

# Assorted Topics (2)

CS772A: Probabilistic Machine Learning

Piyush Rai

# Plan today

- Some classical probabilistic models for sequential data
  - Hidden Markov Models (HMM) and State-Space Models (SSM)
- Another non-Bayesian way to get uncertainty estimates:
  - Conformal Prediction
- Simulation based inference

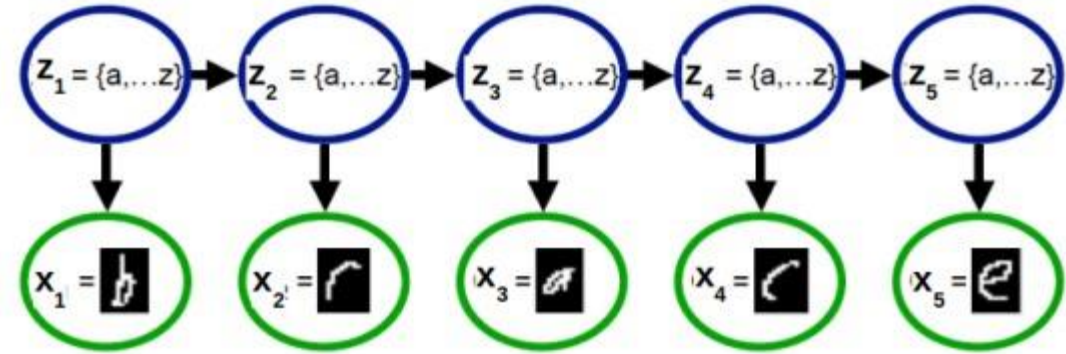
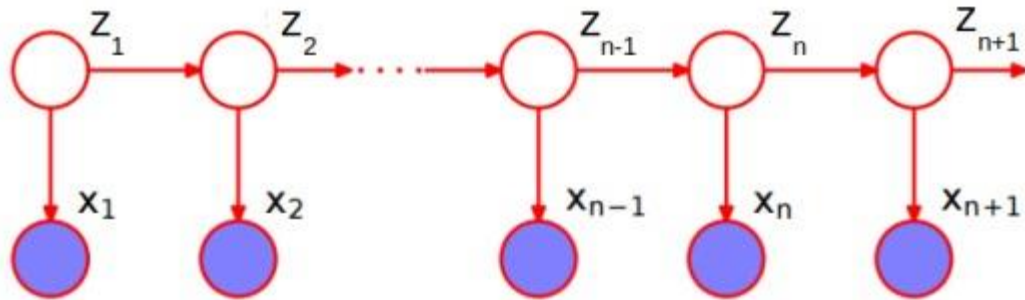


# Probabilistic Models for Sequential Data



# Latent Variable Models for Sequential Data

- Task: Given a sequence of observations, infer the latent state of each observation



Observation  
model

$$\mathbf{x}_n | \mathbf{z}_n \sim p(\mathbf{x}_n | \mathbf{z}_n)$$

(i.i.d. draws of  $\mathbf{x}_n$  given  $\mathbf{z}_n$ )

State-transition  
model

$$\mathbf{z}_n | \mathbf{z}_{n-1} \sim p(\mathbf{z}_n | \mathbf{z}_{n-1})$$

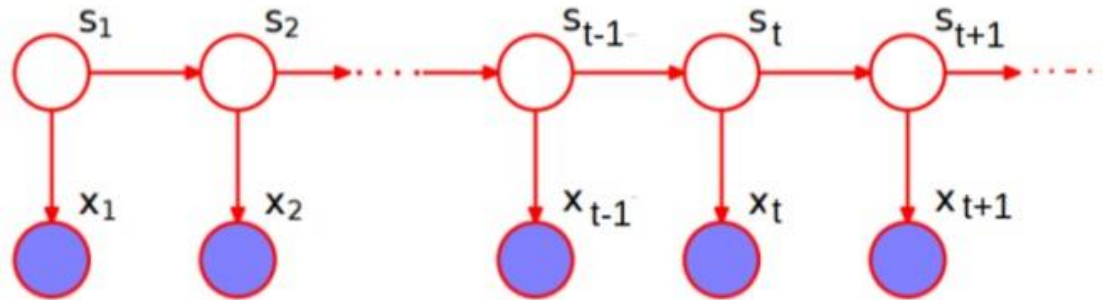
(first-order dependence b/w  $\mathbf{z}_n$ 's)

- If  $\mathbf{z}_n$ 's are discrete, we have a hidden **Markov model (HMM)**  $p(\mathbf{z}_n | \mathbf{z}_{n-1} = \ell) = \text{multinoulli}(\boldsymbol{\pi}_\ell)$
- If  $\mathbf{z}_n$ 's are real-valued, we have a **state-space model (SSM)**  $p(\mathbf{z}_n | \mathbf{z}_{n-1}) = \mathcal{N}(\mathbf{A}\mathbf{z}_{n-1}, \mathbf{I}_K)$



# State-Space Models

- In the most general form, the state-transition and observation models of an SSM



Using 's' instead of 'z' to refer to states

Using 't' to denote the 'time-step'

HMM is similar to SSM except the state-transition model is a discrete distribution

$g_t, h_t$  can be linear or nonlinear functions

$$\begin{aligned} \mathbf{s}_t | \mathbf{s}_{t-1} &= g_t(\mathbf{s}_{t-1}) + \epsilon_t && \text{(must be a cont. dist. over } \mathbf{s}_t) \\ \mathbf{x}_t | \mathbf{s}_t &= h_t(\mathbf{s}_t) + \delta_t && \text{(can be any dist. over } \mathbf{x}_t) \end{aligned}$$

- Assuming Gaussian noise in the state-transition and observation models

This is a **Gaussian SSM**

$$\begin{aligned} \mathbf{s}_t | \mathbf{s}_{t-1} &\sim \mathcal{N}(\mathbf{s}_t | g_t(\mathbf{s}_{t-1}), \mathbf{Q}_t) \\ \mathbf{x}_t | \mathbf{s}_t &\sim \mathcal{N}(\mathbf{x}_t | h_t(\mathbf{s}_t), \mathbf{R}_t) \end{aligned}$$

If  $g_t, h_t, \mathbf{Q}_t, \mathbf{R}_t$  are independent of  $t$  then it is called a **stationary** model

$g_t, h_t, \mathbf{Q}_t, \mathbf{R}_t$  may be known or can be learned



# State-Space Models: A Simple Example

- Consider the linear Gaussian SSM

$$\mathbf{s}_t | \mathbf{s}_{t-1} = \mathbf{A}_t \mathbf{s}_{t-1} + \epsilon_t$$

$$\mathbf{x}_t | \mathbf{s}_t = \mathbf{B}_t \mathbf{s}_t + \delta_t$$

- Suppose  $\mathbf{x}_t \in \mathbb{R}^2$  denotes the (noisy) observed 2D location of an object
- Suppose  $\mathbf{s}_t \in \mathbb{R}^6$  denotes the “state” vector

$$\mathbf{s}_t = [\text{pos1}, \text{vel1}, \text{accel1}, \text{pos2}, \text{vel2}, \text{accel2}]$$

- Here is an example SSM for this problem with pre-defined  $\mathbf{A}_t$  and  $\mathbf{B}_t$  matrices

$$\mathbf{s}_t = \mathbf{A}_t \mathbf{s}_{t-1} + \epsilon_t$$

$$\mathbf{A}_t = \begin{bmatrix} 1 & \Delta t & \frac{1}{2}(\Delta t)^2 & 0 & 0 & 0 \\ 0 & 1 & \Delta t & 0 & 0 & 0 \\ 0 & 0 & e^{-\alpha \Delta t} & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & \Delta t & \frac{1}{2}(\Delta t)^2 \\ 0 & 0 & 0 & 0 & 1 & \Delta t \\ 0 & 0 & 0 & 0 & 0 & e^{-\alpha \Delta t} \end{bmatrix}$$

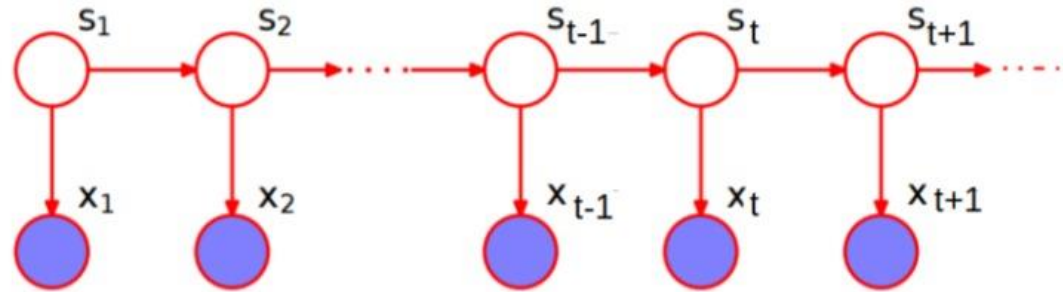
$$\mathbf{x}_t = \mathbf{B}_t \mathbf{s}_t + \delta_t$$

$$\mathbf{B}_t = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix}$$



# Typical Inference Task for Gaussian SSM

- One of the key tasks: Given sequence  $x_1, x_2, \dots, x_T$ , infer latent  $s_1, s_2, \dots, s_T$



- Usually two ways of inferring the latent states

- Infer  $p(s_t | x_1, x_2, \dots, x_t)$ : Called the “filtering” problem

Turns out to be another Gaussian

$$p(s_t | x_1, x_2, \dots, x_t) \propto \underbrace{p(x_t | s_t)}_{\mathcal{N}(x_t | B s_t, R)} \int \underbrace{p(s_t | s_{t-1})}_{\mathcal{N}(s_t | A s_{t-1}, Q)} p(s_{t-1} | x_1, x_2, \dots, x_{t-1}) ds_{t-1}$$

A Gaussian

Kalman Filtering is a popular algorithm for a linear Gaussian SSM

- Infer  $p(s_t | x_1, x_2, \dots, x_t, \dots, x_T)$ : Called the “smoothing” problem

- Some other tasks one can solve for using an SSM

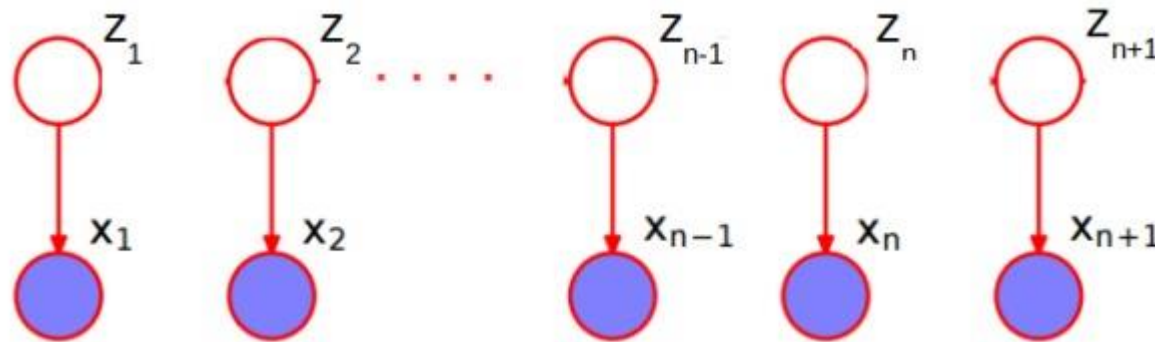
- Predicting future states  $p(s_{t+h} | x_1, x_2, \dots, x_t)$  for  $h \geq 1$ , given observations thus far
- Predicting future observations  $p(x_{t+h} | x_1, x_2, \dots, x_t)$  for  $h \geq 1$ , given observations thus far





# A Special Case

- What if we have i.i.d. latent states, i.e.,  $p(z_n|z_{n-1}) = p(z_n)$ ?



- Discrete case (HMM) becomes a simple mixture model  $p(\mathbf{z}_n|\mathbf{z}_{n-1} = \ell) = p(\mathbf{z}_n) = \text{multinoulli}(\boldsymbol{\pi})$
- Real-valued case (SSM) becomes a PPCA model  $p(\mathbf{z}_n|\mathbf{z}_{n-1}) = p(\mathbf{z}_n) = \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$  or  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Psi})$
- Inference algos for HMM/SSM are thus very similar to that of mixture models/PPCA
  - Only main difference is how the latent variables  $\mathbf{z}_n$ 's are inferred since they aren't i.i.d.
  - E.g., if using EM, only E step needs to change (Bishop Chap 13 has EM for HMM and SSM)





# Conformal Prediction



# Conformal Prediction

- A simple technique to easily obtain confidence intervals
  - In classification, such an interval may refer to the set of highly likely classes for a test input



- For more difficult test inputs, the set would typically be larger
- In a way, conformal prediction gives predictive uncertainty
  - However, unlike Bayesian ML, we don't get model uncertainty
  - Only one model is learned in the standard way and we construct the set of likely classes
  - It's like a black-box method; no change to training procedure for the model



# Conformal Prediction

11

Assume it's a classification model which produces softmax scores

Conformal prediction can be used for regression problems too\*

- Assume we already have a trained model  $\hat{f}$  using some labelled data
- Suppose we get a test input  $X_{test}$  whose true (unknown) label is  $Y_{test}$
- Use  $\hat{f}$  and a **calibration set** of  $n$  examples to generate a prediction set  $\mathcal{C}(X_{test})$  s.t.

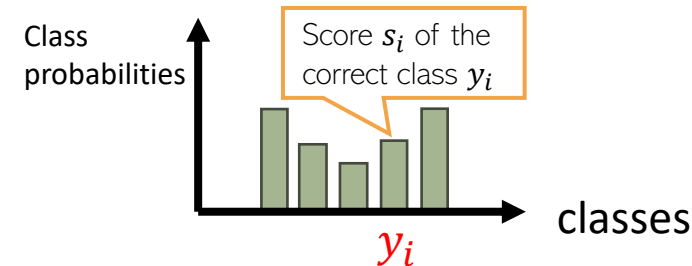
$\alpha$  is a user chosen error rate

$$1 - \alpha \leq p(Y_{test} \in \mathcal{C}(X_{test})) \leq 1 - \alpha + \frac{1}{n + 1}$$

- To construct the set, we first compute, for each example in the calibration set

Probability/score of the correct class  $y_i$  of the input  $x_i$

$$s_i = \hat{f}(x_i)_{y_i}$$

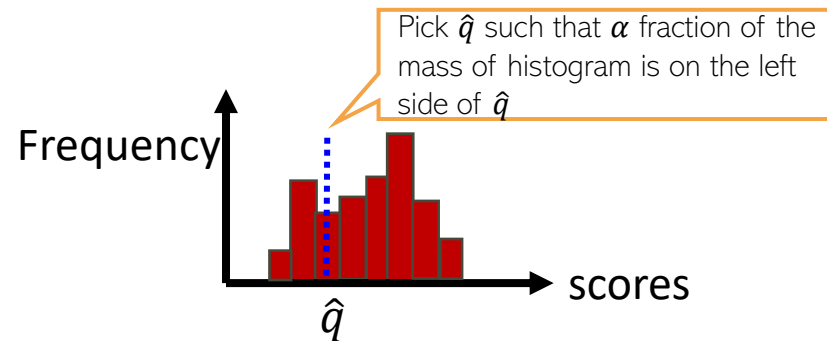


- Use the calibration set scores  $s_1, s_2, \dots, s_n$  to compute their  $\alpha$  quantile

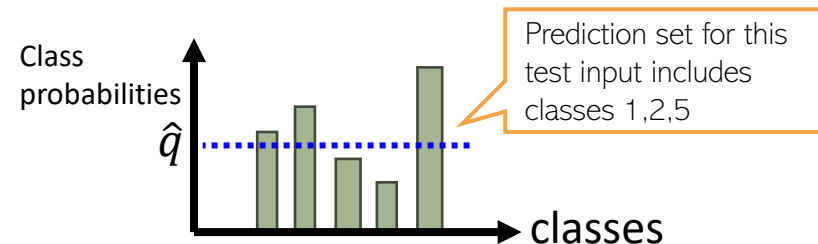


# Conformal Prediction

- Assume the  $\alpha$  (say 0.1) quantile of the calibration set scores is equal to  $\hat{q}$



- Assuming  $n$  is very large, roughly  $(1 - \alpha)$  fraction of inputs will have score higher than  $\hat{q}$
- Given a test input  $X_{test}$ , whose label is unknown, we compute the class probabilities



- Define the prediction set for  $X_{test}$  as

$$\mathcal{C}(X_{test}) = \{y: \hat{f}(X_{test})_y \geq \hat{q}\}$$

Report all the classes whose probability is large enough (the "large enough" value is given by the  $\alpha$  quantile  $\hat{q}$ )



# Conformal Prediction

- A generic black-box method
- Can be easily applied to any already trained classifier
- Predicted set has some nice guarantees

$$1 - \alpha \leq p(Y_{test} \in \mathcal{C}(X_{test})) \leq 1 - \alpha + \frac{1}{n + 1}$$

- Does not make any assumptions on the distribution of the data
  - Thus considered a “distribution-free” approach to uncertainty quantification
- Can also be applied to regression problems\*



# Simulation-based Inference



# Simulation-based Inference

- Suppose we wish to compute the posterior  $p(\theta|D)$
- However, suppose we can't compute the likelihood  $p(D|\theta)$ 
  - Evaluation too expensive, or don't have explicit likelihood
- **Simulation-based Inference (SBI)** approximates  $p(\theta|D)$  as follows

SBI is also known as "Approximate Bayesian Computation" (ABC)

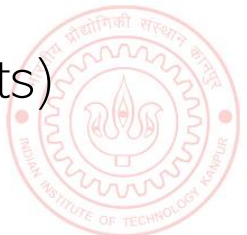
This simulator may be some domain-specific model of the data generation process (e.g., a physics engine, robotics/control simulator, etc)

- For  $i = 1, 2, \dots, S$ 
  - Draw a random  $\theta^{(i)}$  the prior  $p(\theta)$ . Simulate a dataset  $D^{(i)}$  from some simulator using  $\theta^{(i)}$
  - Check how "similar"  $D^{(i)}$  is to  $D$ . Define a suitable distance to measure this, e.g.,

$$d_i = \|s(D^{(i)}) - s(D)\|$$

Here  $s(\cdot)$  denotes a "summary statistics" which provides a summary of the dataset (e.g., its mean and covariance) which makes the comparison easier

- Define the weight of  $\theta^{(i)}$  as inversely proportional to  $d_i$ , e.g.,  $w_i \propto \exp(-d_i)$
- The approximate posterior is  $\{w_i, \theta^{(i)}\}_{i=1}^S$
- The vanilla SBI/ABC can be inefficient in practice (most  $\theta^{(i)}$ 's may have low weights)
  - More efficient versions proposed in recent research, e.g, neural conditional density estimators
  - Check out this package for code and links to other methods: <https://github.com/sbi-dev/sbi>





# Conclusion

- Probabilistic modeling provides a natural way to think about models of data
- Many benefits as compared to non-probabilistic approaches
  - Easier to model and leverage **uncertainty** in data/parameters
  - Principle of **marginalization** while making prediction
  - Easier to encode **prior knowledge** about the problem (via prior/likelihood distributions)
  - Easier to handle **missing data** (by marginalizing it out if possible, or by treating as latent variable)
  - Easier to build complex models can be neatly combining/extending simpler probabilistic models
  - Easier to learn the “right model” (hyperparameter estimation, nonparametric Bayesian models)
- **Bayesian approaches** as well as single model based uncertainty
- Uncertainty is important but proper calibration of uncertainty is also important
- Fast-moving field, lots of recent advances on new models and inference methods



# Thank You!

