

Mixture Density Networks

Sargur Srihari

Mixture Density Networks

- Goal of supervised learning is to model the conditional distribution $p(t|x)$
- In some problems distribution can be multimodal
 - Particularly in inverse problems
- Gaussian assumption can lead to poor results
 - In regression $p(t|x)$ is typically assumed to be Gaussian
 - i.e., $p(t|x) = N(t|y(x, w), b^{-1})$

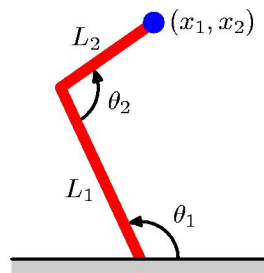
Kinematics of a robot arm

- Robot arm with two links

Forward problem:

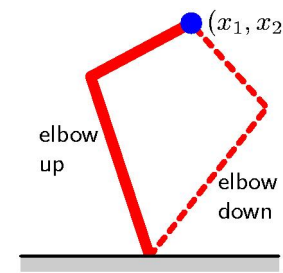
Find end effector position given joint angles

Has a unique solution



Inverse kinematics has two solutions:

Elbow-up and elbow-down



- Inverse problem is a regression problem with
 - two inputs:
 - desired location of arm (x_1, x_2)
 - two outputs:
 - angles for links (θ_1, θ_2)
 - Has two solutions (elbow up, elbow down)

Forward and Inverse Problems

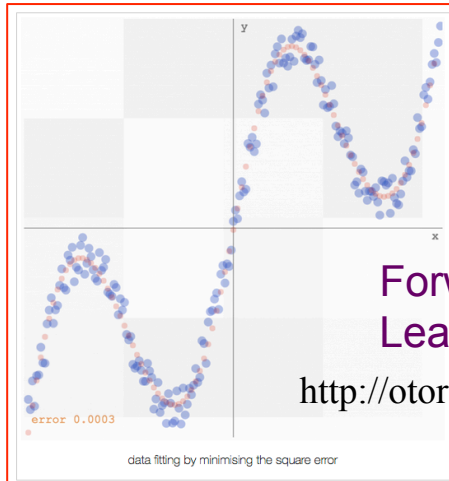
- Forward problems correspond to causality in a physical system
 - Have a unique solution
 - A disease causes a set of symptoms
- If forward problem is a many-to-one mapping,
 - Several diseases have the same symptoms
 - Inverse has multiple solutions
 - Same symptoms caused by several diseases

An Example:

x is disease, y is symptom;

x is elbow position, y is effector position

$$y = x + 0.3\sin(4\pi x) + \varepsilon$$



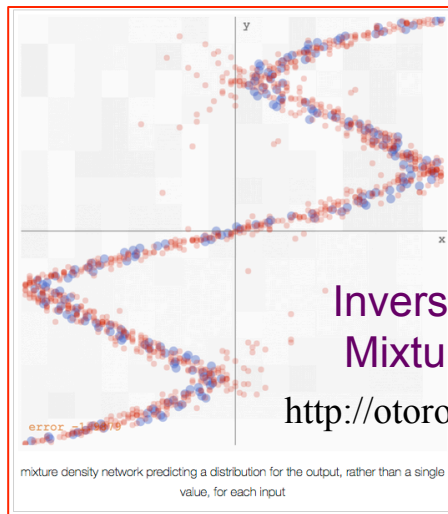
Forward Data and
Least squares regression

<http://otoro.net/ml/mixture/index.html>



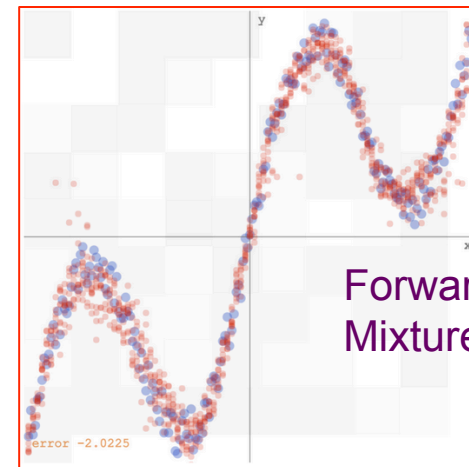
Inverse Data and
Least squares
regression

<http://otoro.net/ml/mixture/inverse.html>



Inverse Data and
Mixture model

<http://otoro.net/ml/mixture/mixture.html>



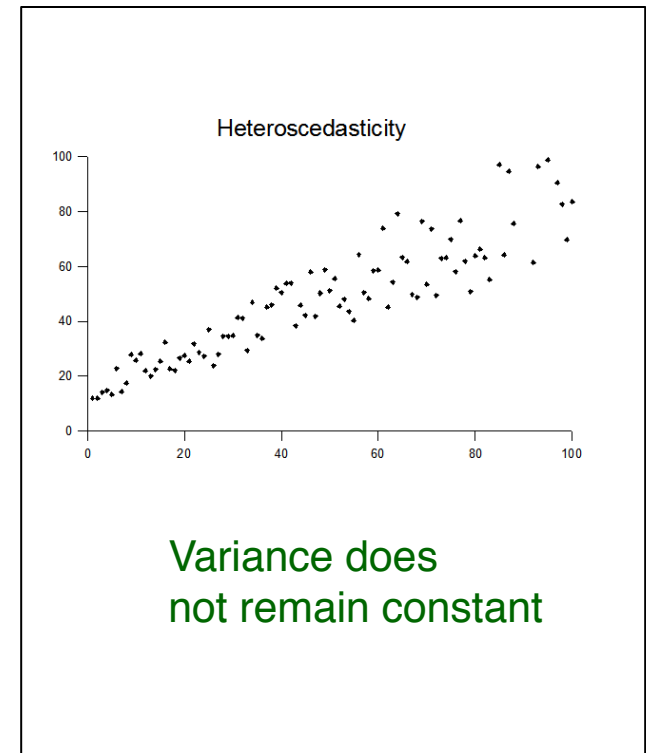
Forward Data and
Mixture model

A Mixture Density

- Generative model with K components

$$p(t | \mathbf{x}) = \sum_{k=1}^K \pi_k(\mathbf{x}) \mathcal{N}(t | \mu_k(\mathbf{x}), \sigma_k^2(\mathbf{x}))$$

- Components can be Gaussian for continuous variables, Bernoulli for binary target variables, etc
- Note that mixing coefficient π is dependent on \mathbf{x} and sums to 1 for each \mathbf{x}
- So also mean and variance
 - An example of a hetero-scedastic model since variance is a function of the input vector \mathbf{x}



Data Set for Forward and Inverse Problems

- Least squares corresponds to
 - Maximum likelihood under a Gaussian assumption
 - Leads to a poor result for highly non-Gaussian inverse problem
- Seek a general framework for modeling conditional probability distributions
- Achieved by using a mixture model $p(t|x)$

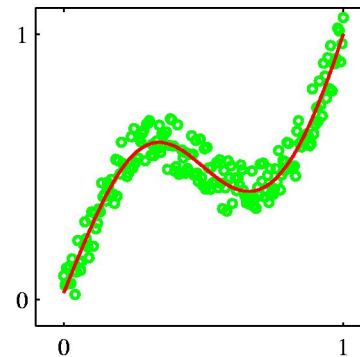
Forward problem data set:

x is sampled uniformly over $(0,1)$ to give values $\{x_n\}$

Target t_n obtained by function

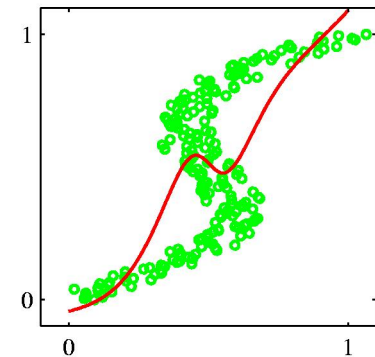
$$x_n + 0.3 \sin(2\pi x_n)$$

Then add noise over $(-0.1, 0.1)$



Red curve is result of fitting a two-layer neural network by minimizing sum-of-squared error

Corresponding inverse problem by reversing x and t



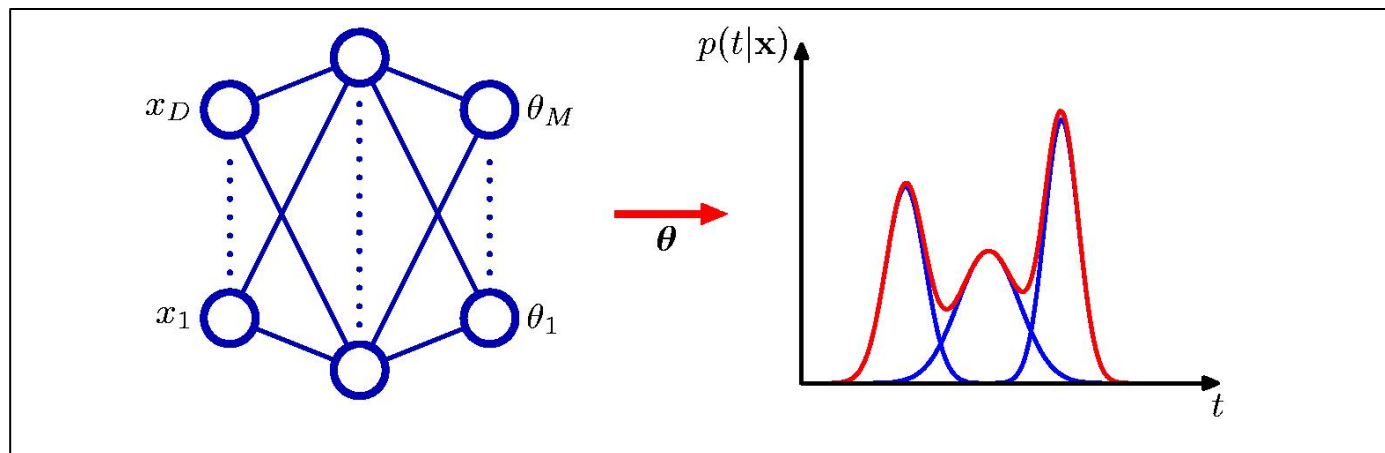
Very poor fit to data

Parameters of Mixture Model

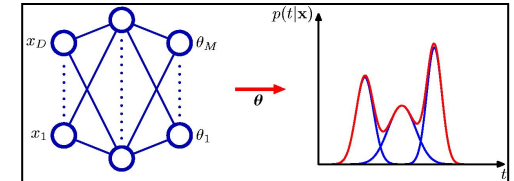
- Parameters of the mixture density:
 1. Mixing coefficients $\pi_k(\mathbf{x})$
 2. Means $\mu_k(\mathbf{x})$
 3. Variances $\sigma_k^2(\mathbf{x})$
- Governed by the outputs of a neural network
 - With \mathbf{x} as input
- A single network predicts the parameters of all the component densities

Mixture density network

- Network represents general conditional probability densities $p(t|\mathbf{x})$ by considering a parametric mixture model
- It takes \mathbf{x} as input and provides the parameters of the distribution as output
 - In effect, it specifies the distribution of $p(t|\mathbf{x})$
 - Takes \mathbf{x} as input vector



No. of Outputs of Neural Network



- Two-layer network with sigmoidal (\tanh) hidden units
- No. of output units calculated as follows:
 - If K components in mixture model then there are K mixing coefficients $\pi_k(x)$ determined by activations a_k^π
 - K outputs a_k^σ that determine kernel widths $\sigma_k(x)$
 - If Target t has L components then there are $K \times L$ outputs denoted a_{kj}^μ that determine components $\mu_{kj}(x)$ of kernel centres $\mu_k(x)$
- Then network will have $(L+2)K$ outputs
 - Instead of usual L outputs of a network, which simply predict the conditional means of target variables

Outputs of Mixture Density Network

1. Mixing coefficients

- must satisfy $\sum_{k=1}^K \pi_k(\mathbf{x}) = 1, \quad 0 \leq \pi_k(\mathbf{x}) \leq 1$
- Achieved using softmax outputs

$$\pi_k(\mathbf{x}) = \frac{\exp(a_k^\pi)}{\sum_{l=1}^K \exp(a_l^\pi)}$$

2. Variances

- Must satisfy $\sigma_k^2(\mathbf{x}) \geq 0$
- Represented as exponentials of activations $\sigma_k(\mathbf{x}) = \exp(a_k^\sigma)$

3. Means

- Real components represented directly by output activations

$$\mu_{kj}(\mathbf{x}) = a_{kj}^\mu$$

Error Function for Mixture Density Network

- Can be set by maximum likelihood
- From distribution

$$p(t | x) = \sum_{k=1}^K \pi_k(x) N(t | \mu_k(x), \sigma_k^2(x))$$

This is what will be used to predict the output for a given input

- Negative logarithm of Likelihood function is

$$E(w) = - \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k(x_n, w) N(t | \mu_k(x_n, w), \sigma_k^2(x_n, w)) \right\}$$

Minimization of Error Function

- Need to calculate derivatives of error $E(w)$ wrt component
- Can be evaluated provided we find suitable expressions for derivatives of error wrt output unit activations
 - They represent error signals for each pattern and output unit
- Derivative terms are easily obtained due to summation of terms one for each data point

View of mixing coefficients

- Convenient to view mixing coefficients $\pi_k(\mathbf{x})$ as \mathbf{x} -dependent prior probabilities
- Corresponding posterior probabilities are

$$\gamma_k(\mathbf{t} | \mathbf{x}) = \frac{\pi_k N_{nk}}{\sum_{l=1}^K \pi_l N_{nl}}$$

- where N_{nk} denotes $N(\mathbf{t}_n | \mu_k(\mathbf{x}_n), \sigma_k^2(\mathbf{x}_n))$

Derivatives with respect to Network Output Activations

1. Mixing coefficients

$$\frac{\partial E_n}{\partial a_k^\pi} = \pi_k - \gamma_k$$

2. Component means

$$\frac{\partial E_n}{\partial a_{kl}^\mu} = \gamma_k \left\{ \frac{\mu_{kl} - t_l}{\sigma_k^2} \right\}$$

3. Component variances

$$\frac{\partial E_n}{\partial a_{kl}^\sigma} = -\gamma_k \left\{ \frac{\|t - \mu_k\|^2}{\sigma_k^2} - \frac{1}{\sigma_k} \right\}$$

Training Data for Mixture Density Network

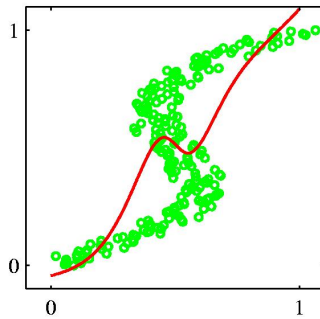
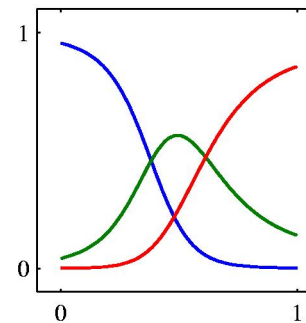
- Obtained easily from forward data by exchanging roles of x and t
 - For different joint angles, the position of end effector



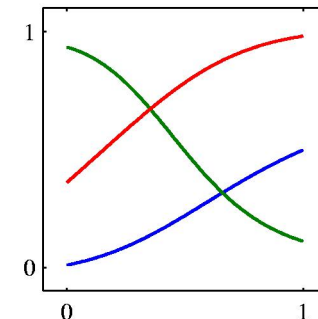
- Note that we are using x as input and t as output after data exchange

Output of Mixture Density Network

Data Set

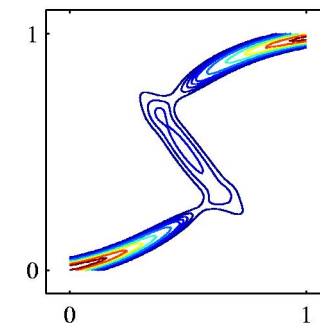
Mixing Coefficients
 $\pi_k(x)$ versus x 

(a)

Means $\mu_k(x)$ 

(b)

Three components have to
sum to Unity

Contours of
Conditional
probability density
of target data

(c)

While the outputs of the neural network
(and hence the parameters) are necessarily single valued,
The model is able to produce a conditional density that is
unimodal for some values of x and trimodal for other values

Use of Mixture Density Network

- Once mixture density network has been trained
 - can predict conditional density function of the target data for given value of input vector
- From this density can calculate more specific quantities of interest in applications
 - e.g., mean of the target data

Predicting value of output vector

- Conditional distribution represents complete description of generator of data
- From this density we can calculate the mean
 - Which is the conditional average of target data

$$E[t | x] = \int t p(t | x) dt = \sum_{k=1}^K \pi_k(x) \mu_k(x)$$

This is expected value,
Not Error!

- This is same as least squared solution and is limited value
 - Average of two solutions is not a solution
- Variance of density function about the conditional average

$$s^2(x) = E[\|t - E[t | x]\|^2 | x]$$

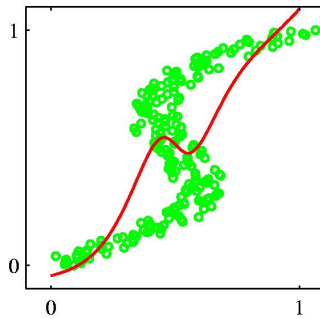
$$= \sum_{k=1}^K \pi_k(x) \left\{ \sigma_k^2(x) + \left\| \mu_k(x) - \sum_{l=1}^K \pi_l(x) \mu_l(x) \right\|^2 \right\}$$

Mode as the solution

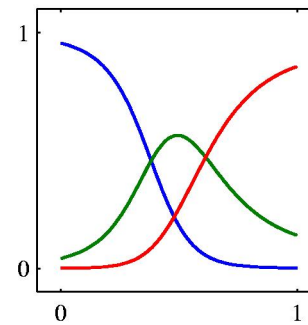
- Each of the modes of mixture density is a better solution than the single mean
- Does not have a simple analytical solution
 - Need numerical iteration
- Simple alternative:
 - Take mean of the most probable component
 - One with largest mixing coefficient for each value of x

Example of Mixture Density Network

Data Set

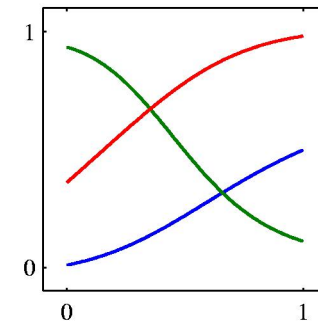


Mixing Coefficients
 $\pi_k(x)$ versus x



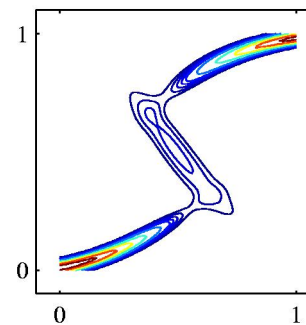
(a)

Means $\mu_k(x)$



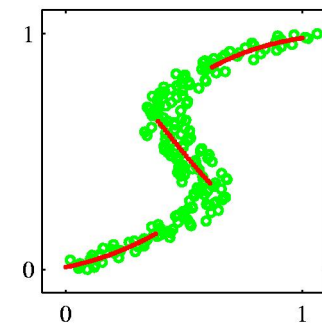
(b)

Contours of
Conditional
Probability density



(c)

Approximate
Conditional mode
(red points of
Conditional density)



(d)