

Adversarial Attacks on Probabilistic Autoregressive Forecasting Models

ICML 2020



Raphaël Dang-Nhu



Gagandeep Singh



Pavol Bielik

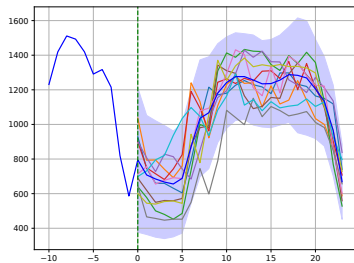


Martin Vechev

Department of Computer Science, ETH Zürich

dangnhur@ethz.ch

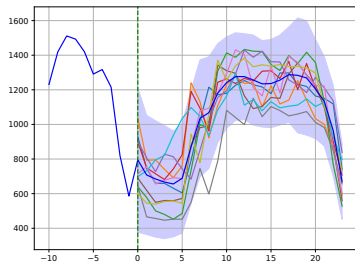
Neural architectures with stochastic behavior



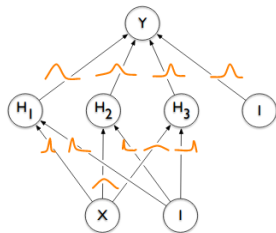
(i) Probabilistic forecasting model

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

Neural architectures with stochastic behavior



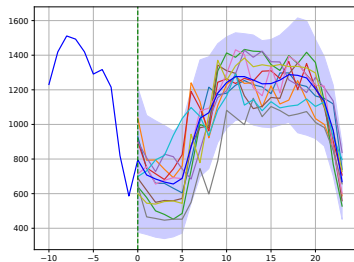
(i) Probabilistic forecasting model



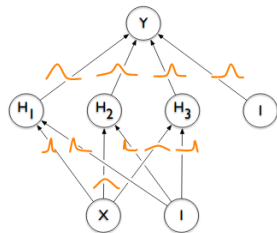
(ii) Bayesian neural network

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

Neural architectures with stochastic behavior



(i) Probabilistic forecasting model

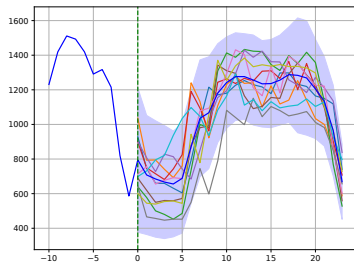


(ii) Bayesian neural network

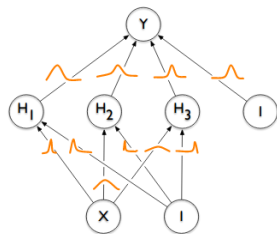
- Multiple sources of noise: (i) each timestep, (ii) each weight¹

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

Neural architectures with stochastic behavior



(i) Probabilistic forecasting model



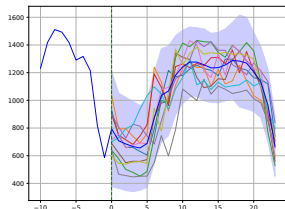
(ii) Bayesian neural network

- 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

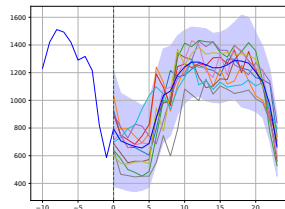
Focus of this work: probabilistic forecasting models

- Stochastic sequence model
- Generates several prediction traces

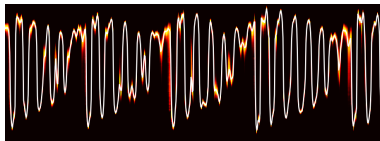


Focus of this work: probabilistic forecasting models

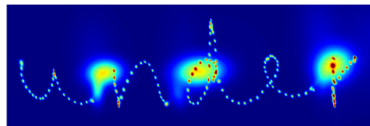
- Stochastic sequence model
- Generates several prediction traces



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

- 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

- 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

- 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

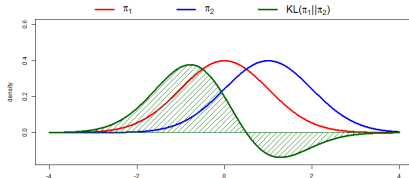
Stochastic model with input x , and output $y \sim q_x(\cdot)$.
Previously considered attack objectives:

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

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

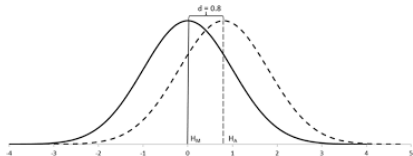
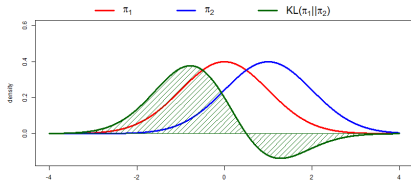


Stochastic model with input x , and output $y \sim q_x(\cdot)$.

Previously considered attack objectives:

Untargeted attacks on information divergence D with the original predicted distribution

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



Untargeted/Targeted attacks on the mean of the distribution

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

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

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

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

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

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

Motivation 1: option pricing in finance

Consider a stock with

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

Name	$\chi(y)$	Observation z
European call option	$\max(0, y_h)$	
Asian call option	$\text{average}_i(y_i)$	
Limit sell order	$\mathbb{1}[\max_i y_i \geq \text{threshold}]$	
Barrier option	y_h	$\max_i y_i \geq \text{threshold}$

Motivation 1: option pricing in finance

Consider a stock with

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

Name	$\chi(y)$	Observation z
European call option	$\max(0, y_h)$	
Asian call option	$\text{average}_i(y_i)$	
Limit sell order	$\mathbb{1} [\max_i y_i \geq \text{threshold}]$	
Barrier option	y_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.

Motivation 2: attacking model uncertainty

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

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

Motivation 2: attacking model uncertainty

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

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

Gradient-based attacks require computing

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

The expectation and its gradient have no analytical closed form

Technical challenge

Gradient-based attacks require computing

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

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.

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

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 $y \sim g(x, \eta)$, where g is deterministic.

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

Reparametrization estimator

Comparison

Respective advantages of gradient estimators.

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

Comparison

Respective advantages of gradient estimators.

Method	REINFORCE	Reparametrization
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

Experiments: stock prices

Algorithmic trading scenario, standard additive threat model, maximum Euclidean norm of 0.1^3 for the perturbation.

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

Experiments: stock prices

Algorithmic trading scenario, standard additive threat model, maximum Euclidean norm of 0.1^3 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%.

Experiments: stock prices

Algorithmic trading scenario, standard additive threat model, maximum Euclidean norm of 0.1^3 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%.

Experiments: stock prices

Algorithmic trading scenario, standard additive threat model, maximum Euclidean norm of 0.1^3 for the perturbation.

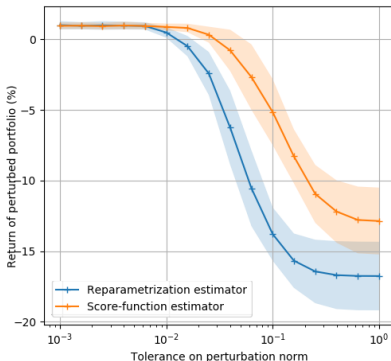
- Attack is successful on 90% of test inputs.
- The network incurs a daily financial loss of -13% .

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

Experiments: stock prices

Algorithmic trading scenario, standard additive threat model, maximum Euclidean norm of 0.1^3 for the perturbation.

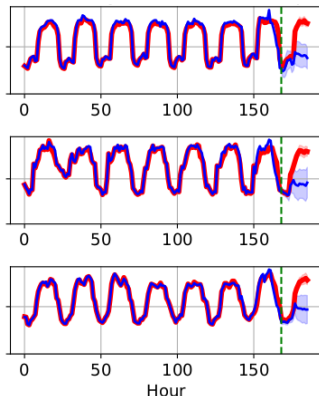
- Attack is successful on 90% of test inputs.
- The network incurs a daily financial loss of -13% .



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

Experiments: electricity

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



Thanks for listening

Code and trained models are available at

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

Contact at

`dangnhur@ethz.ch`