

ISI Informal Workshop on Learning, Algorithms and Networks
February 18, 2020

Adversarial Games for generative modeling of Temporally-Marked Event Sequences

Giuseppe Manco

Institute of High Performance Computing and Networks



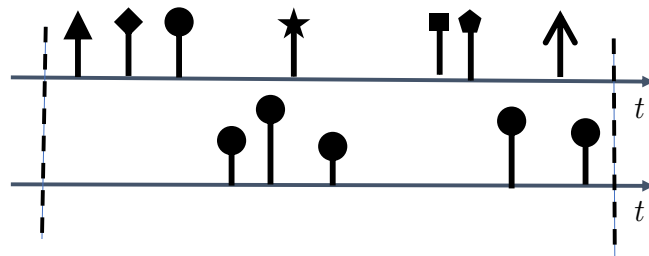
Temporally Marked Point Processes (TMPPs)

- Data = (noisy) observations of events generated by some complex process
 - Each event refers to a single process (execution) instance and stores at least:
 - (1) a time-stamp and (2) a categorical attribute representing the type of the event
 - **Trace** \mathcal{H} = sequence of all events linked to a process instance (stores instance's history)

$$\mathcal{H} = \{ \mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_m \}$$

$$\mathbf{e}_i = (\mathbf{a}_i, t_i)$$

- Marked Temporal Point Processes' view



Discrete events in
continuous time

intra-event times may vary over a wider
range of temporal scales
(differently from “classic” time-series data)

Application contexts

- **Business Process (or Workflow) Management events**
 - execution traces of a process, storing information on performed activities
- **Medical events**
 - acute incidents, doctor's visits, tests, diagnoses, and medications
- **Consumer behavior**
 - purchasing patterns
- **"Quantified self" data**
 - Wearable devices and apps to record eating, traveling, working, sleeping, waking
- **Social media actions**
 - previous posts, shares, comments, messages,...
- **Smart cities and mobility patterns**
 - trajectories, taxi/car/public transportation adoptions, etc.
- **Smart industry**
 - Optimization of the production, predictive maintenance...

Goal

- **General goal:** Understand the **structural and temporal dynamics** of process traces
 - can provide insights on the complex patterns that govern the process
 - can be used to forecast future events
- **Specific objectives:**
 - **Data generation:** Generate new realistic traces from scratch
 - to have a surrogate of real data (due to privacy/scalability constraints), or for simulation analyses
 - **Predict future events:** Given an incomplete trace (partial history of a process instance)

$$\mathcal{H}_{<h} = \{e_1, \dots, e_{h-1}\} \quad \text{with} \quad h \leq m$$

Make forecasts on the subsequent events (structured-prediction / conditional-generation task)

Approach: Learn a Probabilistic Generative Model

- Probability distribution to learn

$$P(\mathbf{e}_h | \mathcal{H}_{<h}) = P(\mathbf{a}_h, t_h | \mathcal{H}_{<h})$$

- Possible decompositions
 - Independence

$$P(\mathbf{a}_h, t_h | \mathcal{H}_{<h}) = P(\mathbf{a}_h | \mathcal{H}_{<h})P(t_h | \mathcal{H}_{<h})$$

- Time is context – dependent

$$P(\mathbf{a}_h, t_h | \mathcal{H}_{<h}) = P(\mathbf{a}_h | \mathcal{H}_{<h})P(t_h | \mathbf{a}_h, \mathcal{H}_{<h})$$

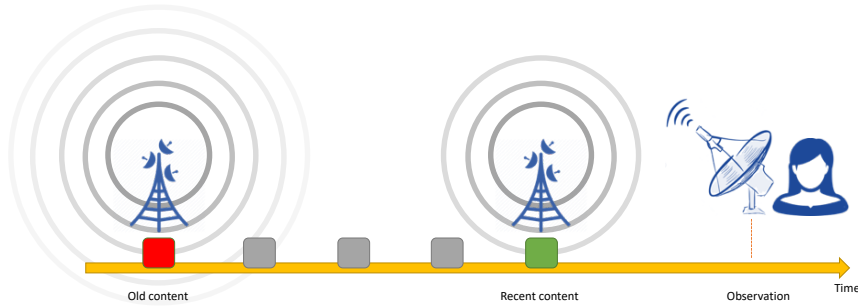
Intermezzo: ML estimation

- Given a sample $X = \{x_1, x_2, \dots, x_n\}$
 - sampled from a true distribution \mathbb{P}_r
- Given a proposal distribution \mathbb{P}_θ parametrized by θ
- Find the parameter $\hat{\theta}$ that optimizes the likelihood:

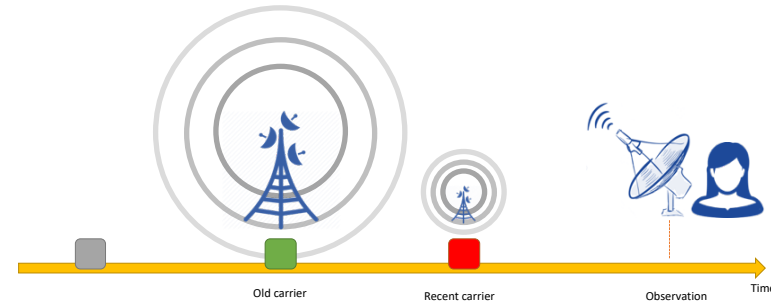
$$\begin{aligned}\hat{\theta} &= \operatorname{argmax}_\theta P_\theta(X) \\ &= \operatorname{argmax}_\theta \prod_i P_\theta(x_i) \\ &= \operatorname{argmax}_\theta \sum_i \log P_\theta(x_i) \\ &= \operatorname{argmax}_\theta \mathbb{E}_{x \sim \mathbb{P}_r} \log P_\theta(x)\end{aligned}$$

Parameterized Models

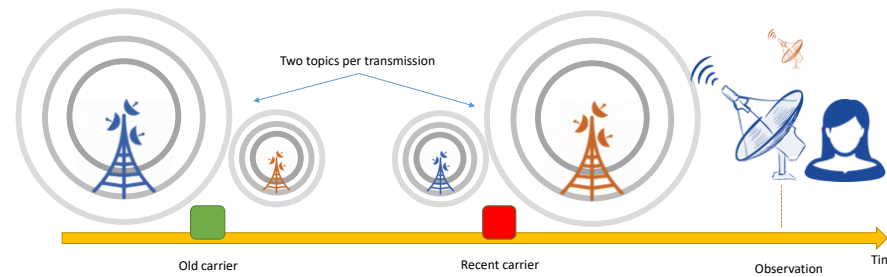
- Three key elements for modeling contagion



Time Dependency



Carrier Dependency



Topic Dependency

Auto-Regressive generative models

- Explicit probability model
 - factorized as a product of conditional per-step distributions (chain rule)

$$P(\mathcal{H}) = \prod_{h=1}^m P(\mathbf{e}_h | \mathcal{H}_{<h})$$

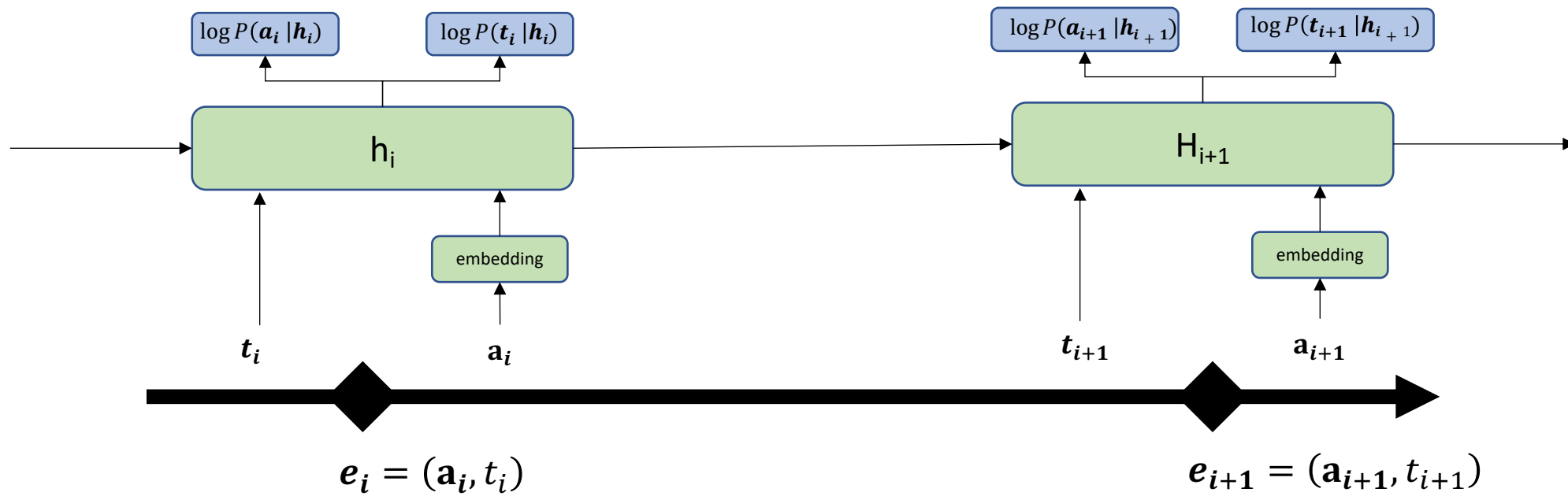
- Conditional probs. are approximated with an NN (usually an RNN)

$$P(\mathbf{e}_h | \mathcal{H}_{<h}) \approx f_{\theta}(\mathbf{s}_h)$$
$$\mathbf{s}_h = \text{RNN}(\mathbf{e}_{h-1}, \mathbf{s}_{h-1})$$

- **Learning**: tractable **maximum-likelihood (ML)** training
 - optimizes exact likelihood
- Inference:
 - generates a suffix (or an entire sequence) via incremental auto-regression

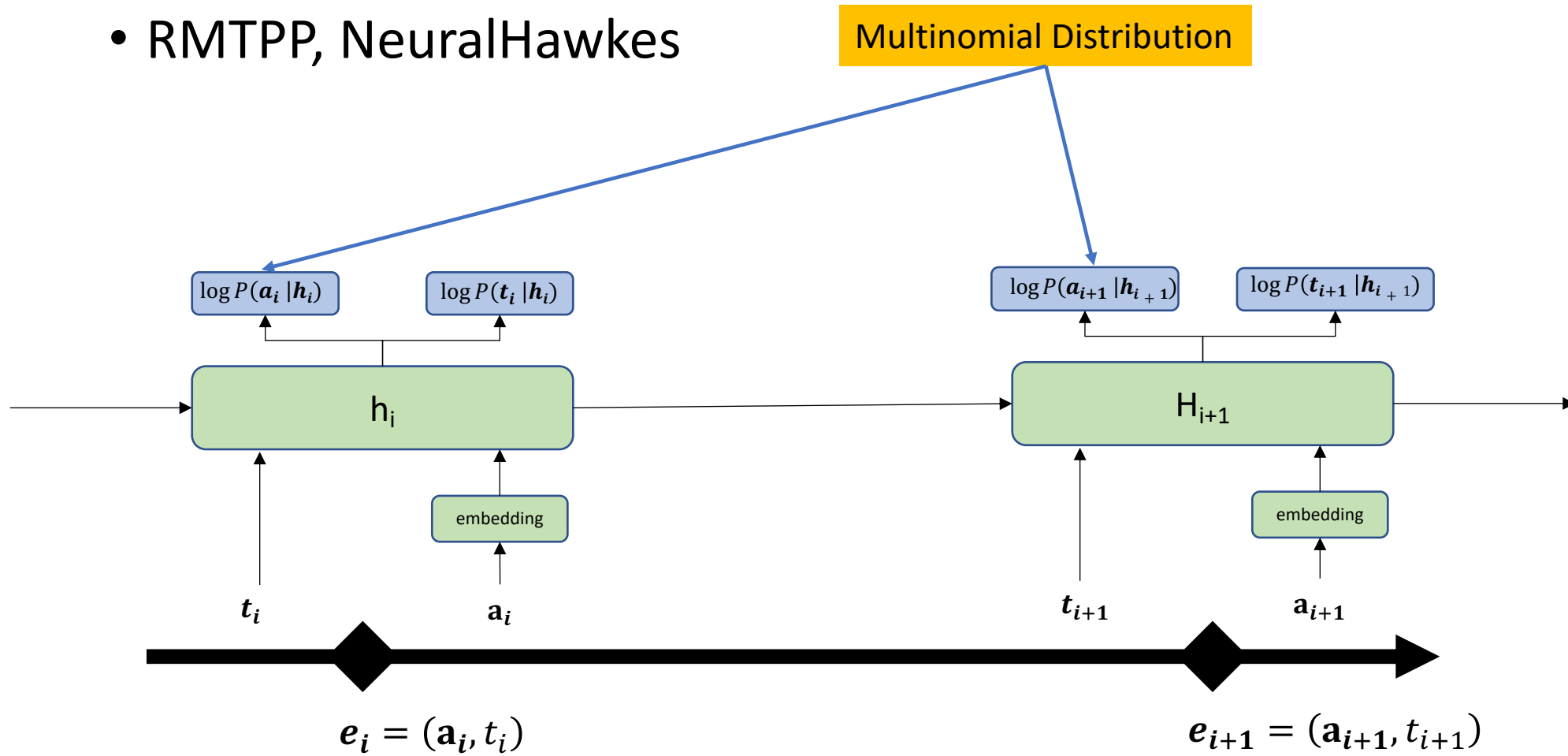
Examples

- RMTTPP, NeuralHawkes [Du et al 2016] [Mei 2017]



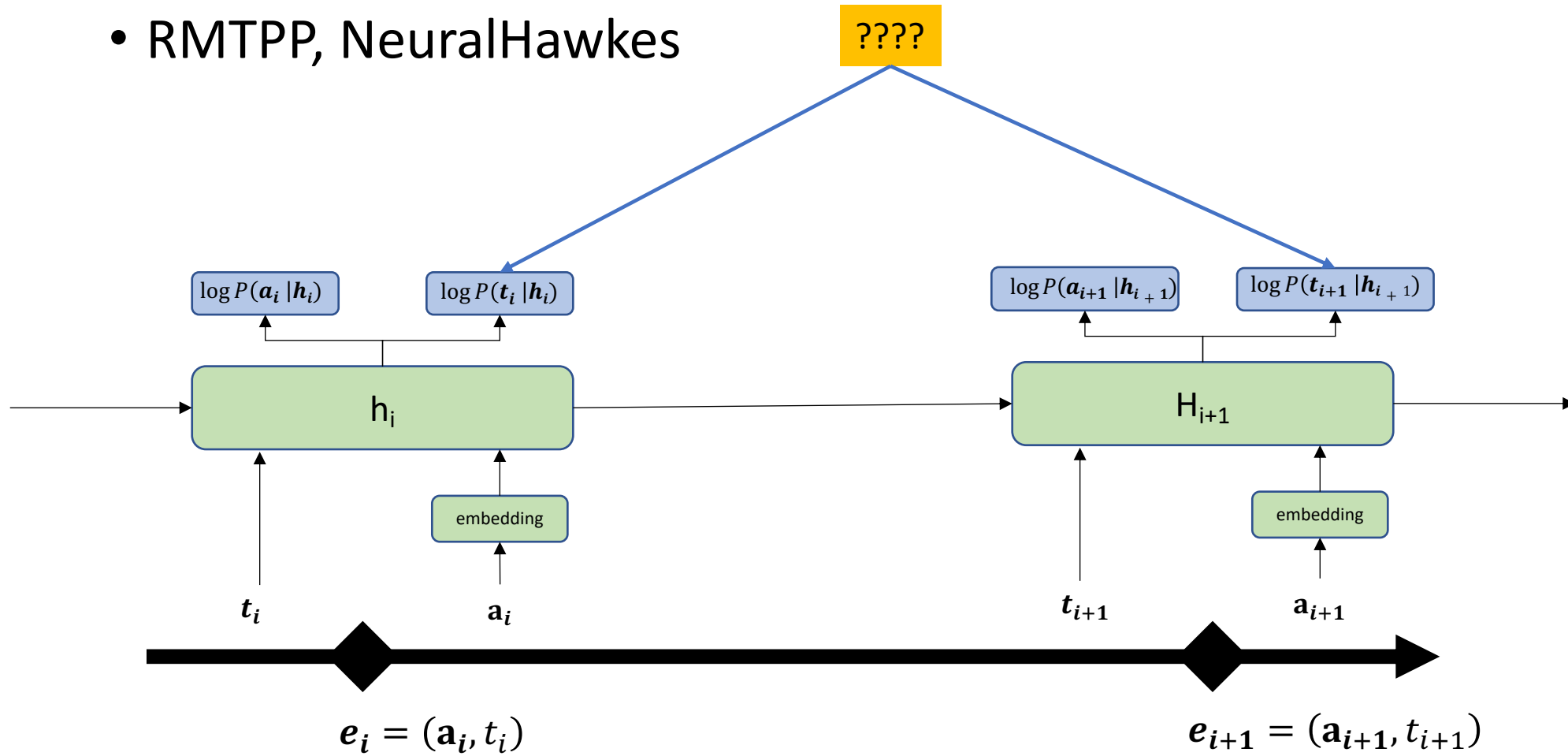
Examples

- RMTTPP, NeuralHawkes

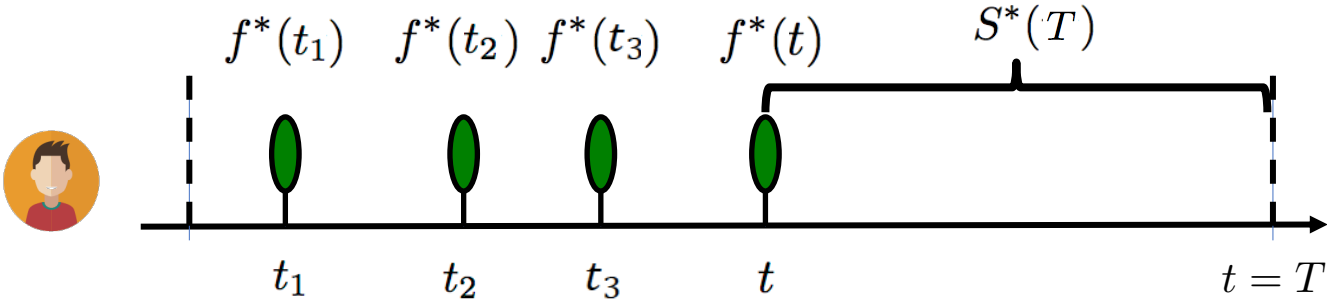
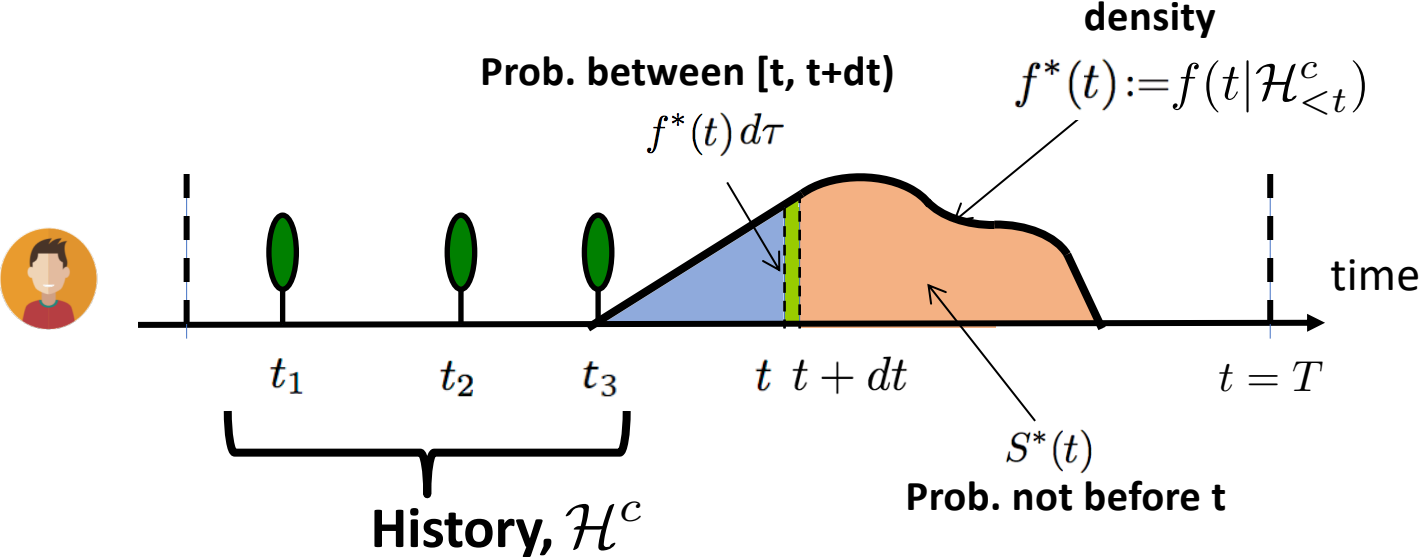


Examples

- RMTTPP, NeuralHawkes

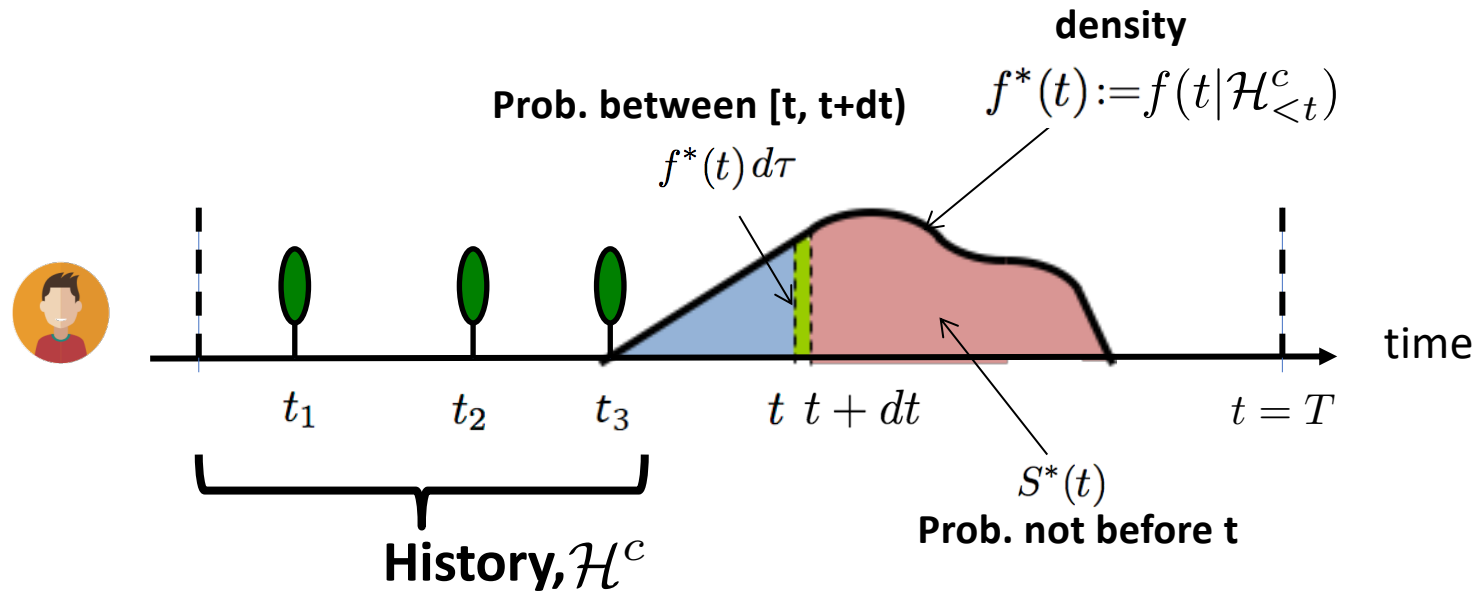


Model time as a random variable



Likelihood of a timeline: $f^*(t_1) f^*(t_2) f^*(t_3) f^*(t) S^*(T)$

Intensity function



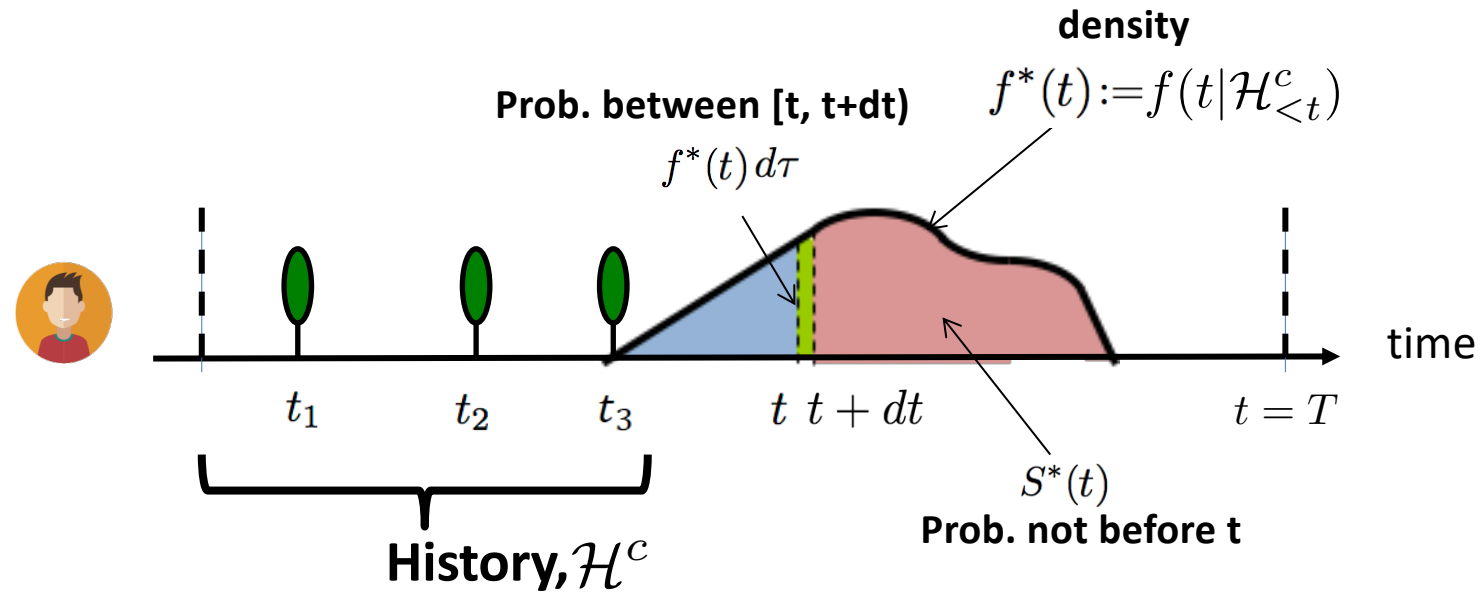
Intensity:

Probability between $[t, t+dt)$ but not before t

$$\lambda^*(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t < T \leq t + \Delta t | T > t)}{\Delta t} = \frac{f^*(t)}{S^*(t)}$$

$\lambda^*(t)$ It is a rate = # of events / unit of time

Intensity function



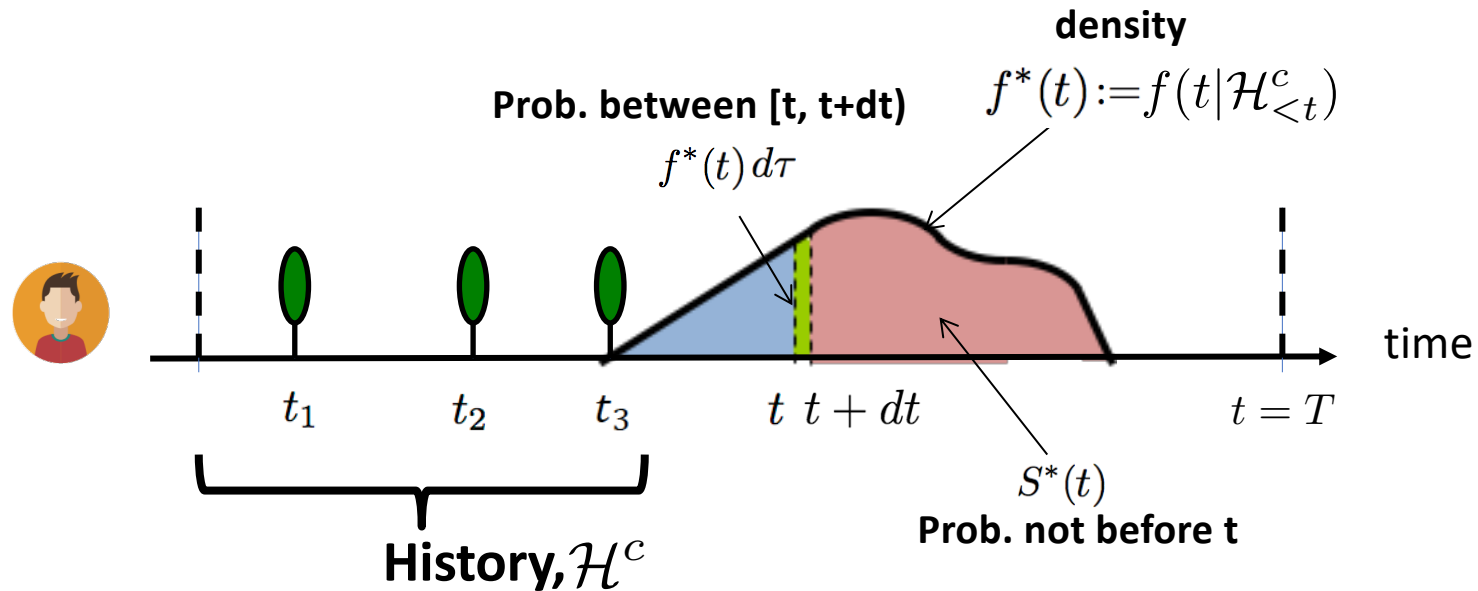
Intensity:

Probability between $[t, t+dt)$ but not before t

$$f^*(t) = \lambda^*(t) S^*(t)$$

$\lambda^*(t)$ It is a rate = # of events / unit of time

Intensity function



Intensity:

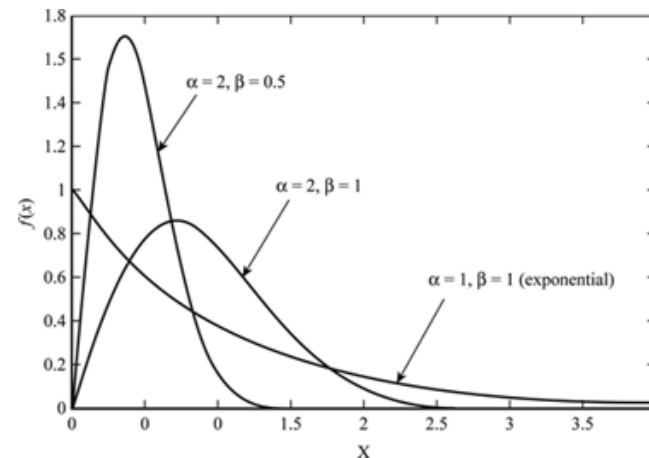
Probability between $[t, t+dt)$ but not before t

$$S^*(t) = \exp\left(-\int_0^t \lambda^*(u) du\right)$$

$\lambda^*(t)$ It is a rate = # of events / unit of time

Appropriate distributions for parametric estimation

Name	$S^*(t)$	$\lambda^*(t)$	$f^*(t)$
Weibull	$\exp(-at^\beta)$	$\alpha\beta t^{\beta-1}$	$\alpha\beta t^{\beta-1}\exp(-at^\beta)$
Log-Logistic	$\frac{1}{1+(at)^\beta}$	$\frac{\beta\alpha^\beta t^{\beta-1}}{1+(at)^\beta}$	$\frac{\beta\alpha^\beta t^{\beta-1}}{(1+(at)^\beta)^2}$
Log-normal	$1 - \Phi\left(\frac{\log t + \log \alpha}{\sigma}\right)$...	$\frac{1}{t\sqrt{2\pi\sigma}} \exp\left(-\frac{1}{2\sigma^2}(\log t + \log \alpha)^2\right)$
Exponential	$\exp(-at)$	α	$\alpha\exp(-at)$
Hawkes	...	$\mu + \alpha \sum_{t_i \in \mathcal{H}_{<t}} \beta_i \exp(t - t_i)$...



Why it doesn't work

- ML is consistent: in principle, it can learn any distribution, provided that it is given infinite data and perfect model space.
 - Minimizing the ML is equivalent to minimize the Kullback-Leibler (KL) divergence between the true distribution \mathbb{P}_r and the parametric distribution \mathbb{P}_θ

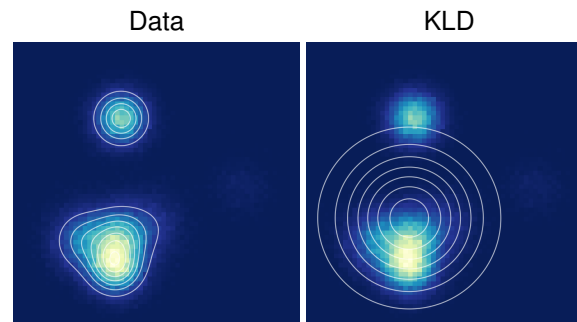
$$KL[\mathbb{P}_r | \mathbb{P}_\theta] = \int P_r(x) \log \frac{P_r(x)}{P_\theta(x)} dx$$

- However, in real settings (due to model mis-specification and finite data), it tends to produce **overgeneralized models**.

Why it doesn't work

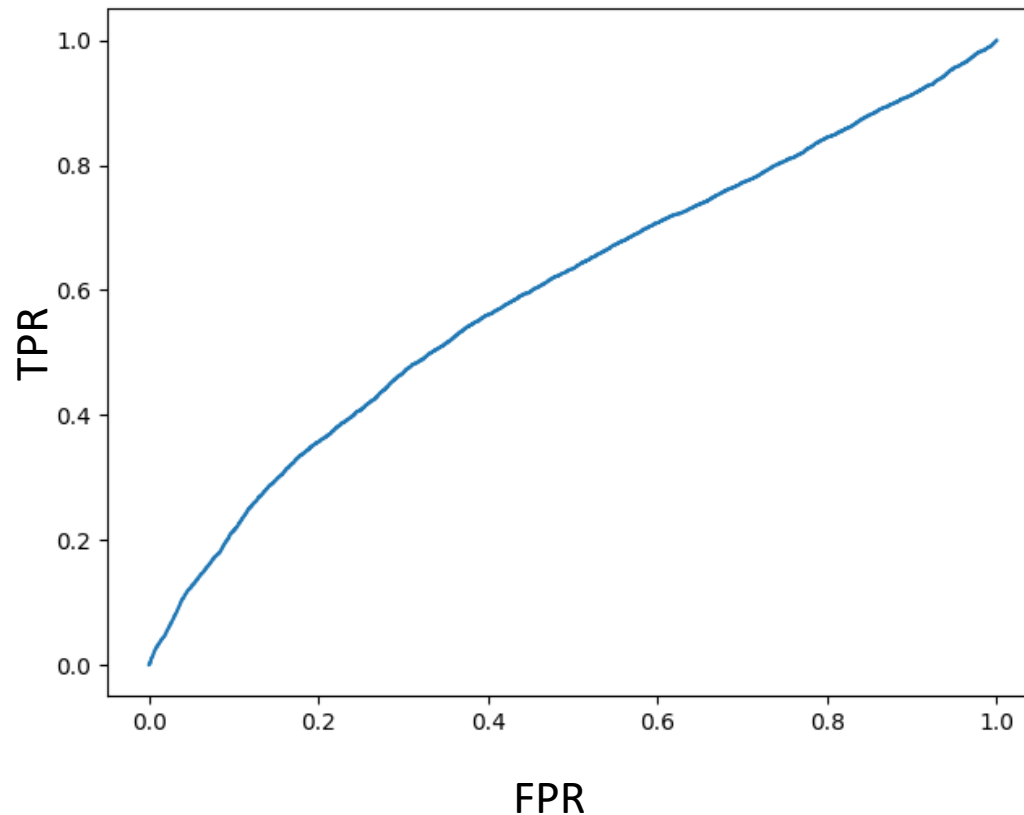
$$KL[\mathbb{P}_r|\mathbb{P}_\theta] = \int P_r(x) \log \frac{P_r(x)}{P_\theta(x)} dx$$

- When $P_r(x) > P_\theta(x)$, large regions of \mathbb{P}_r get low values in \mathbb{P}_θ . Their contribution in $KL[\mathbb{P}_r|\mathbb{P}_\theta]$ tends to infinity.
- However, when $P_r(x) < P_\theta(x)$ then x has a low (true) probability, but high probability of being generated by the model. The contribution to $KL[\mathbb{P}_r|\mathbb{P}_\theta]$ tends to 0.
 - [Arjovsky 2017]



Source: [Theis et al 2016]

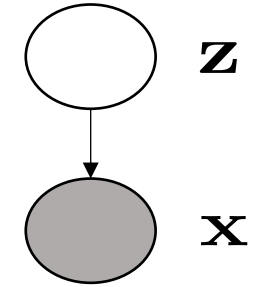
Example: activations on Twitter, Weibull model



- Noisy temporal patterns, data collection bias

Latent generative models

- Assume a stochastic generative process for observed data, ruled by latent/hidden variables \mathbf{z}
 - Can capture high-level structure in event sequences and multiple sources of variability
 - E.g.: for business process instances, \mathbf{z} can capture variations due to context factors, alternative process configurations/variants, organizational changes
- Deep latent generative models:
 - learn $P(\mathbf{x}|\mathbf{z})$ via a neural network: **no assumption** on distribution's shape
 - (Sequential) **Variational Autoencoders**



$$\mathbf{z} \sim P_{\phi}(\cdot)$$

$$\mathbf{x} \sim P_{\theta}(\cdot|\mathbf{z})$$

$$P(\mathbf{x}) = \int P(\mathbf{x}|\mathbf{z})P(\mathbf{z})d\mathbf{z}$$

Latent Generative Models

- Introduce a proposal distribution Q_ϕ parameterized by ϕ
- Approximate the likelihood with

$$\log P(x) \geq \underbrace{\mathbb{E}_{z \sim Q_\phi} [\log P_\theta(x|z)] - \text{KL}[Q_\phi(z)|P(z)]}_{\text{Evidence Lower Bound (ELBO)}}$$

Evidence Lower Bound (ELBO)

- Optimize the ELBO with respect to ϕ and θ

Latent Generative Models

- Pros w.r.t. autoregressive models:
 - More robust to overfitting (regularization effect of latent modeling)
 - Useful latent representations (estimated via inference queries $P(z|x)$)
- Con: Imprecise generated samples, mainly due to two sources:
 - Still ML-oriented training
 - combined to approximate (ELBO) optimization, and assumption on z 's prior/posterior

What happens if we use Discriminative Learning instead?

- PATH: Predicting Activation Time Horizon
 - Focuses on predicting the activation of (groups of) users within a time horizon T_h
 - Focuses on specific users and predicts their behavior based on the effect of the partial cascading behavior of specific users
 - Captures cumulative history via the embedding
- PATH learns an embedding of the past event history via Recurrent Neural Networks that also cater for the diffusion memory
 - (the representation of) users that frequently become active in different streams within a small time interval are closer

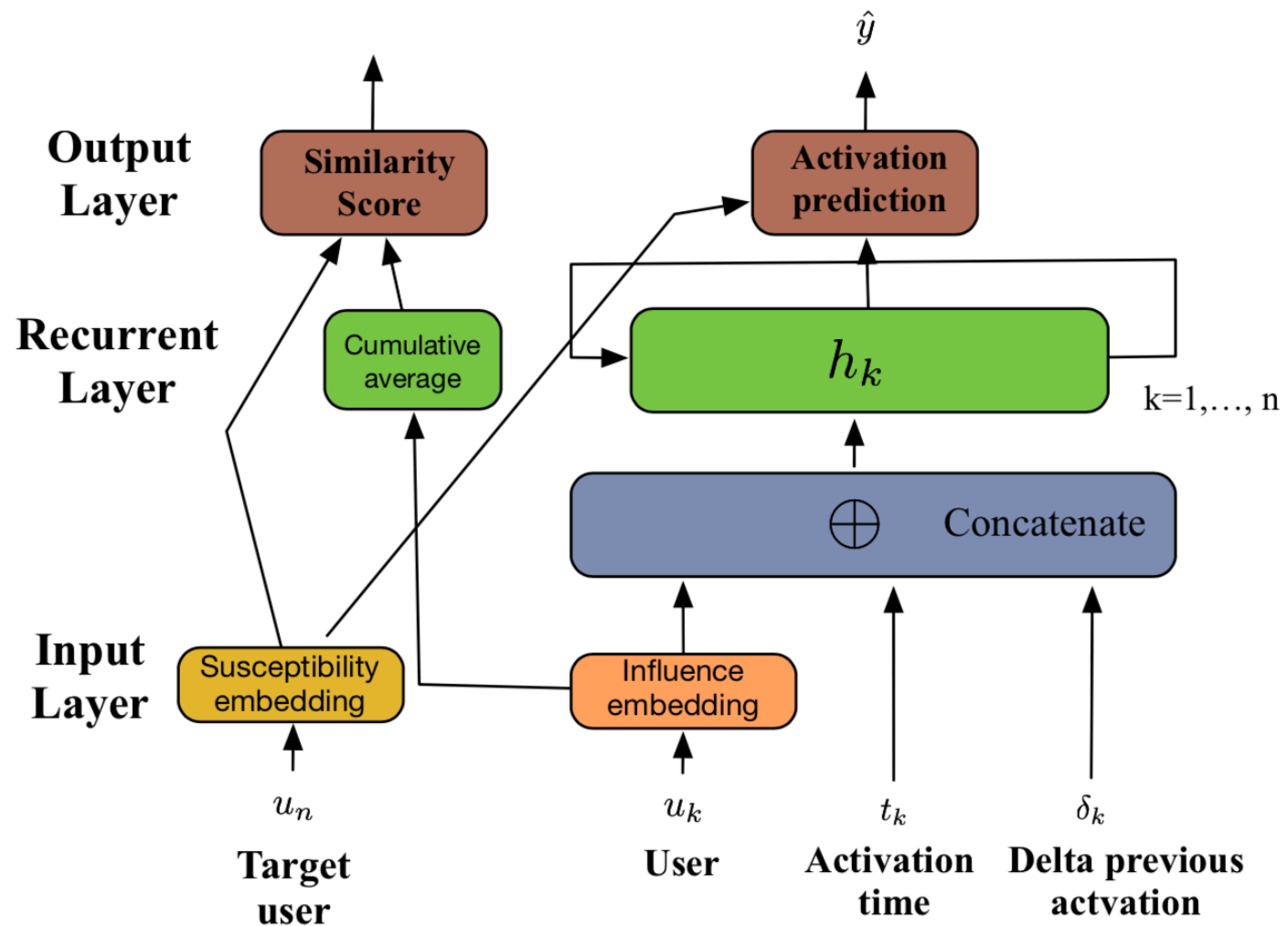
$$\mathbf{a}_k = \mathbf{W}_e \mathbf{u}_k$$

$$\mathbf{s}_k = \mathbf{V}_e \mathbf{u}_k$$

$$\mathbf{h}_k = \text{LSTM}([\mathbf{a}_k, t_k, \delta_k], \mathbf{h}_{k-1})$$

$$\hat{y}_i = \sigma(\mathbf{W}_o[\mathbf{s}_n, \mathbf{h}_n])$$

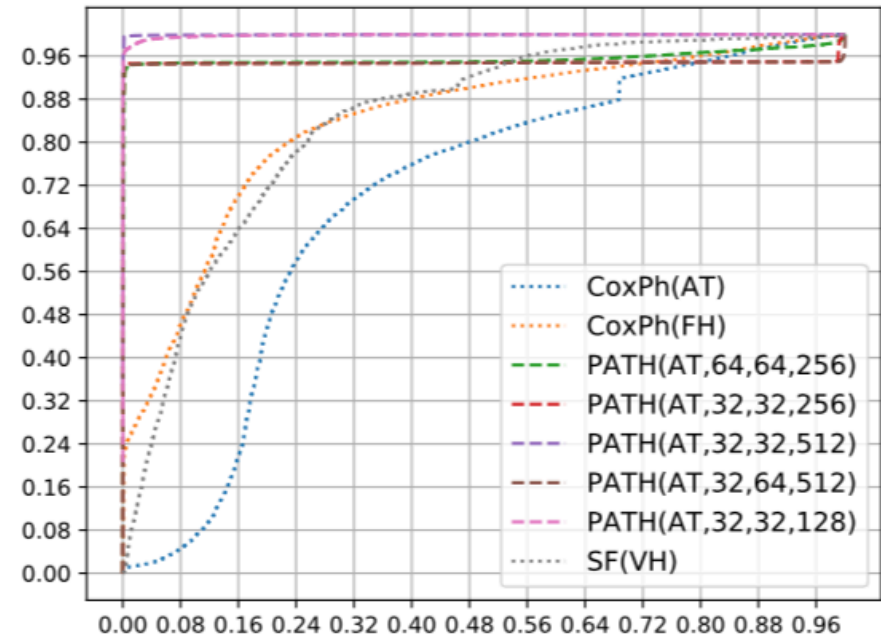
$$\tilde{y}_i = \exp \left\{ - \left\| \mathbf{s}_n - \frac{1}{n} \sum_{k=1}^{n-1} \mathbf{a}_k \right\|^2 \right\}$$



$$\mathcal{L} = \sum_{\substack{\langle \mathcal{H}, y \rangle \in \mathcal{T}_{\mathcal{C}} \\ |\mathcal{H}|=n}} \{y(\gamma \log(\hat{y}) + \beta \log \tilde{y}) + (1-y)(\gamma \log(1-\hat{y}) + \beta \log(1-\tilde{y}))\}$$

Why it works

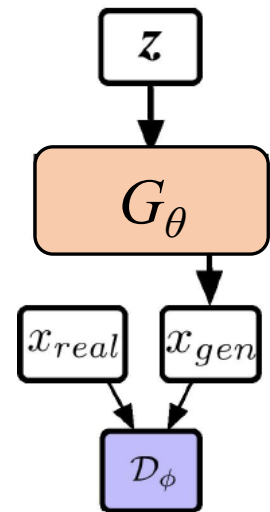
- Checking is easier than generating
 - Autoregressive models try to parameterize \mathbb{P}_r
 - By contrast, discriminative models only compute «closeness» to examples within \mathbb{P}_r



(b) Twitter

Generative Adversarial Networks (GANs)

- **Model:** a deterministic function G_θ is learnt to transform random z into data
 - no assumption on data/latent probability distributions
 - allows for sampling from $P_\theta(x)$ efficiently
 - can produce any distribution $P(x)$ with powerful enough G_θ
- **Learning as a two-player game**
 - *Discriminator* D_ϕ : Trained to optimally discriminate real data from generated samples
 - *Generator* G_θ : Trained to generate realistic samples fooling the discriminator



Is it better than ML?

- Discriminator Loss:

$$L_D(\phi, \theta) = \mathbb{E}_{x \sim \mathbb{P}_r} [\log D_\phi(x)] + \mathbb{E}_{x \sim \mathbb{P}_\theta} [\log(1 - D_\phi)]$$

- Generator loss

$$L_G(\phi, \theta) = L_D(\phi, \theta)$$

- Adversarial Game

$$\max_{\phi} \min_{\theta} \mathbb{E}_{x \sim \mathbb{P}_r} [\log D_\phi(x)] + \mathbb{E}_{x \sim \mathbb{P}_\theta} [\log(1 - D_\phi)]$$

Is it better than ML?

- Optimal discriminator:

$$L_D(\phi, \theta) = \int p_r(x) \log D(x) + p_\theta(x) \log(1 - D(x)) dx$$

- Maximizing the integrand with respect to $D(x)$ gives

$$D^*(x) = \frac{P_r(x)}{P_r(x) + P_\theta(x)}$$

Is it better than ML?

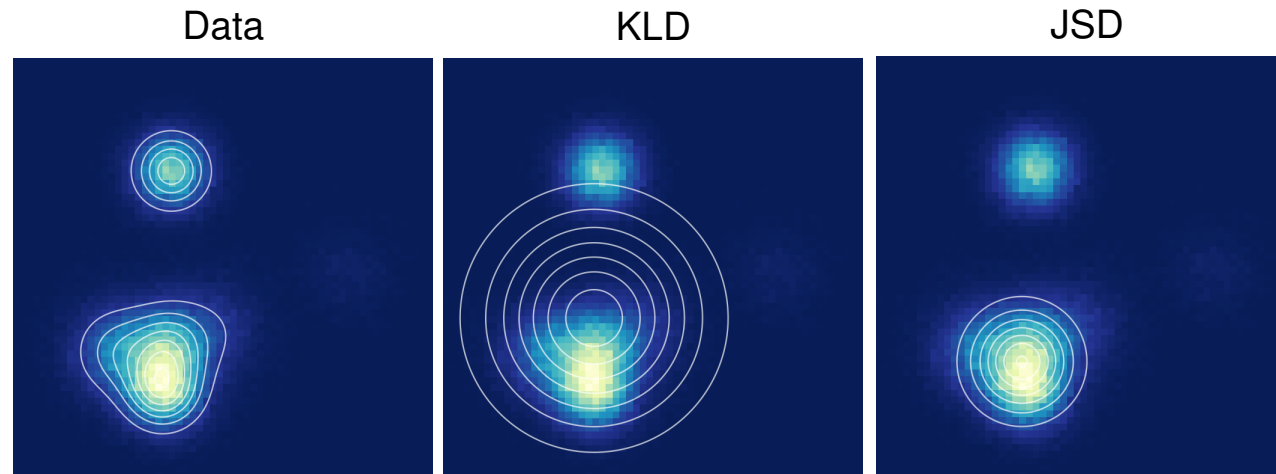
$$\begin{aligned}L_D(D^*, \theta) &= \int p_r(x) \log D(x) + p_\theta(x) \log(1 - D(x)) dx \\&= \int p_r(x) \log \frac{P_r(x)}{P_r(x) + P_\theta(x)} + p_\theta(x) \log \frac{P_\theta(x)}{P_r(x) + P_\theta(x)} dx \\&= KL \left[\mathbb{P}_r \mid \frac{\mathbb{P}_\theta + \mathbb{P}_r}{2} \right] + KL \left[\mathbb{P}_\theta \mid \frac{\mathbb{P}_\theta + \mathbb{P}_r}{2} \right] - 2 \log 2 \\&= 2JS[\mathbb{P}_r \mid \mathbb{P}_\theta] - 2 \log 2\end{aligned}$$

Is it better than ML?

$$L_G(D^*, \theta) \approx JS[\mathbb{P}_r | \mathbb{P}_\theta]$$

- Minimizing with regards to θ is equivalent to minimize the Jensen-Shannon Divergence

$$JS[\mathbb{P}_r | \mathbb{P}_\theta] = \frac{1}{2} KL[\mathbb{P}_r | \mathbb{P}_\theta] + \frac{1}{2} KL[\mathbb{P}_\theta | \mathbb{P}_r]$$



Source: [Theis et al 2016]

GAN learning scheme

Algorithm 1 Inference algorithm.

- 1 Initialize ϕ and θ
- 2 **for** number of epochs **do**
- 3 **for** k steps **do**
- 4 Sample $\{\tilde{x}_\theta^{(1)}, \dots, \tilde{x}_\theta^{(m)}\}$ from \mathbb{P}_θ ;
- 5 Sample $\{x^{(1)}, \dots, x^{(1)}\}$ from \mathbb{P}_r ;
- 6 Update ϕ by ascending its stochastic gradient:

$$\nabla_\phi \frac{1}{m} \sum_{i=1}^m \left[\log \left(D_\phi \left(x^{(i)} \right) \right) + \log \left(1 - D_\phi \left(\tilde{x}_\theta^{(1)} \right) \right) \right]$$

- 7 **end for**
- 8 Sample $\{\tilde{x}_\theta^{(1)}, \dots, \tilde{x}_\theta^{(m)}\}$ from \mathbb{P}_θ ;
- 9 Update θ by descending its stochastic gradient:

$$\nabla_\theta \frac{1}{m} \sum_{i=1}^m \log \left(1 - D_\phi \left(\tilde{x}_\theta^{(i)} \right) \right)$$

- 10 **end for**
 - 11 Return ϕ and θ .
-

GAN learning scheme

- Critical: Backpropagation from samples

Algorithm 1 Inference algorithm.

- 1 Initialize ϕ and θ
- 2 **for** number of epochs **do**
- 3 **for** k steps **do**
- 4 Sample $\{\tilde{x}_\theta^{(1)}, \dots, \tilde{x}_\theta^{(m)}\}$ from \mathbb{P}_θ ;
- 5 Sample $\{x^{(1)}, \dots, x^{(1)}\}$ from \mathbb{P}_r ;
- 6 Update ϕ by ascending its stochastic gradient:

$$\nabla_\phi \frac{1}{m} \sum_{i=1}^m \left[\log \left(D_\phi \left(x^{(i)} \right) \right) + \log \left(1 - D_\phi \left(\tilde{x}_\theta^{(1)} \right) \right) \right]$$

- 7 **end for**
- 8 Sample $\{\tilde{x}_\theta^{(1)}, \dots, \tilde{x}_\theta^{(m)}\}$ from \mathbb{P}_θ ;
- 9 Update θ by descending its stochastic gradient:

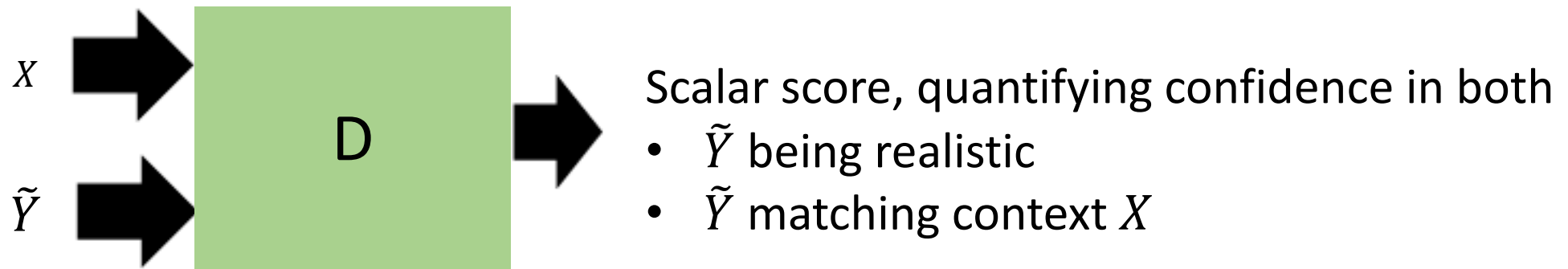
$$\nabla_\theta \frac{1}{m} \sum_{i=1}^m \log \left(1 - D_\phi \left(\tilde{x}_\theta^{(i)} \right) \right)$$

- 10 **end for**
 - 11 Return ϕ and θ .
-

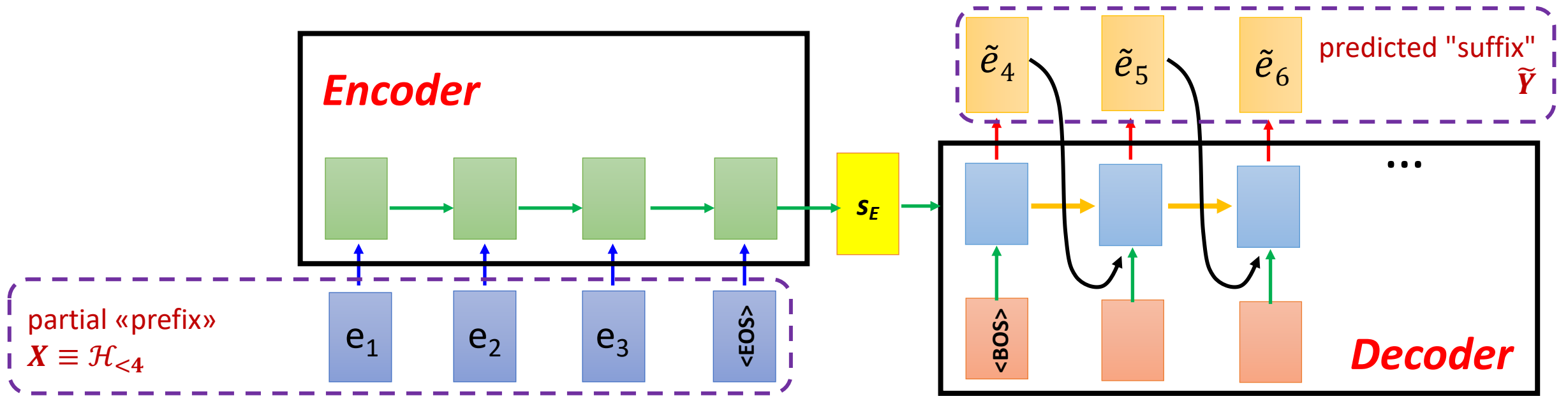
Context: Conditional GAN

Real future events (to be predicted)

$$Y = e_h, e_{h+1}, \dots$$



Basic seq2seq conditional Generator



- Two components:
 - RNN Encoder: maps a partial trace $\mathcal{H}_{<h}$ into a condensed representation s_E
 - RNN Decoder: generates a suffix step-by-step, using s_E as its initial state (or as an additional input for all steps)
- Advantages w.r.t. a single RNN
 - Prefix as a whole has a semantics
 - Correlations within the prefix can influence the prediction of the suffix

GAN learning scheme

- Critical: Backpropagation from samples
 - How do we sample events while preserving backpropagation?

Algorithm 1 Inference algorithm.

- 1 Initialize ϕ and θ
- 2 **for** number of epochs **do**
- 3 **for** k steps **do**
- 4 Sample $\{\tilde{x}_\theta^{(1)}, \dots, \tilde{x}_\theta^{(m)}\}$ from \mathbb{P}_θ ;
- 5 Sample $\{x^{(1)}, \dots, x^{(1)}\}$ from \mathbb{P}_r ;
- 6 Update ϕ by ascending its stochastic gradient:

$$\nabla_\phi \frac{1}{m} \sum_{i=1}^m \left[\log \left(D_\phi \left(x^{(i)} \right) \right) + \log \left(1 - D_\phi \left(\tilde{x}_\theta^{(1)} \right) \right) \right]$$

- 7 **end for**
- 8 Sample $\{\tilde{x}_\theta^{(1)}, \dots, \tilde{x}_\theta^{(m)}\}$ from \mathbb{P}_θ ;
- 9 Update θ by descending its stochastic gradient:

$$\nabla_\theta \frac{1}{m} \sum_{i=1}^m \log \left(1 - D_\phi \left(\tilde{x}_\theta^{(i)} \right) \right)$$

- 10 **end for**
 - 11 Return ϕ and θ .
-

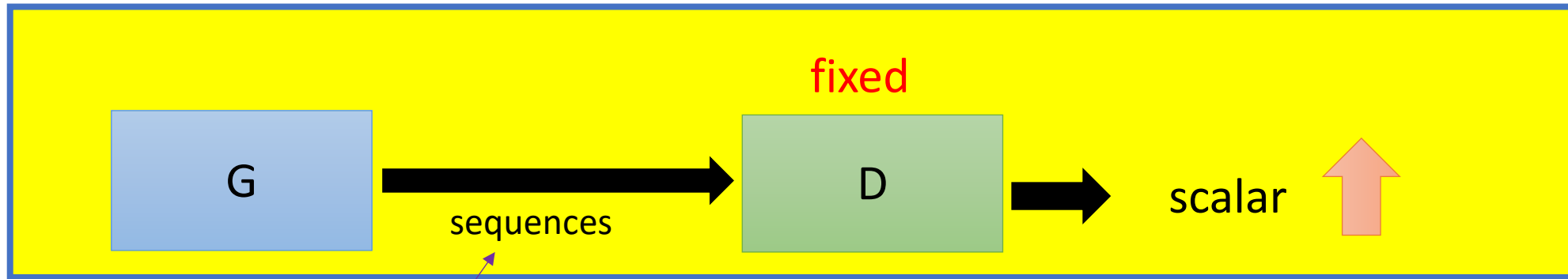
GAN for Sequence Data: Challenges (1)

- Sample $t \sim P_{\theta}(\cdot | \mathcal{H})$
 - Need to make explicit the relationship with θ
 - E.g., Weibull distribution can be reparameterized

$$\begin{aligned} & \text{If } u \sim U(0,1) \\ \text{Then } t &= \alpha(-\log u)^{1/\beta} \sim \text{Weibull}(\alpha, \beta) \end{aligned}$$

- Explicit dependency of t both α and β
- Issue with generalized (neural-based) intensity function $\lambda^*(t)$

GAN for Sequence Data: Challenges (2)



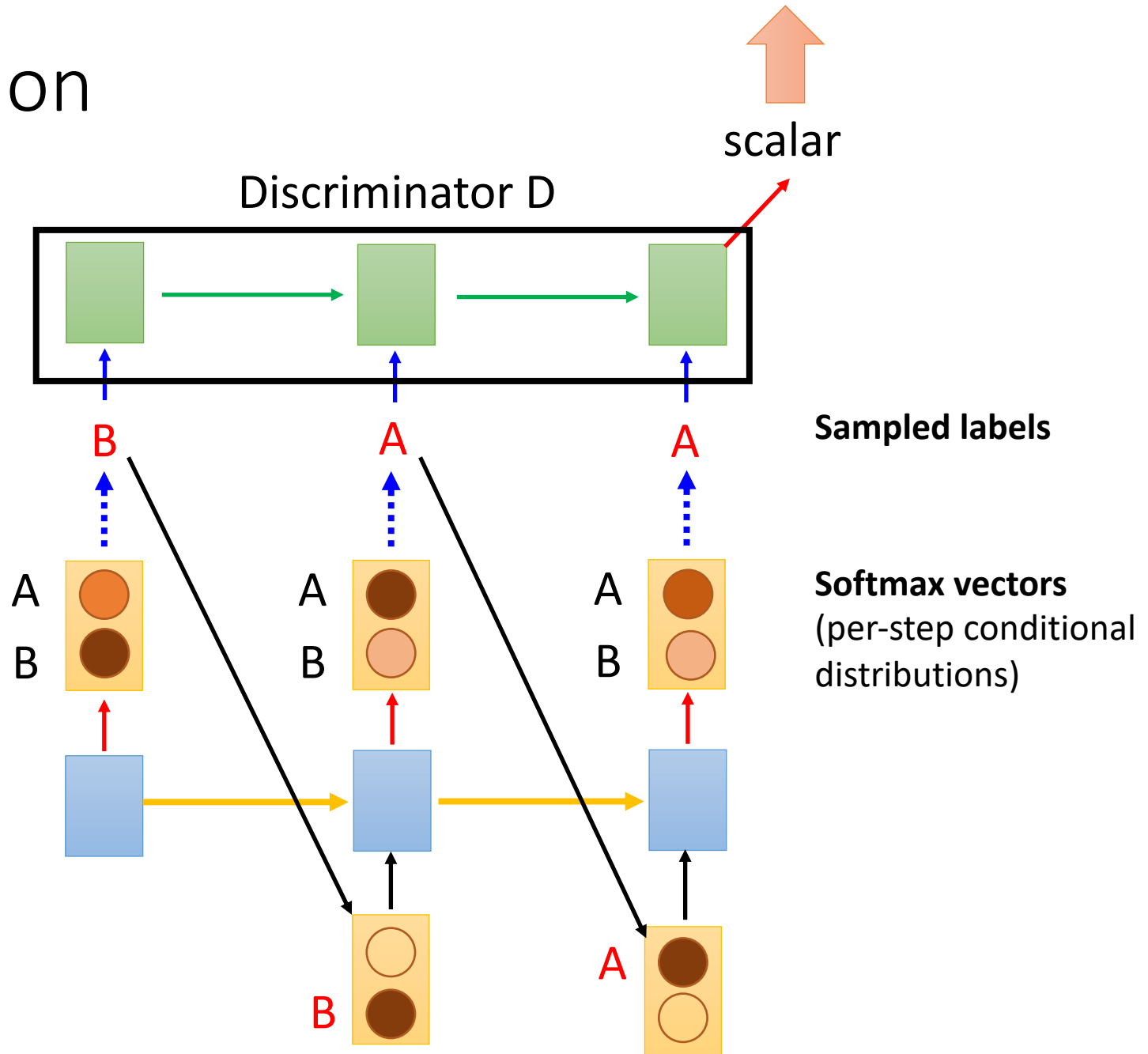
Sequence generation with an RNN decoder entails sampling from a discrete distribution: at every time step the most probable event type is chosen from a softmax output.

loss gradients cannot backpropagate to the Generator!

Discrete Event Generation

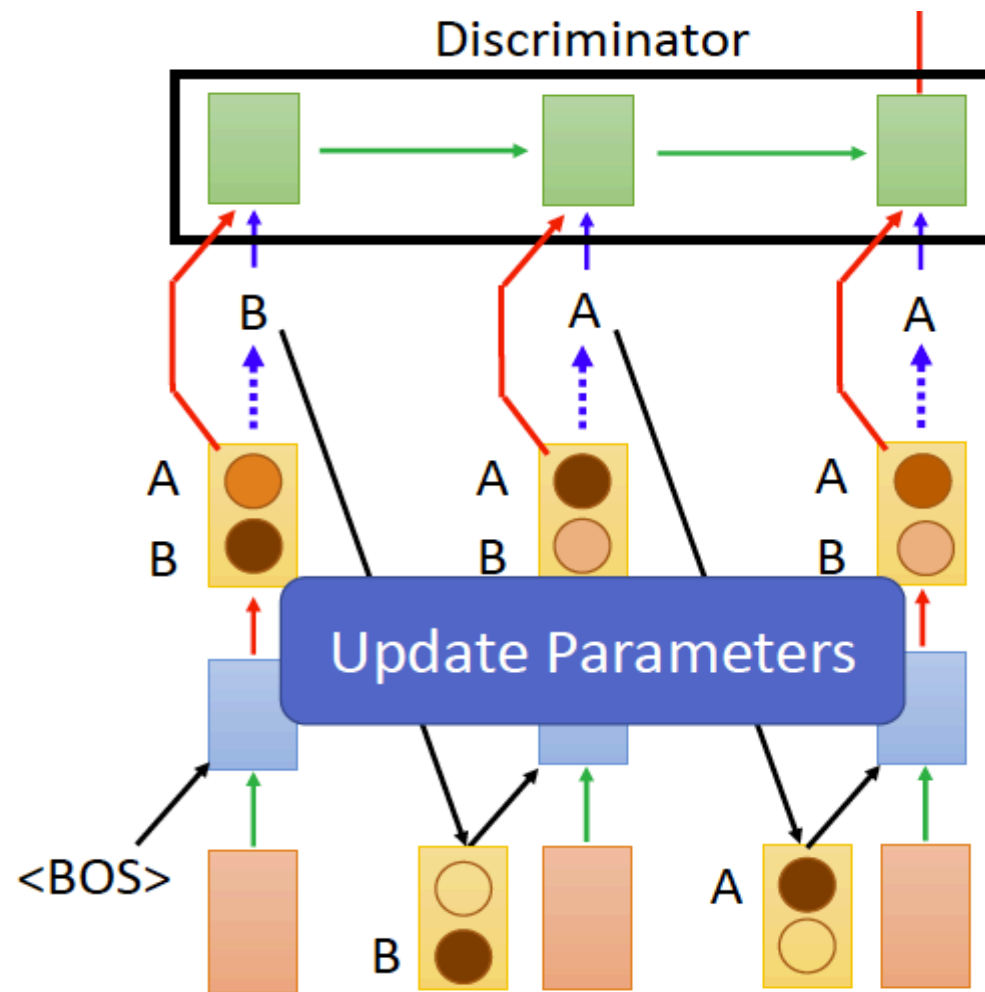
Discrete event attributes a_j
predicted for future events

Assume possible event types are
represented by labels **A, B, ...**



Dealing with continuous inputs

- Use a softmax distribution as input to the Discriminator
- No direct sampling channel



... but sequences of softmax vectors are easily recognized

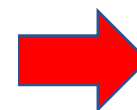
Real
sequence

1	0	0	0	0
0	1	0	0	0
0	0	1	0	0
0	0	0	1	0
0	0	0	0	1

Generated
sequence

0.9	0.1	0.1	0	0
0.1	0.9	0.1	0	0
0	0	0.7	0.1	0
0	0	0.1	0.8	0.1
0	0	0	0.1	0.9

Discriminator can
easily tell them apart



The GAN gets stuck in an
equilibrium that is not
good for the Generator

- WGAN and/or adversarial training over abstract features can reduce this risk
- Also using an (Approximated) Embedding layer to map each events to a dense representation
 - softmax components are used as weights for combining the embedding vector of the corresponding event labels [Xu et al., ENMLP'17]

Approximated Embedding Layer (AEL)

Embedding Layer

Real
sequence

1	0	0	0	0
0	1	0	0	0
0	0	1	0	0
0	0	0	1	0
0	0	0	0	1

Generated
sequence

0.9	0.1	0.1	0	0
0.1	0.9	0.1	0	0
0	0	0.7	0.1	0
0	0	0.1	0.8	0.1
0	0	0	0.1	0.9

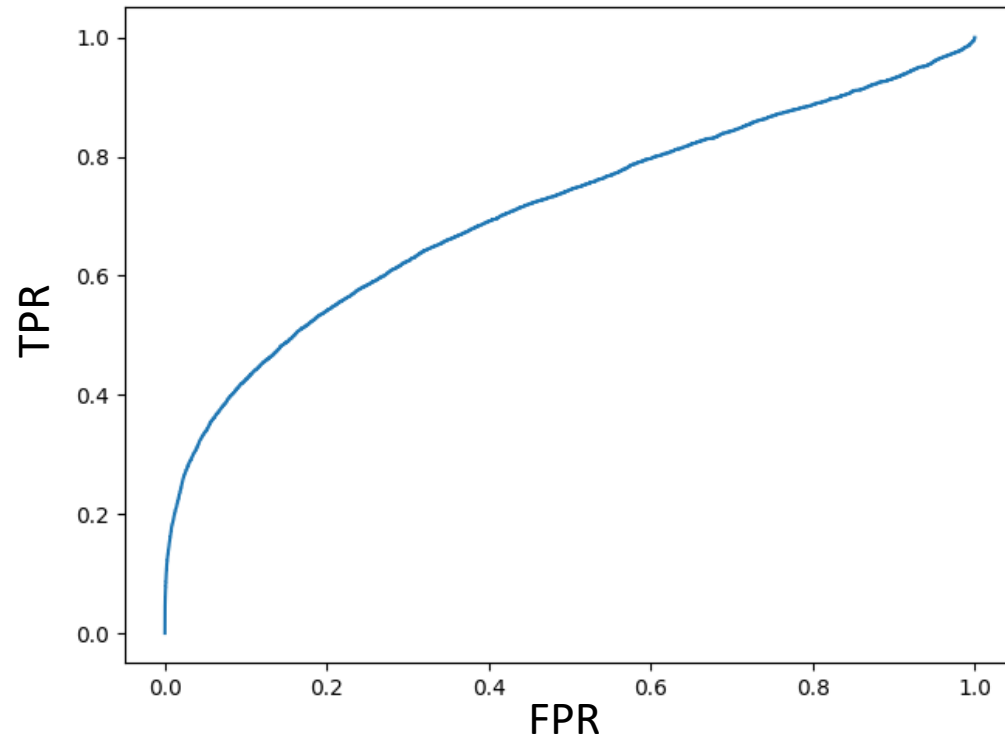
$\begin{bmatrix} 0.81 \\ 0.74 \\ 0.10 \\ 0.40 \\ 0.59 \\ 0.57 \\ 0.18 \end{bmatrix}$

$\begin{bmatrix} 0.81 & 0.82 & 0.59 & 0.57 & 0.18 \\ 0.74 & 0.10 & 0.40 & 0.42 & 0.10 \end{bmatrix}$

$\begin{bmatrix} 0.81 & 0.82 & 0.53 & 0.63 & 0.13 \\ 0.82 & 0.02 & 0.38 & 0.47 & 0.06 \end{bmatrix}$

Harsh time for
the discriminator

Activations on Twitter, Weibull model

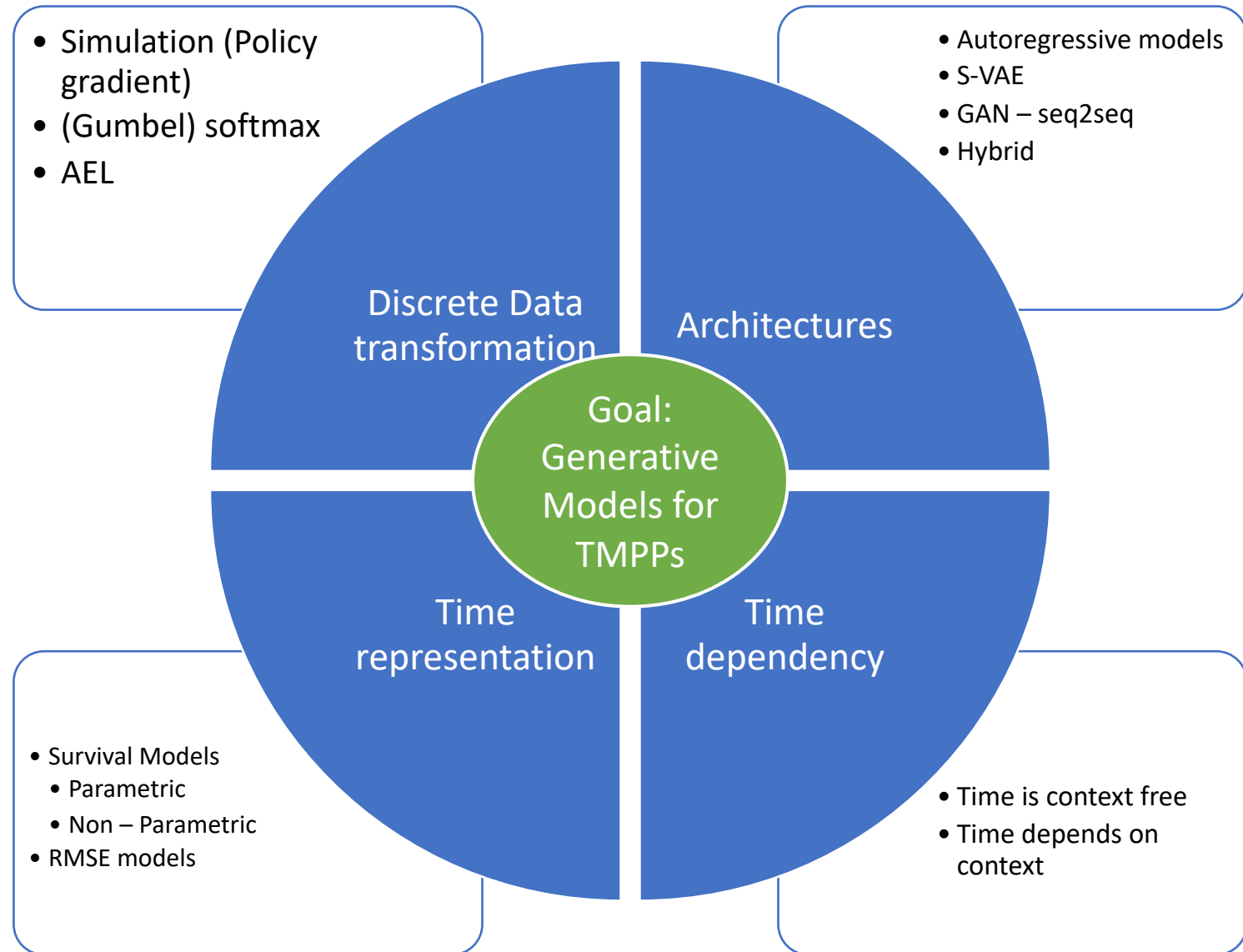


- Noisy temporal patterns, data collection bias

Conclusion and open issues

- Research goal: learn generative models for complex processes (regarded as TMPPs)
 - Two tasks: generate new traces, and predict future events for unfinished traces
- Limitations of ML-based approaches (autoregressive models, VAEs)
- GANs for MMPs: basic scheme
 - Generator **G**: RNN Generator, Seq2Seq Generator
 - Discriminator **D** as a ranker/critic
 - **G** passes to **D** differentiable approximations for the generated event sequences
- Open challenges (and other design choices)...

Research line: Summary and Open Issues



References

- H. Mei, J. Eisner. The Neural Hawkes Process: A Neurally Self-Modulating Multivariate Point Process. NeurIPS 2017.
- N. Du, H. Dai, R. Trivedi, U. Upadhyay, M. Gomez-Rodriguez, L. Song. Recurrent Marked Temporal Point Processes: Embedding Event History to Vector. KD 2016
- Z. Xu, B. Liu, B. Wang, C. Sun, X. Wang, Z. Wang, C. Qi. Neural Response Generation via GAN with an Approximate Embedding Layer. EMNLP 2017
- L. Theis, A. van der Oord, M. Bethge. A note on the evaluation of generative models. ICLR 2016
- M. Arjovsky, L. Bottou. Towards Principled Methods for Training Generative Adversarial Networks. ICLR 2017
- F. Huszar. How (not to) Train your generative model: Scheduled sampling, Likelihood, adversary?. arXiv:1511.05101v1 [stat.ML] 16 Nov 2015
- I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio. Generative Adversarial Networks. NeurIPS 2014.

Thank you

Questions?



giuseppe.manco@icar.cnr.it



@beman

