Adversarial Attacks on Probabilistic Autoregressive Forecasting Models

ICML 2020



Raphaël Dang-Nhu



Gagandeep Singh



Pavol Bielik



Martin Vechev

Department of Computer Science, ETH Zürich

dangnhur@ethz.ch



(i) Probabilistic forecasting model

¹Blundell et al., Weight Uncertainty in Neural Networks, ICML 2015



(i) Probabilistic forecasting model



(ii) Bayesian neural network

¹Blundell et al., Weight Uncertainty in Neural Networks, ICML 2015



(i) Probabilistic forecasting model

(ii) Bayesian neural network

 $\cdot\,$ Multiple sources of noise: (i) each timestep, (ii) each weight^1

¹Blundell et al., Weight Uncertainty in Neural Networks, ICML 2015



(i) Probabilistic forecasting model



(ii) Bayesian neural network

- Multiple sources of noise: (i) each timestep, (ii) each weight¹
- Complex resulting output distribution, approximated via Monte-Carlo sampling

¹Blundell et al., Weight Uncertainty in Neural Networks, ICML 2015

Focus of this work: probabilistic forecasting models

- Stochastic sequence model
- Generates several prediction traces



Focus of this work: probabilistic forecasting models

- Stochastic sequence model
- Generates several prediction traces



Traditionally used as a generative model



WaveNet for raw audio



Handwriting generation

Probabilistic forecasting models for decision-making²

- Allows to predict volatility of the time-series.
- Useful with low signal-to-noise ratio.

Key idea: use generated traces as Monte-Carlo samples to estimate the evolution of the time-series

²Salinas et al., DeepAR: Probabilistic forecasting with autoregressive recurrent networks, International Journal of Forecasting, 2020

Probabilistic forecasting models for decision-making²

- Allows to predict volatility of the time-series.
- Useful with low signal-to-noise ratio.

Key idea: use generated traces as Monte-Carlo samples to estimate the evolution of the time-series



Stock prices



Electricity consumption



Business sales

Integrated in Amazon Sagemaker (DeepAR architecture)

²Salinas et al., DeepAR: Probabilistic forecasting with autoregressive recurrent networks, International Journal of Forecasting, 2020

 $\cdot\,$ New class of attack objectives based on output statistics

- New class of attack objectives based on output **statistics**
- Adaptation of gradient-based adversarial attacks to these new attack objectives for stochastic models

- $\cdot\,$ New class of attack objectives based on output statistics
- Adaptation of gradient-based adversarial attacks to these new attack objectives for stochastic models
- Main technical aspect: developing estimators for propagating the objective gradient through the Monte-Carlo approximation

- $\cdot\,$ New class of attack objectives based on output statistics
- Adaptation of gradient-based adversarial attacks to these new attack objectives for stochastic models
- Main technical aspect: developing estimators for propagating the objective gradient through the Monte-Carlo approximation

- New class of attack objectives based on output statistics
- Adaptation of gradient-based adversarial attacks to these new attack objectives for stochastic models
- Main technical aspect: developing estimators for propagating the objective gradient through the Monte-Carlo approximation

We aim at providing an off-the-shelf methodology for these attacks

Class of attack objectives

Stochastic model with input *x*, and output $y \sim q_x(\cdot)$. Previously considered attack objectives: Stochastic model with input *x*, and output $y \sim q_x(\cdot)$. Previously considered attack objectives:

Untargeted attacks on information divergence *D* with the original predicted distribution

 $\max_{\delta} D\left(q_{x+\delta} \| q_x\right)$



Stochastic model with input *x*, and output $y \sim q_x(\cdot)$. Previously considered attack objectives:

Untargeted attacks on information divergence *D* with the original predicted distribution

 $\max_{\delta} D\left(q_{x+\delta} \| q_x\right)$





Untargeted/Targeted attacks on the mean of the distribution

 $\min_{\delta} \operatorname{distance} \left(\mathbb{E}_{q_{x+\delta}}[y], \operatorname{target} \right)$

We perform a targeted attack on a **statistic** $\chi(y)$ of the output.

We perform a targeted attack on a **statistic** $\chi(y)$ of the output. This corresponds to minimizing

distance $(\mathbb{E}_{q_{x+\delta}}[\chi(y)], \text{target})$

We perform a targeted attack on a **statistic** $\chi(y)$ of the output. This corresponds to minimizing

distance $(\mathbb{E}_{q_{x+\delta}}[\chi(y)], \text{target})$

Extensions:

- Bayesian setting $q_x(y|z)$.
- Generalization to simultaneous attack of several statistics.
- Statistics depending on *x*.

Motivation 1: option pricing in finance

Consider a stock with

- past prices $x = (p_1, \ldots, p_{t-1})$
- predicted future prices $y = (p_t, \ldots, p_T)$.

Motivation 1: option pricing in finance

Consider a stock with

- past prices $x = (p_1, \ldots, p_{t-1})$
- predicted future prices $y = (p_t, \ldots, p_T)$.

Name	$\chi(y)$	Observation z
European call option	$max(0, y_h)$	
Asian call option	$average_i(y_i)$	
Limit sell order	$\mathbb{1}\left[\max_{i} y_{i} \geq \text{threshold}\right]$	
Barrier option	У _h	$\max_i y_i \geq \text{threshold}$

Motivation 1: option pricing in finance

Consider a stock with

- past prices $x = (p_1, \ldots, p_{t-1})$
- predicted future prices $y = (p_t, \ldots, p_T)$.

Name	$\chi(y)$	Observation z
European call option	$max(0, y_h)$	
Asian call option	$average_i(y_i)$	
Limit sell order	$\mathbb{1}\left[\max_{i} y_{i} \geq \text{threshold}\right]$	
Barrier option	У _h	$\max_i y_i \geq \text{threshold}$

Our framework allows to specifically target one of these options

Motivation 2: attacking model uncertainty

Some defenses use **prediction uncertainty** to detect adversarial examples.

Some defenses use **prediction uncertainty** to detect adversarial examples.

New attacks bypass these defenses by enforcing **uncertainty constraints** for the adversarial example.

Some defenses use **prediction uncertainty** to detect adversarial examples.

New attacks bypass these defenses by enforcing **uncertainty constraints** for the adversarial example.

Our framework allows to express these constraints, with

- The entropy $\mathbb{E}_{q_x}[-\log(q[y|x])].$
- The distribution's moments $\mathbb{E}_{q_x}[y^k]$.

Details about the estimators

Gradient-based attacks require computing

$$abla_{\delta} \mathbb{E}_{q[\boldsymbol{y}|\boldsymbol{x}+\boldsymbol{\delta},z]}[\chi(\boldsymbol{y})]$$

Gradient-based attacks require computing

$$abla_{\delta} \mathbb{E}_{q[\boldsymbol{y}|\boldsymbol{x}+\boldsymbol{\delta},z]}[\chi(\boldsymbol{y})]$$

The expectation and its gradient have no analytical closed form

Gradient-based attacks require computing

$$abla_{\delta} \mathbb{E}_{q[\boldsymbol{y}|\boldsymbol{x}+\boldsymbol{\delta},z]}[\chi(\boldsymbol{y})]$$

The expectation and its gradient have no analytical closed form

We provide two different estimators to approximate the gradient

Approach 1: REINFORCE

- A.k.a as log-derivative trick and score-function estimator.
- Based on interversion of expectation and derivative.

Approach 1: REINFORCE

- A.k.a as log-derivative trick and score-function estimator.
- Based on interversion of expectation and derivative.

$$\nabla_{\boldsymbol{\delta}} \mathbb{E}_{q[\boldsymbol{y}|\boldsymbol{x}+\boldsymbol{\delta},z]}[\chi(\boldsymbol{y})] \\ \simeq \frac{\sum_{l=1}^{L} \chi(\boldsymbol{y}^{l}) q[z|\boldsymbol{x}+\boldsymbol{\delta},\boldsymbol{y}^{l}] \nabla_{\boldsymbol{\delta}} \log(q[\boldsymbol{y}^{l}|\boldsymbol{x}+\boldsymbol{\delta},z])}{\sum_{l=1}^{L} q[z|\boldsymbol{x}+\boldsymbol{\delta},\boldsymbol{y}^{l}]}$$

REINFORCE estimator

Approach 2: Reparametrization

- Mitigates the high-variance of REINFORCE.
- Typically used for variational inference.
- Assumes a reparametrization $y \sim g(x, \eta)$, where g is deterministic.

Approach 2: Reparametrization

- Mitigates the high-variance of REINFORCE.
- Typically used for variational inference.
- Assumes a reparametrization y ~ g(x, η), where g is deterministic.

$$\nabla_{\boldsymbol{\delta}} \mathbb{E}_{q[\boldsymbol{y}|\boldsymbol{x}+\boldsymbol{\delta},z]}[\chi(\boldsymbol{y})] \\ \simeq \nabla_{\boldsymbol{\delta}} \left(\frac{\sum_{l=1}^{L} \chi(g_{\boldsymbol{x}}(\boldsymbol{\delta},\boldsymbol{\eta}^{l}))q[z|\boldsymbol{x}+\boldsymbol{\delta},g_{\boldsymbol{x}}(\boldsymbol{\delta},\boldsymbol{\eta}^{l})]}{\sum_{l=1}^{L} q[z|\boldsymbol{x}+\boldsymbol{\delta},g_{\boldsymbol{x}}(\boldsymbol{\delta},\boldsymbol{\eta}^{l})]} \right)$$

Reparametrization estimator

Respective advantages of gradient estimators.

Method	REINFORCE	Reparametrization
Applies to non-differentiable statistics Requires no reparametrization Applies to Bayesian setting Yields best gradient estimates	22	22

Respective advantages of gradient estimators.

Method REINFORCE Reparametrization Applies to non-differentiable statistics ✓ Requires no reparametrization ✓			
Applies to non-differentiable statistics	Method	REINFORCE	Reparametrization
Applies to Bayesian setting	Applies to non-differentiable statistics Requires no reparametrization Applies to Bayesian setting Yields best gradient estimates	~	~

Detailed comparison and conditions in the paper!

Experimental evaluation

Algorithmic trading scenario, standard additive threat model, maximum Euclidean norm of 0.1³ for the perturbation.

³Corresponds to perturbing one value by 10%, 10 values by 3.3%, 100 values by 1%.

Algorithmic trading scenario, standard additive threat model, maximum Euclidean norm of 0.1³ for the perturbation.

• Attack is successful on 90% of test inputs.

³Corresponds to perturbing one value by 10%, 10 values by 3.3%, 100 values by 1%. ¹⁴

Algorithmic trading scenario, standard additive threat model, maximum Euclidean norm of 0.1³ for the perturbation.

- Attack is successful on 90% of test inputs.
- The network incurs a daily financial loss of -13%.

³Corresponds to perturbing one value by 10%, 10 values by 3.3%, 100 values by 1%. ¹⁴

Algorithmic trading scenario, standard additive threat model, maximum Euclidean norm of 0.1³ for the perturbation.

- Attack is successful on 90% of test inputs.
- The network incurs a daily financial loss of -13%.

³Corresponds to perturbing one value by 10%, 10 values by 3.3%, 100 values by 1%.

Algorithmic trading scenario, standard additive threat model, maximum Euclidean norm of 0.1³ for the perturbation.

- Attack is successful on 90% of test inputs.
- The network incurs a daily financial loss of -13%.



³Corresponds to perturbing one value by 10%, 10 values by 3.3%, 100 values by 1%. ¹⁴

Original test samples (red) and adversarial examples (blue) for prediction of electricity consumption.



Code and trained models are available at

github.com/eth-sri/
probabilistic-forecasts-attacks

Contact at

dangnhur@ethz.ch