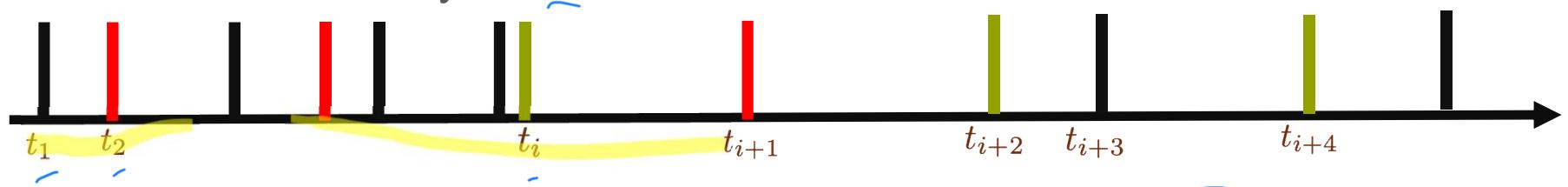
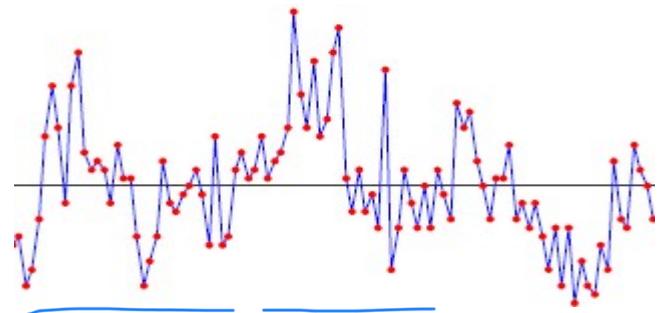


Temporal Sequences

Sunita Sarawagi

Temporal Sequences

- Regular time-series
 - Daily traffic on individual webpages from different regions in Wikipedia
 - Hourly load on various servers of different services in a Data center
 - Monthly demand for products from different regions in Flipkart
- Irregular Event sequences
 - User visits to a music service and the song played
 - Event logs of a system
 - Attacks on a system



Temporal Sequences

- A temporal sequence

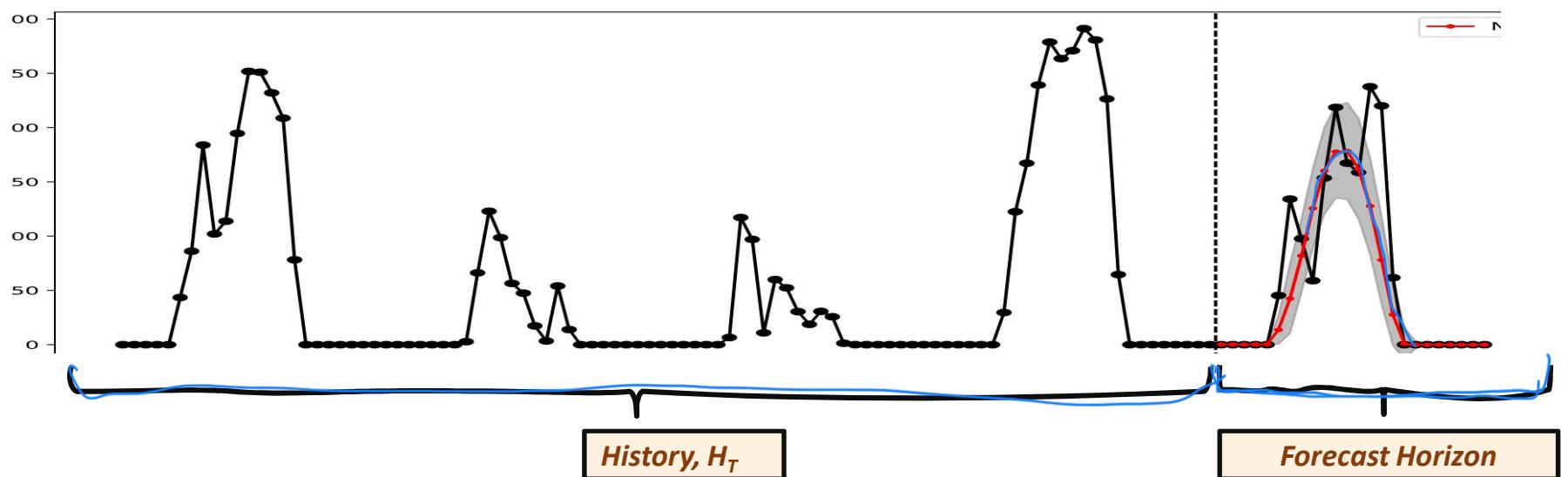
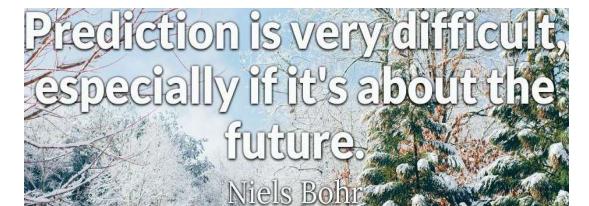
$$(t_1^i, x_1^i, y_1^i), (t_2^i, x_2^i, y_2^i), \dots, (t_n^i, x_n^i, y_n^i)$$

Time Input features Value

- Several such sequences indexed by i
 - Sequences are related to each other
 - The index i could span a multi-dimensional space e.g. server and application, or product-id and region.

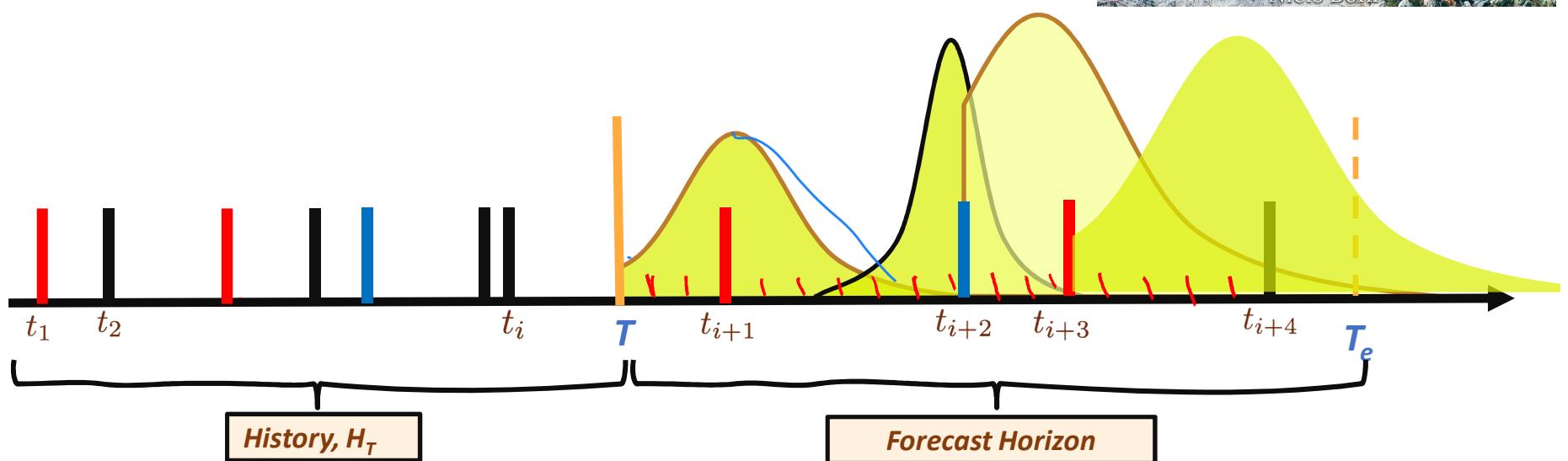
Tasks on Temporal sequences

- Probabilistic Forecasting:
 - short Vs long-term.



Tasks on Temporal sequences

- Forecasting in event sequences:
 - Predict both time of event and type of event.



Tasks on Temporal Sequences

- Outlier detection
- Missing value imputation

Our focus here is on forecasting.

Key modelling questions often carry over to other tasks.

Conventional methods

- Statistical time-series methods e.g. ARIMA, Box-Jenkins
 - Local, cannot handle features x_j , non-stationarity, interactions
- State space models e.g. Kalman filters
 - Simplifying assumption (linear Gaussian) may not hold in practice.
- Classical machine learning: e.g. Regression methods
 - Domain experts featurize history e.g. wavelets, Fourier coefficient
 - Powerful regressors like Boosted regression trees to predict y
- Can deep learning provide higher gains than these?

Demand Forecasting in Flipkart

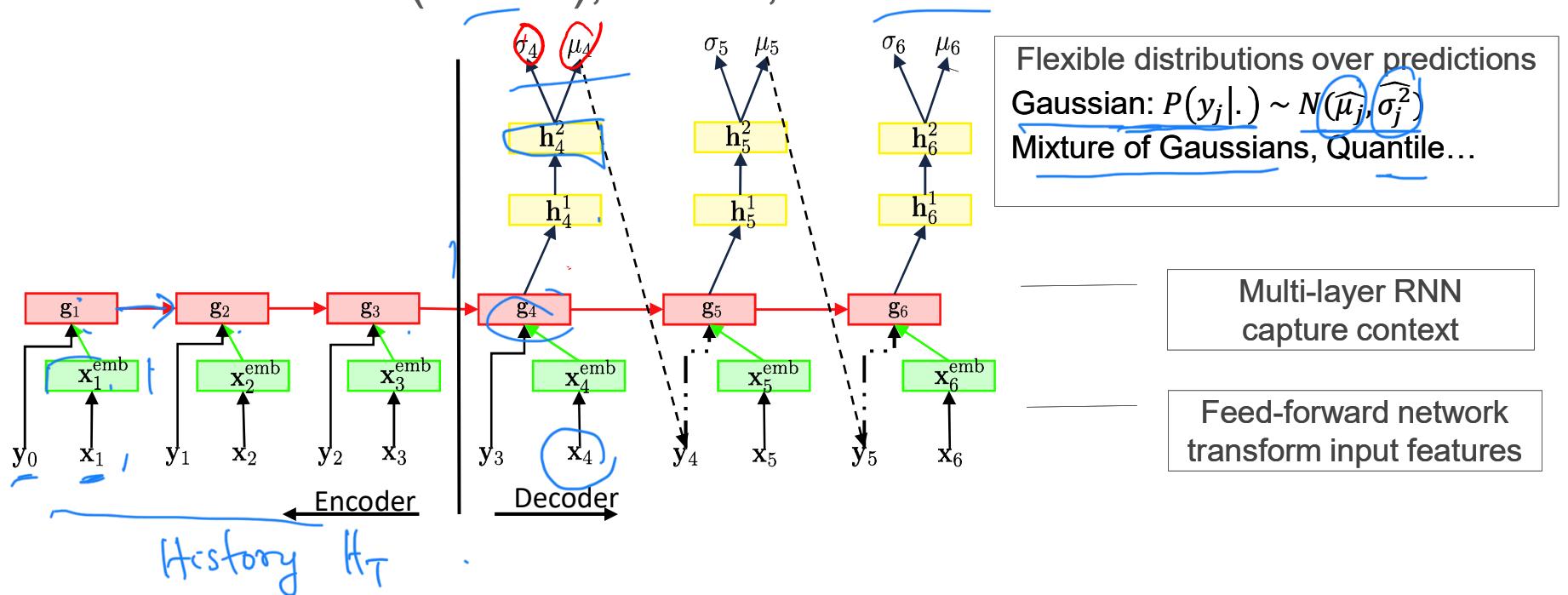
- Eight weeks out inventory management
- Huge catalogue of products --- 100 millions.
- Sold across different regions.
- Large number of confounding factors.¹
- Large divergence in scale and variance in the output values

Large number and variety of features

Demand Factor	Feature Name	Description	Data Type
Product	Product ID Product Tier Historical Out-of-Stock %	Captures “latent” factors of the product Tiers are based on popularity of products Assumed to be always In-Stock during Testing phase	1-Hot Categorical Numerical
Product Visibility	Sale event type Deal Card Banner on Homepage	Category-level of the Sale event	Categorical Binary Binary
Price	Listed Price Discounted Price Effective Price Difference in price from historical mean Historical Min Price Historical Max Price Avg. discount on similar products	Govt. mandated Max. Retail Price Price offered by Flipkart Final price after cash-back, product exchange, freebie, etc. Captures “stable” price of the product Products in the same category are considered similar	Numerical Numerical Numerical Numerical Numerical Numerical Numerical
Convenience	No-cost EMI Product Exchange Exclusive to Flipkart		Binary Binary Binary
Time	Week of the month Lag feature - Price Lag feature - Sale Time elapsed since last price change Time since last event or upcoming event Time since first order	Captures seasonality Mean of past n weeks’ price Mean of past n weeks’ sale units Sales peak when price changes, but reverts quickly to previous volume Units sold dips immediately before and after a Sale event Captures “Newness” of Product	Numerical Numerical Numerical Numerical Numerical Numerical

Neural Model for Time Series Forecasting

Encode history using neural sequence model: Recurrent Neural Network (RNNs), CNNs, Transformer w. self-attention



Predicting a Gaussian distribution on forecasts

$$\mathbf{g}_n^i = \underline{RNN}([y_{j-1}^i, \mathbf{x}_j^i] : j = 1 \dots n | \theta_{\text{enc}})$$

$$\mathbf{h}_j^i = \underline{FF}([\mathbf{g}_n^i, \mathbf{x}_j^i] : j = n+1 \dots n+K | \theta_{\text{dec}})$$

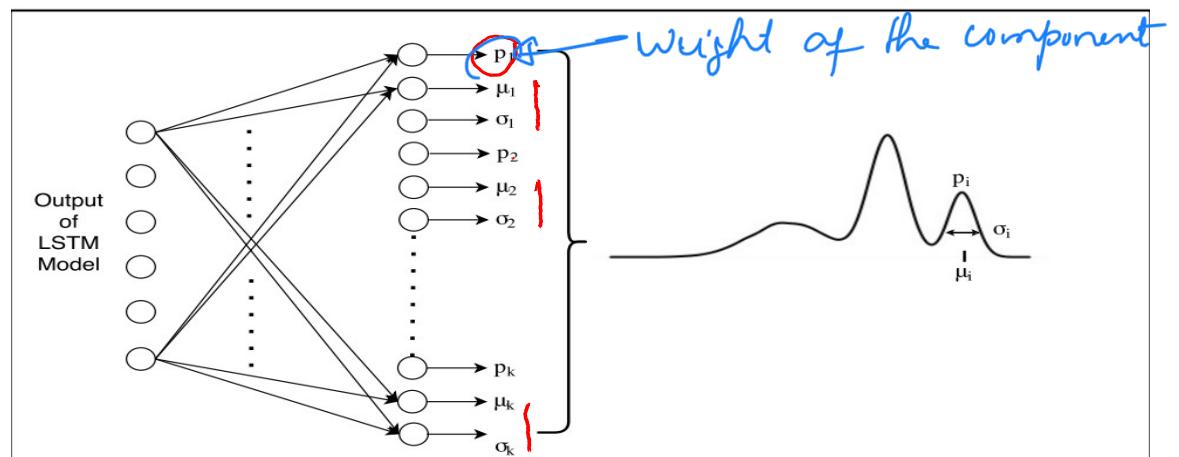
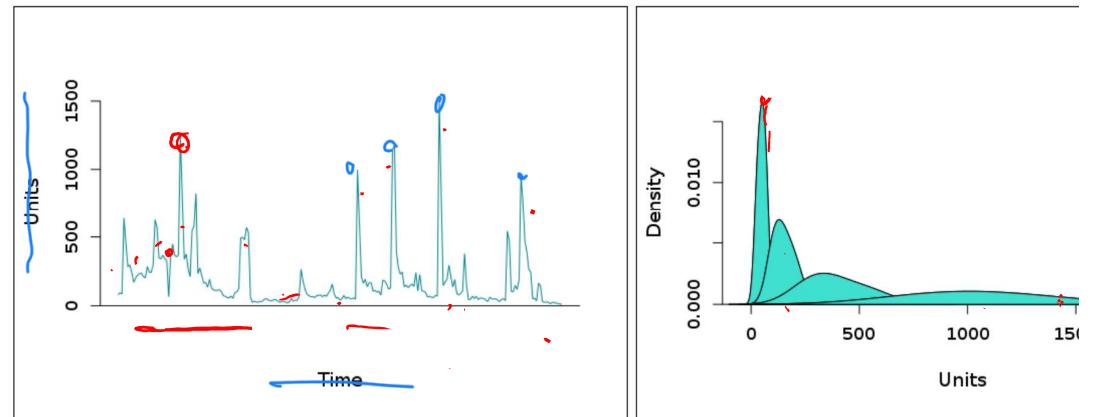
$$\mu_j^i = \theta_\mu[\mathbf{h}_j^i, 1], \quad \sigma_j^i = \log(1 + \exp(\tilde{\theta_\sigma}[\mathbf{h}_j^i, 1]))$$

$$\Pr(y_j^i | \mathbf{x}_j^i, (\mathbf{x}_1^i, y_1^i), \dots, (\mathbf{x}_{j-1}^i, y_{j-1}^i)) = \mathcal{N}(\mu_j^i, \sigma_j^i)$$

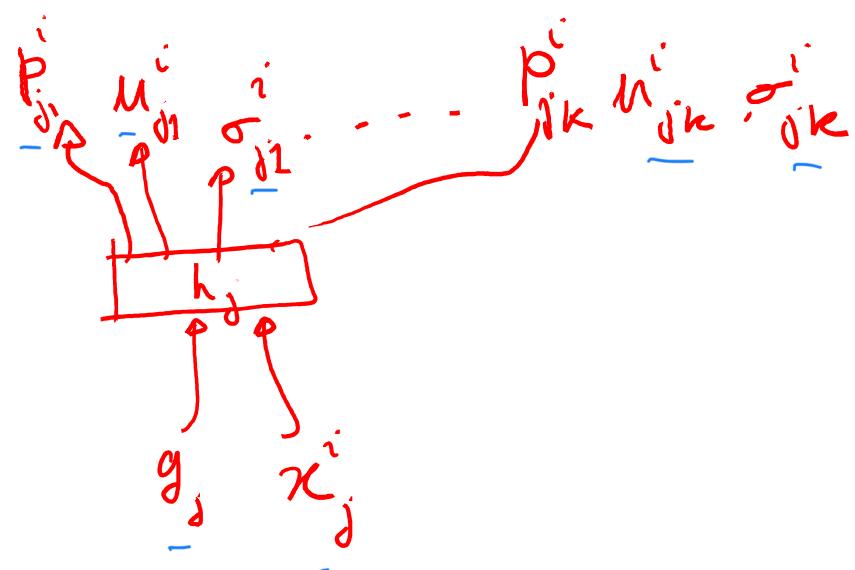
Output Values highly multi-modal

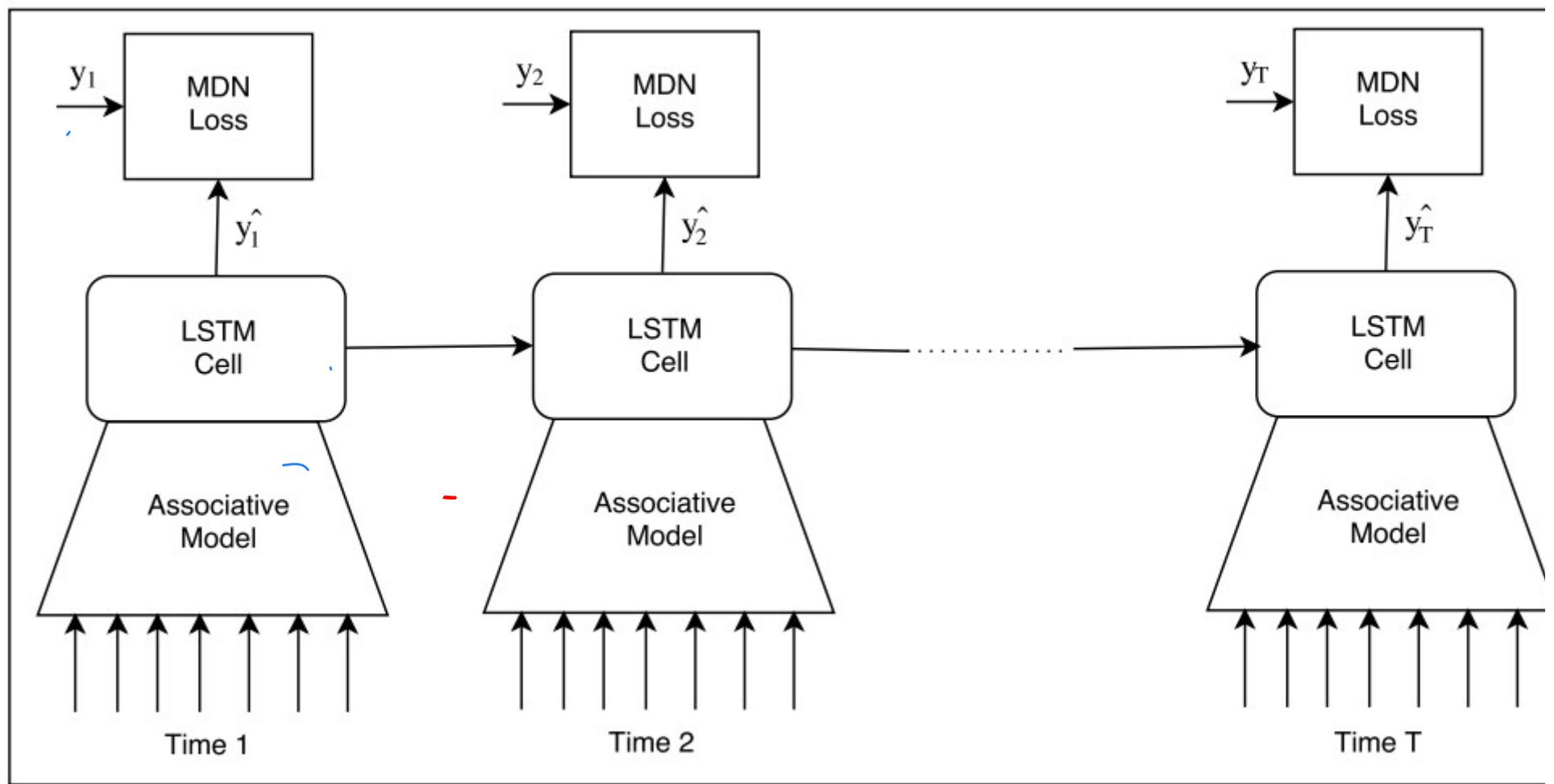
- Single Gaussian insufficient
- Use mixture of Gaussians
- Last layer outputs k mixture components, means, variance.

$0, 1, 0.1, 0.2, 100, 200,$
 $1 \leftrightarrow 0$



$$P(y_d^i | H_i, x_d^i) = \sum_{k=1}^K p_{jk}^i N(y_d^i | \mu_{jk}^i; \sigma_{jk}^i)$$





Training

- Full joint distribution modelled and maximized

$$\sum_{i:\text{series}} \sum_{j:\text{positions}} \log P(y_j^i | \theta(x_j^i, (x_1^i, y_1^i), \dots, (x_{j-1}^i, y_{j-1}^i)))$$

$\sim \mathcal{N}(y_j^i | \mu_j^i; \sigma_j^{i,2})$.

- Parameters θ trained end to end jointly over all series
- Training is efficient and easily parallelizable

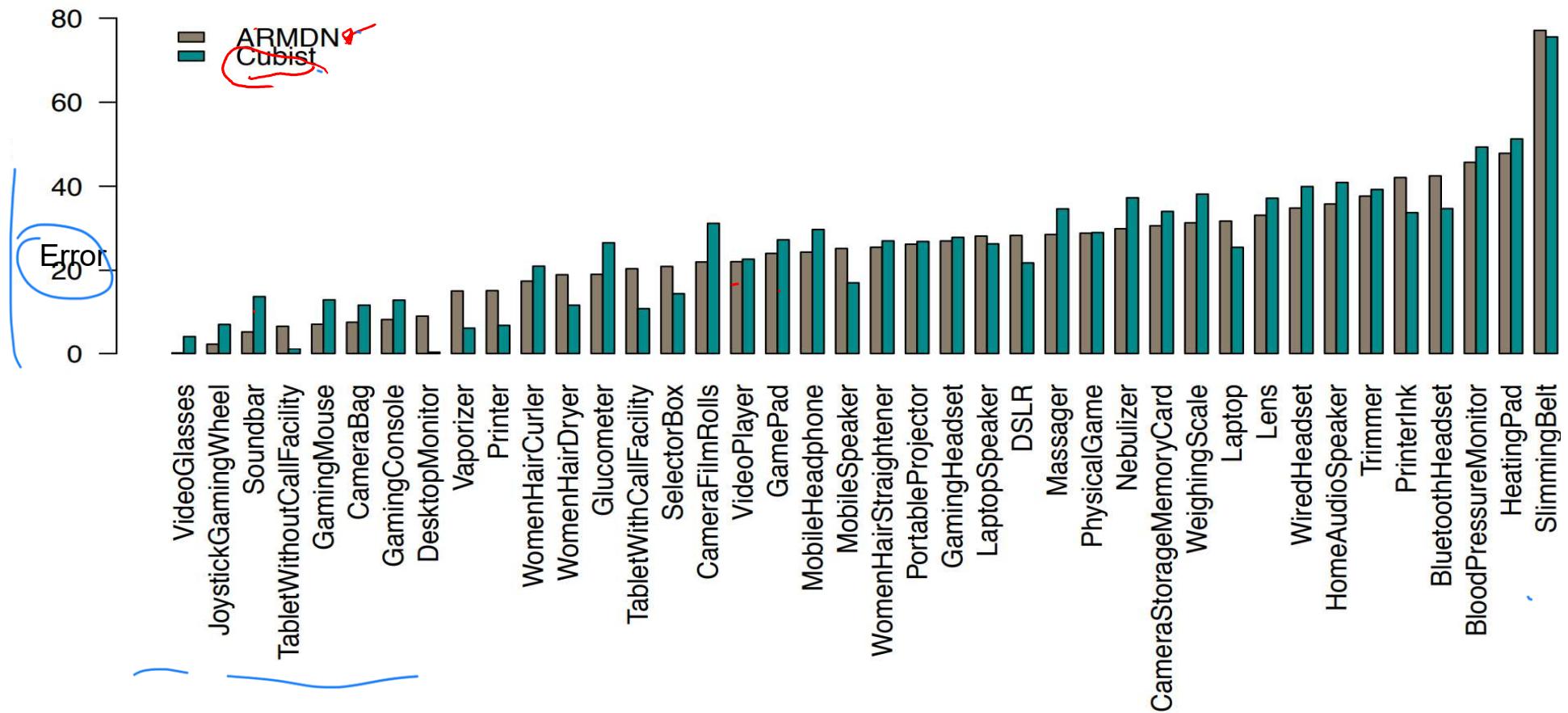
$$\nabla_{\theta} \log P(y_j^i | \hat{u}_j^i, \hat{\sigma}_{j^i}^2) = \frac{-\frac{1}{2} (y_j^i - \hat{u}_j(\theta))^2}{\hat{\sigma}_{j^i}^2(\theta)} - \log \hat{\sigma}_-(\theta)$$

RNN parameters
 input encoding layers
 output layers.

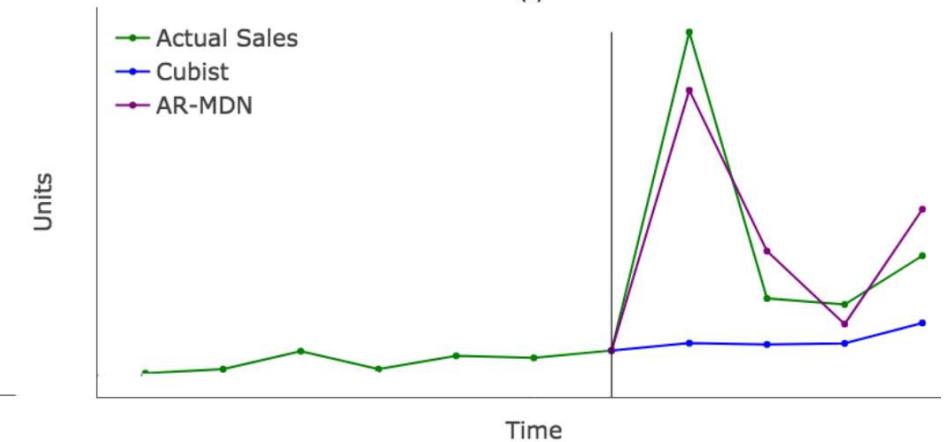
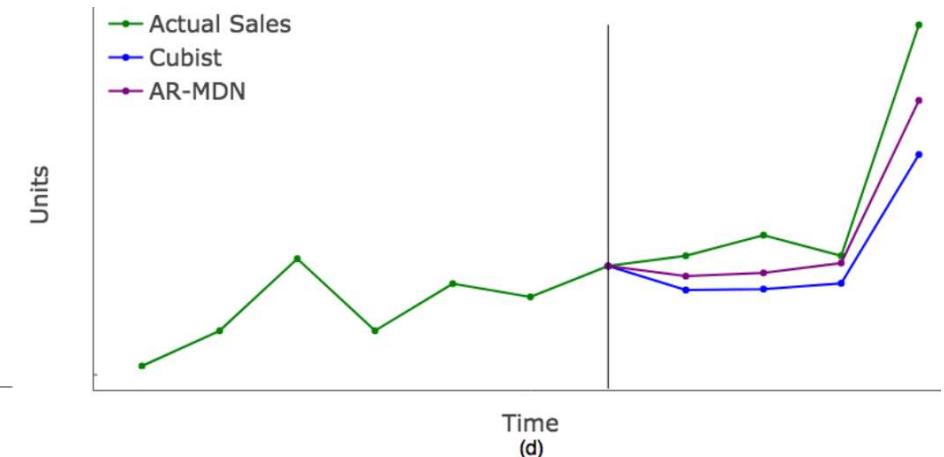
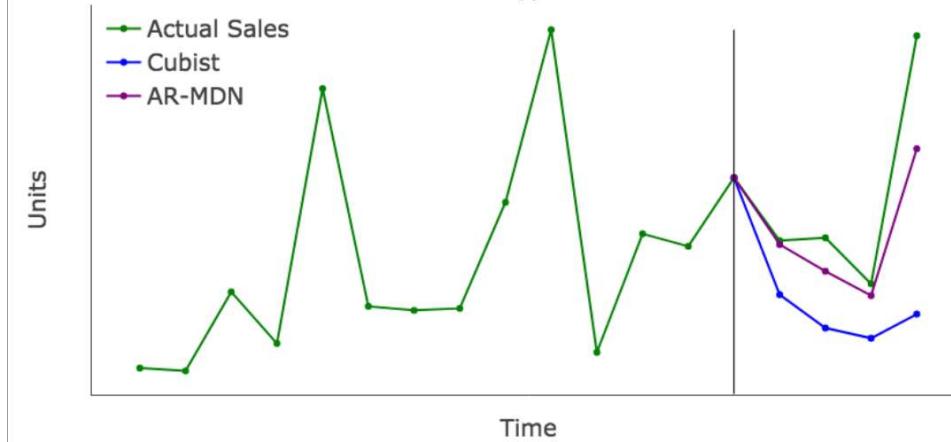
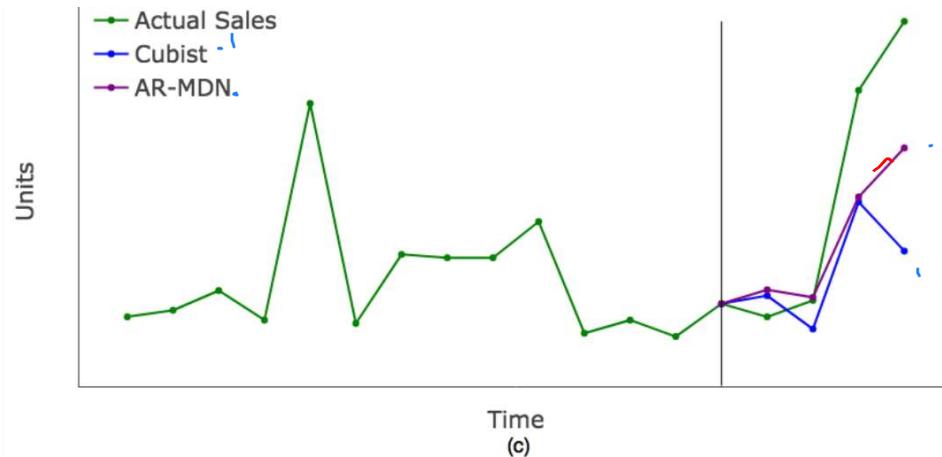
M66:

$$\nabla_{\theta} \log \sum_{k=1}^K p_{jik}(\theta) \frac{e^{-\frac{1}{2} (y_j^i - u_{jk}(\theta))^2}}{\sqrt{\pi} \sigma_k(\theta) \hat{\sigma}_{jk}^2(\theta)}$$

Neural method better than SOTA non-neural method



Some anecdotes



Limitations of this model

- Does not output the joint distribution over multiple predictions it makes both along time and along different items.
- Tools to address these limitations
 - Gaussian Processes
 - Gaussian Copula

How do you detect if your data is a time series?
if " data follows the distribution
-al assumption.

Check whether the
standardized values:

$$\frac{z_j^i}{\sigma_j^i} = \frac{y_j^i - \hat{u}_j^i}{(\sigma_j^i)}$$

used for modeling likelihood.
Do $\{z_j^i\}$: values look

like samples from a $\mathcal{N}(0, 1)$ distribution?

