

DBN's : IntroductionBoltzmann Machine:

$$P(\underline{v}, \underline{h}) = \frac{1}{Z} e^{-E(\underline{v}, \underline{h})}$$

observed nodes \nearrow \nearrow hidden nodes.

$$E(\underline{v}, \underline{h}) = -\sum_{ij} v_i w_{ij} h_j - \sum_i c_i v_i - \sum_{ij} h_i w_{ij}^h h_j$$

motivated by modeling the brain, but just an MRF.

$$v_i \in \{0, 1\}$$

$$h_i \in \{0, 1\}$$

obvious example. $\{\underline{v}^M : \mu = 1 \text{ to } N\}$

$$\underline{v} = (v_1, \dots, v_n)$$

states of observed nodes.

Task: learn the weights of

$$\text{the model} - \{w_{ij}, c_i, w_{ij}^h\}$$

can be learnt, in principle, by the EM algorithm
or stochastic variant (e.g. data augmentation)

Difficult, because it is hard to do inference
on this type of model in general.

Why learn this model? It learns a generative
model for data. It learns "receptive" fields

- e.g. the w_{ij} , how a hidden unit " h_i " is activated
by the inputs.

To build a DBN, you first consider a simplified
"Restricted Boltzmann Machine" RBM.

The connection terms w_{ij}^h between the hidden units
are removed. \rightarrow replace $-\sum_{ij} h_i w_{ij}^h h_j$ by $-\sum_i b_i h_i$.

This has the property that $P(\underline{h}|\underline{v})$ and $P(\underline{v}|\underline{h})$ have
simple forms:

$$P(\underline{h}|\underline{v}) = \prod_i P(h_i|\underline{v})$$

$$P(h_i|\underline{v}) = e^{-h_i(\sum_j v_j w_{ij} + b_i)}$$

similarly

$$P(\underline{v}|\underline{h}) = \prod_i P(v_i|\underline{h})$$

similar form for $P(v_i|\underline{h})$

$$P(v_i|\underline{h}) = e^{-v_i(\sum_j h_j w_{ij} + b_i)}$$

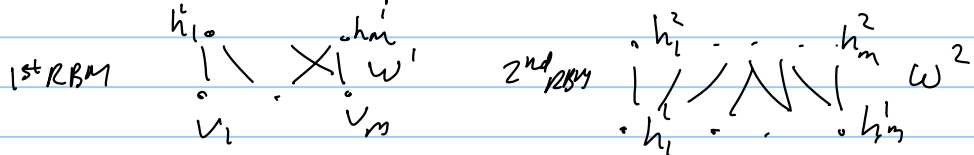
(2) Key point RBM's.

Easy to sample from $P(\underline{v}|\underline{h})$ i.e. sample from

Easy to sample from $P(\underline{h}|\underline{v})$. $P(\underline{v}|\underline{h})$ independent

This enables us to do inference (stochastically) and makes it practical to learn the parameter $\{W\}$ of an RBM.

But RBM's are too simple \rightarrow so add another RBM



use the hidden states $(h_1^1 \dots h_m^1)$ of the 1st RBM as input to the 2nd RBM.

Then add another RBM and so on n^{th} order RBM.

Train the model layer by layer.

Intuition: the 'receptive fields'

the w^1 gives hidden units (h^1) which capture image features - like a dictionary in the previous two lectures.

\rightarrow then the receptive fields w^2, w^3, \dots of the higher levels learn better dictionaries

The hidden states at the top levels can be used as inputs to a classifier level.

This can be trained in

supervised mode - i.e. the $\{o_i^l\}$ and $\{v_i^l\}$ are provided by the training data.

Effective for learning how to recognize digits - Hinton et al.

For other variants see LeCun, Ng, Bengio.

Hinton's Talk: 4 March (Thurs) 4:15 p.m. CS.

(3)

Active Appearance Models

AAM's

Object



Variations in geometry.
Variations in appearance.

First, suppose objects are aligned (i.e. remove spatial variation).

Perform PCA principle component analysis.

$$\{I^m(x) : m \in \Lambda\}$$

$$\bar{I} = \frac{1}{|\Lambda|} \sum_{m \in \Lambda} I^m(x)$$

$$K(x, y) = \frac{1}{|\Lambda|} \sum_{m \in \Lambda} \langle I^m(x) - \bar{I}(x) \rangle \langle I^m(y) - \bar{I}(y) \rangle$$

$$\text{solve } \sum_y K(x, y) B^m(y) = \lambda^m B^m(x)$$

solve for eigenvalues & eigenvectors.

represent objects appearance.

$$I(x) = \bar{I}(x) + \sum_{a=1}^m \alpha_a B_a(x) + \epsilon(x)$$

efficient representation if

coefficient basis functions.

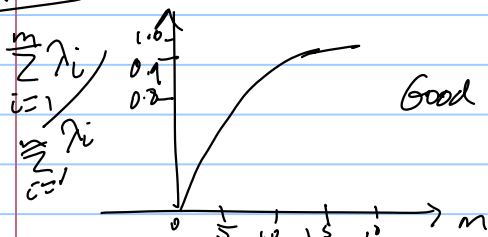
we can approximate $I(x)$ using a small no. basis functions i.e. m small

$$\text{Square Error } \frac{1}{|\Lambda|} \sum_{m \in \Lambda} |I(x) - \bar{I}(x) - \sum_{a=1}^m \alpha_a B_a(x)|^2$$

$$= \frac{\lambda^1 + \dots + \lambda^m}{\sum_{i=1}^m \lambda_i}$$

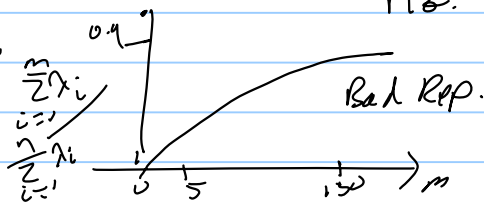
all the eigenvalues.

For faces



Good Representation

For text - eg. ABC
GOT
PIG.

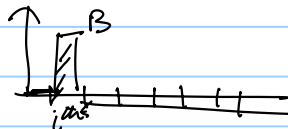


Poor Rep.

Intuition: PCA assumes that the data lies in a linear space. This may be true (approximately) for faces, but certainly is not true for text.

Example:

data



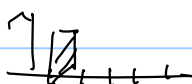
N -states

$$I(x) = B \delta(x, x_i)$$

i th bin

$$\sum_{i=1}^N \lambda_i / \sum_{i=1}^N \lambda_i$$

Suppose the data is



Then PCA gives $(N-1)$ principal components, all with the same small eigenvalues.

(4)

Second: fix appearance, model geometry

$$I(x) \rightarrow I(s(x))$$

where $s(x)$ is a spatial warp

Also represent the spatial warp by PCA also.

i.e. take training data $\{s^\mu(x) : \mu \in \Lambda\}$

performs PCA on $\{s^\mu\}$

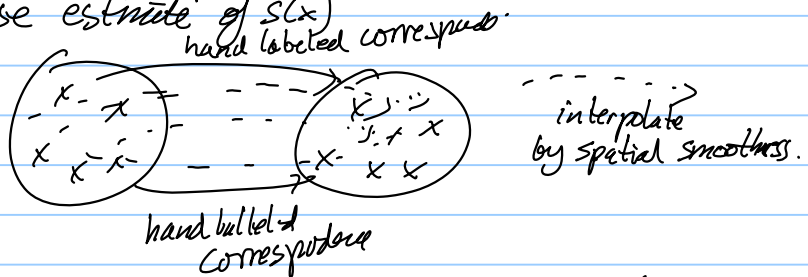
express:
$$s(x) = \sum_{b=1}^q \beta_b e_b(x) + \epsilon(x)$$

where q is no of coefficients.
 $e_b(x)$ are the eigenvectors.

How to get the training data?

Typically handlabeled.

Can label a set of sparse points and interpolate to get dense estimate of $s(x)$



In principle, can learn the spatial and basis functions by EM.

Generative Model $P(I | \{e, \beta\}, \{\alpha, \beta\}) = \frac{1}{Z} e^{-E(I; \{e, \beta\}, \{\alpha, \beta\})}$

$$E(I; \{e, \beta\}, \{\alpha, \beta\}) = \sum_x \left(I(x) - \sum_a \alpha_a \beta_a \left(\sum_b \beta_b e_b(x) \right) \right)^2$$

set of images $\{I^\mu : \mu \in \Lambda\}$

parameter $\{\alpha^\mu, \beta^\mu\}$ - different coefficients for each image.

$$P(\langle I^\mu \rangle | \{e, \beta\}, \{\alpha^\mu, \beta^\mu\})$$

$$= \prod_{\mu} P(I^\mu | \{e, \beta\}, \alpha^\mu, \beta^\mu)$$

EM

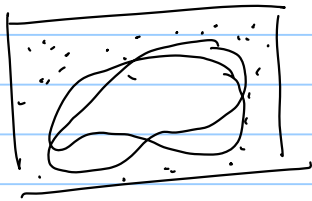
M-step estimate $\{e, \beta\}$

E-step estimate probabilities $q(\alpha^\mu, \beta^\mu)$ for the coefficients

Problem \rightarrow many local minima \rightarrow

(5)

Recent work: eg. J. Kolkinos.



Start with objects only
approximately located.
→ e.g. box round object

Background clutter

To simplify the problem and make it tractable:

— filter the images to detect edges and other features like ridges — this remove the appearance variation so need to estimate the β 's and α 's

Simply need to estimate the spatial warps and their coefficients:

$$I^*(x) = T \left(\sum_{b=1}^m \beta_b^m S_b(x) \right) + \epsilon(x)$$

filled x Task: estimate $T(\cdot)$ the edge/ridge image of the object.

Still require EM, E-step to estimate T, β^m
M-step to estimate $\{S_b\}$

EM does converge (Kolkinos & Yuille) with reasonable initial conditions.

Complication: what is m ? how many basis functions?

Strategy (KEY) → greedy search.

Assume $m=1$, learn $S_1(x)$

then set $m=2$, learn $S_2(x)$...

Note: there were limitations using PCA to represent the appearance (i.e. approx. linear assumption).

what are the limitations of PCA to represent spatial warps?

Roughly true for faces, torso/body of a cow.

But bad for representing the spatial variability of the legs of a cow

Better to treat the legs as separate parts connected to the torso/body

