COMS30035, Machine learning: Sequential Data 1: Markov Models

Edwin Simpson edwin.simpson@bristol.ac.uk

Department of Computer Science, SCEEM University of Bristol

November 6, 2023

Agenda

Markov Models

- Hidden Markov Models
- EM for HMMs
- Linear Dynamical Systems

Textbook

We will follow Chapter 13 of the Bishop book: Bishop, C. M., Pattern recognition and machine learning (2006). Available for free <u>here</u>.

Up to now, we have considered the data points in our datasets to be independent and identically distributed (i.i.d.)

- Up to now, we have considered the data points in our datasets to be independent and identically distributed (i.i.d.)
- Independent: the value of one data point does not affect the others, p(x1, x2) = p(x1)p(x2)

- Up to now, we have considered the data points in our datasets to be independent and identically distributed (i.i.d.)
- Independent: the value of one data point does not affect the others, p(x1, x2) = p(x1)p(x2)
- Identically distributed: all data points have the same distribution, p(x_i) = p(x_j), ∀i, ∀j

- So, once you have trained a classifier or regressor, you can predict the output for each data point independently.
- Can you think of situations where the i.i.d. assumption does not apply?

Sequential Data

- The i.i.d. assumption ignores any ordering of the data points.
- Data points often occur in a sequence, such as words in a sentence, frames in a video, sensor observations over time, stock prices...
- This can be generalised to more than one dimension: object in different parts of an image, geographical data on a map... (not covered in this lecture).

Sequential Data

- The i.i.d. assumption ignores any ordering of the data points.
- Data points often occur in a sequence, such as words in a sentence, frames in a video, sensor observations over time, stock prices...
- This can be generalised to more than one dimension: object in different parts of an image, geographical data on a map... (not covered in this lecture).
- Can you think of some classification or regression tasks for these types of data?

How have we modelled relationships between data points so far?

How have we modelled relationships between data points so far? – Through their input features.

- How have we modelled relationships between data points so far? Through their input features.
- Can we model sequential relationships by simply making *time* or *position in the sequence* into another feature?

- How have we modelled relationships between data points so far? Through their input features.
- Can we model sequential relationships by simply making *time* or *position in the sequence* into another feature?
- No The timestamp or positional index is not in itself an informative feature
- But the data observed at other points in the sequence tells us about our current data point

- Look at the following two texts from Bishop's book, both with a missing word:
 - "later termed Bayes' ____ by Poincarré"
 - "The evaluation of this conditional can be seen as an example of Bayes"

- Look at the following two texts from Bishop's book, both with a missing word:
 - "later termed Bayes' ____ by Poincarré"
 - "The evaluation of this conditional can be seen as an example of Bayes"
- Can you guess the missing words? How did you guess them?

- Look at the following two texts from Bishop's book, both with a missing word:
 - "later termed Bayes' ____ by Poincarré"
 - "The evaluation of this conditional can be seen as an example of Bayes"
- Can you guess the missing words? How did you guess them?
- You can guess that the missing word in both cases is "theorem" or maybe "rule", because of the word "Bayes" right before it.
- The first missing word is at position 3, the second is at position 13, but these position indexes don't help to identify the missing word.

How Can We Model the Dependencies?

i.i.d.,
$$p(\boldsymbol{x}_n | \boldsymbol{x}_1, ..., \boldsymbol{x}_{n-1}) = p(\boldsymbol{x}_n)$$

 $(\boldsymbol{x}_1) \quad (\boldsymbol{x}_2) \quad (\boldsymbol{x}_3) \quad (\boldsymbol{x}_4)$

How Can We Model the Dependencies?

i.i.d.,
$$p(\boldsymbol{x}_n | \boldsymbol{x}_1, ..., \boldsymbol{x}_{n-1}) = p(\boldsymbol{x}_n)$$

 $(\boldsymbol{x}_1) \quad (\boldsymbol{x}_2) \quad (\boldsymbol{x}_3) \quad (\boldsymbol{x}_4)$

Modelling all connections, $p(\mathbf{x}_n | \mathbf{x}_1, ..., \mathbf{x}_{n-1}) - intractable$



How Can We Model the Dependencies?

i.i.d.,
$$p(\boldsymbol{x}_n | \boldsymbol{x}_1, ..., \boldsymbol{x}_{n-1}) = p(\boldsymbol{x}_n)$$

 $(\boldsymbol{x}_1) \quad (\boldsymbol{x}_2) \quad (\boldsymbol{x}_3) \quad (\boldsymbol{x}_4)$

Modelling all connections, $p(\mathbf{x}_n | \mathbf{x}_1, ..., \mathbf{x}_{n-1}) - intractable$



1st order Markov chain, $p(x_n | x_1, ..., x_{n-1}) = p(x_n | x_{n-1})$

$$(X_1) \rightarrow (X_2) \rightarrow (X_3) \rightarrow (X_4)$$

$$p(\mathbf{x}_1, ..., \mathbf{x}_N) = p(\mathbf{x}_1) \prod_{n=2}^{N} p(\mathbf{x}_n | \mathbf{x}_{n-1})$$

Homogeneous Markov Chains

- Stationary distribution: the probability distribution remains the same over time.
- This leads to a homogeneous Markov chain.
- E.g., the parameters of the distribution remain the same while the data evolves.
- Contrast with non-stationary distributions that change over time.

Higher-Order Markov Models

- Sometimes it is necessary to consider earlier observations using a higher-order chain.
- However, the number of parameters increases with the order of the Markov chain, meaning higher-order models are often impractical.

1st order Markov chain, $p(\boldsymbol{x}_n | \boldsymbol{x}_1, ..., \boldsymbol{x}_{n-1}) = p(\boldsymbol{x}_n | \boldsymbol{x}_{n-1})$

$$(\mathbf{X}_1) \rightarrow (\mathbf{X}_2) \rightarrow (\mathbf{X}_3) \rightarrow (\mathbf{X}_4)$$

2nd order Markov chain, $p(x_n | x_1, ..., x_{n-1}) = p(x_n | x_{n-1}, x_{n-2})$

$$(X_1 \rightarrow (X_2 \rightarrow (X_3 \rightarrow (X_4)$$

- What if we don't directly observe the states we want to model?
- E.g., we want to predict the state of the weather (raining, sunny, cloudy, rainfall)
- We observe noisy measurements of temperature, wind, rainfall over a period of time

- What if we don't directly observe the states we want to model?
- E.g., we want to identify different actions in a video of a game of tennis, such as backhand volley
- We observe the frames in a video, each one of which is a tensor of pixel values

- What if we don't directly observe the states we want to model?
- E.g., we want to identify different actions in a video of a game of tennis, such as backhand volley
- We observe the frames in a video, each one of which is a tensor of pixel values
- We encounter the same problem as we do in i.i.d. classification and regression: the sequential variable we wish to predict is not directly observed.

- lntroduce latent variables, \boldsymbol{z}_n that form a Markov chain;
- Each observation \boldsymbol{x}_n depends on \boldsymbol{z}_n ;
- This means we do not need to model the dependencies between observations x_n directly;
- Latent variables model the state of the system, while observations may be of different types, contain noise...



Does this look similar to any classifiers you have come across before?



- ► Hidden Markov Models (HMMs): Discrete state z, observations may be continuous or discrete according to any distribution. → next part of this lecture
- ► Linear Dynamical Systems (LDS): Continuous state z, observations are continuous, both have Gaussian distributions → after reading week
- ▶ We will consider both supervised and unsupervised settings.

