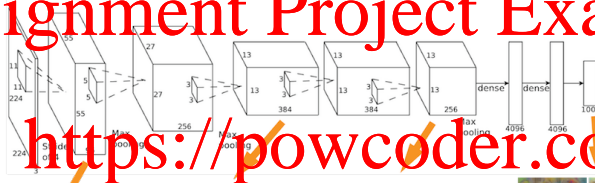




Motivation

Autoencoder

Assignment Project Exam Help



Conv 1: Edge+Blob



Conv 3: Texture



Conv 5: Object Parts



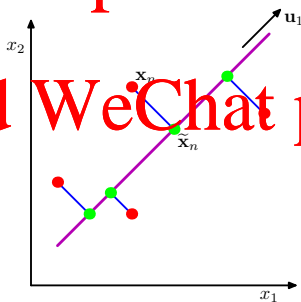
Fc8: Object Classes

<https://powcoder.com>
Add WeChat powcoder

AlexNet / VGG-F network visualized by mNeuron.

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 formulations lead to the same result.
- Need to find the lower dimensional subspace, called the principal subspace.



Multiple PCA layers? - linear transforms



Motivation

Autoencoder

Assignment Project Exam Help

- Principle Component Analysis is a linear transformation (because it is a projection)

- The composite of two linear transformations is linear

- Linear transformations $M: \mathbb{R}^m \rightarrow \mathbb{R}^n$ are matrices

- Let S and T be matrices of appropriate dimension such that ST is defined

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

- Similarly for multiplication with a scalar

⇒ multiple PCA layers pointless

Add WeChat powcoder

Multiple PCA layers? - projection



Motivation

Autoencoder

- 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

- Consider the features formed by projecting onto the principal components

$$Z = XU_k$$

- We perform PCA a second time, $Z^T Z = V \Lambda_Z V^T$.
- By the definition of Z and $X^T X$ and the orthogonality of U

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

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

⇒ again, multiple PCA layers pointless



Motivation

Autoencoder

- 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
- With one linear hidden layer and the mean squared error criterion, the k hidden units learn to project the input in the span of the first k principal components of the data
- If the hidden layer is nonlinear, the autoencoder behaves differently from PCA, with the ability to capture multimodal aspects of the input distribution
- Let f be the **decoder**. We want to minimise the reconstruction error

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



Motivation

Autoencoder

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

<https://powcoder.com>

- If $x|c(x)$ is Gaussian, we recover the familiar squared error
- If the inputs x_i are either binary or considered to be binomial probabilities, then the loss function would be the cross entropy

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

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



Assignment Project Exam Help

- Consider a small number of hidden units.
- $c(x)$ is viewed as a lossy compression of x
- Cannot have small loss for all x , so focus on training examples
- Hope code $c(x)$ is a distributed representation that captures the main factors of variation in the data

<https://powcoder.com>

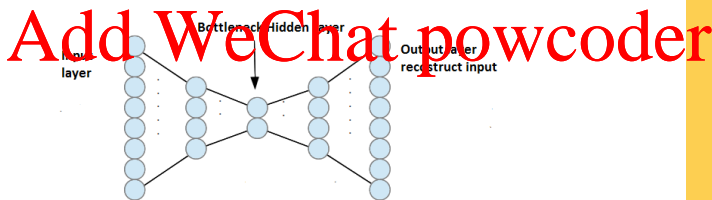
Add WeChat powcoder



Motivation

Autoencoder

- 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 example, let $z_1 = c_1(x)$, and $z_2 = c_2(z_1)$, the corresponding decoding is given by $f_2(f_1(z_2))$
- Because of non-linear activation functions, the latent feature z_2 can capture more complex patterns than z_1 .



Higher level image features - faces

Statistical Machine
Learning

© 2020
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University



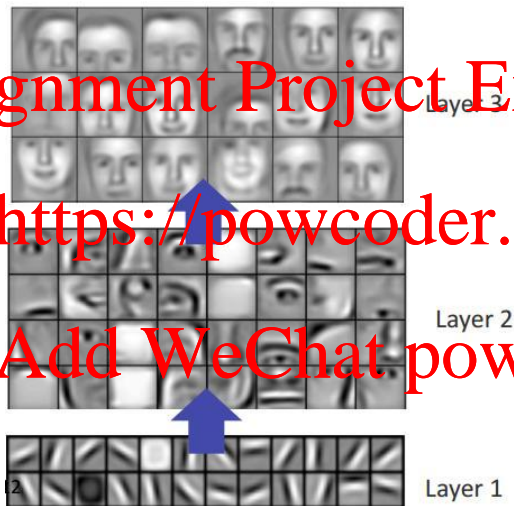
Motivation

Autoencoder

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



codingplayground.blogspot.com



Motivation

Autoencoder

Assignment Project Exam Help

- Latent features z_j in layer i can capture high level patterns

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

- These features may also be useful for supervised learning tasks.
- In contrast to the feed forward network, the features z_j are constructed in an unsupervised fashion.
- Discard the decoding layers, and directly use z_j with a supervised training method, such as logistic regression.
- Various such pre-trained networks are available on-line, e.g [VGG-19](https://openai.org/research/vgg19).

<https://powcoder.com>

Add WeChat powcoder



- 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 the inputs $x_i \sim \mathcal{N}(0, 1)$, weights $w_i \sim \mathcal{N}(0, \sigma^2)$ and activation $z = \sum_{i=1}^m x_i w_i$. Then

$$\text{VAR}[z] = \mathbb{E}[(z - \mathbb{E}[z])^2] = \mathbb{E}[z^2] = \mathbb{E}\left[\left(\sum_{i=1}^m x_i w_i\right)^2\right]$$

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

- 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 -dimensional input and an encoding of dimension at least d could potentially just learn the identity function

- Avoid by:

- Regularisation
- Early stopping of stochastic gradient descent
- Add noise in the encoding
- Sparsity constraint on code z

<https://powcoder.com>

Add WeChat powcoder



Motivation

Autoencoder

Assignment Project Exam Help

- Add noise to input, keeping perfect example as output
- Autoencoder tries to:
 - 1. preserve information of input
 - 2. undo stochastic corruption process
- Reconstruction log likelihood

<https://powcoder.com>

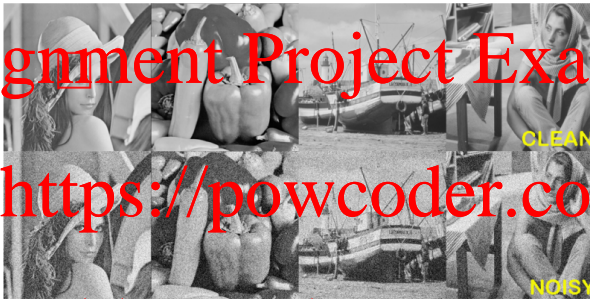
$-\log P(x|c(\hat{x}))$
Add WeChat powcoder
where x noise free, \hat{x} corrupted



Motivation

Autoencoder

- Images with Gaussian noise added.



- Denoised using Stacked Sparse Denoising Autoencoder



Images from Xie et. al. NIPS 2012



Motivation

Autoencoder

Assignment Project Exam Help

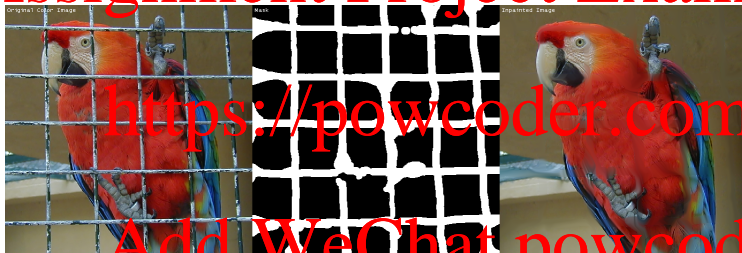


Image from [http:](http://cimng.eu/greycstoration/demonstration.shtml)

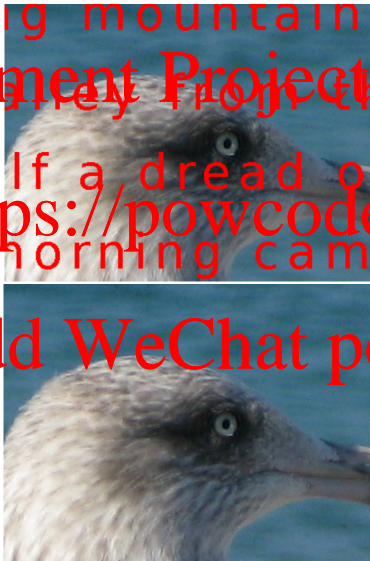
[//cimng.eu/greycstoration/demonstration.shtml](http://cimng.eu/greycstoration/demonstration.shtml)

Undo text over image

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder





Motivation

Autoencoder

Assignment Project Exam Help

- For fixed basis functions $\phi(x)$, we use domain knowledge for encoding features
- Neural networks use data to learn a set of transformations. $\phi_i(x)$ is the i^{th} component of the feature vector, and is learned by the network.
- The transformations $\phi_i(\cdot)$ for a particular dataset may no longer be orthogonal, and furthermore may be minor variations of each other.
- We collect all the transformed features into a matrix Φ .

<https://powcoder.com>
Add WeChat powcoder



Motivation

Autoencoder

- 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

- Given bases in matrix Φ , look for codes by choosing α such that input signal x is reconstructed with low ℓ_2 reconstruction error, while w is sparse

$$\min_{\alpha} \sum_{n=1}^N \frac{1}{2} \|x_n - \Phi \alpha_n\|_2^2 + \lambda \|\alpha\|_1$$

- 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 $p(\alpha_i) = \frac{\lambda}{2} \exp(-\lambda|\alpha_i|)$.

The image denoising problem

Statistical Machine
Learning

© 2020
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University



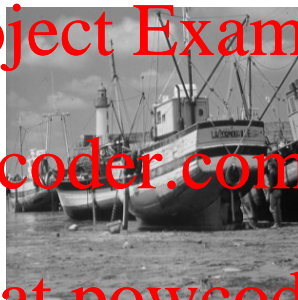
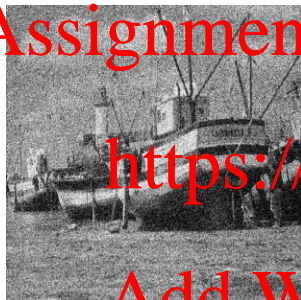
Motivation

Autoencoder

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



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



- Only have noisy measurements

$y = x_{\text{orig}} + w$
measurements original image noise

- Given $\Phi \in \mathbb{R}^{m \times p}$, find α such that

$\|\alpha\|_0$ is small for $x \approx \Phi\alpha$
<https://powcoder.com>

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

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

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

Add WeChat powcoder



Motivation

Autoencoder

- Want to minimise number of components

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

but $\|\cdot\|_0$ is hard to optimise

- Relax to a convex norm

$$\min_{\alpha} \sum_{n=1}^N \frac{1}{2} \|x_n - \Phi \alpha_n\|_2^2 + \lambda \|\alpha\|_1$$

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)?

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Motivation

Autoencoder

Assignment Project Exam Help

- Expect to be asked when columns of Φ “not too parallel”
- Assume columns of Φ are normalised to unit norm
- Let $K = \Phi\Phi^T$ be the Gram matrix, then $K(i,j)$ is the value of the inner product between ϕ_i and ϕ_j .
- Define the **mutual coherence**

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

- 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 $M \approx 1$.

Add WeChat powcoder



Assignment Project Exam Help

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

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

then it is the unique sparsest solution.

- If α^* satisfies the stronger condition

$$\|\alpha^*\|_0 < \frac{1}{2} \left(1 + \frac{1}{M} \right)$$

Add WeChat powcoder

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



Motivation

Autoencoder

Assignment Project Exam Help

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

- <http://deeplearning.net/tutorial/>

- <http://www.deeplearningbook.org/contents/autoencoders.html>

- Fuchs, "On Sparse Representations in Arbitrary Redundant Bases", IEEE Trans. Info. Theory, 2004

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

<https://powcoder.com>

Add WeChat powcoder