

# Recent Works on Feature Interaction

Xingzhi Sun

<https://github.com/xingzhis/XAI>

Wednesday 30<sup>th</sup> September, 2020

# Table of Contents

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

## 1 Definitions of feature interaction

- Post-hoc statistic
- Model-specific

## 2 Models

- Tree-based ensembles
- GA<sup>2</sup>M
- Neuron Networks
- Factorization Machines
- Hybrid models with neuron networks

# Table of Contents

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

## 1 Definitions of feature interaction

- Post-hoc statistic
- Model-specific

## 2 Models

- Tree-based ensembles
- GA<sup>2</sup>M
- Neuron Networks
- Factorization Machines
- Hybrid models with neuron networks

# Non-additiveness

Post-hoc statistic for interaction

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

## Non-additiveness

$F(\mathbf{x})$  cannot be written in the form of  $F(\mathbf{x}) = f_j(x_j) + f_{\setminus j}(\mathbf{x}_{\setminus j})$

# Non-additiveness

Post-hoc statistic for interaction

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

## Non-additiveness

$F(\mathbf{x})$  cannot be written in the form of  $F(\mathbf{x}) = f_j(x_j) + f_{\setminus j}(\mathbf{x}_{\setminus j})$

the effect of both variables

vs

the sum of effects of each variable

# Non-additiveness

Post-hoc statistic for interaction

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization  
Machines

Hybrid models with  
neuron networks

Summary

## Non-additiveness

$F(\mathbf{x})$  cannot be written in the form of  $F(\mathbf{x}) = f_j(x_j) + f_{\setminus j}(\mathbf{x}_{\setminus j})$

the effect of both variables

vs

the sum of effects of each variable

let  $F_s(\mathbf{x}_s) = \mathbb{E}_{\mathbf{x}_{\setminus s}} [F(\mathbf{x}_s, \mathbf{x}_{\setminus s})]$

## Predictive learning via rule ensembles

$$H_{jk}^2 = \sum_{i=1}^N \left[ \hat{F}_{jk}(x_{ij}, x_{ik}) - \hat{F}_j(x_{ij}) - \hat{F}_k(x_{ik}) \right]^2 / \sum_{i=1}^N \hat{F}_{jk}^2(x_{ij}, x_{ik})$$

# Model Expressiveness

Post-hoc statistic for interaction

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

the **accurate** model that contains the interaction

VS

the **accurate** model that bans the interaction

# Model Expressiveness

Post-hoc statistic for interaction

the **accurate** model that contains the interaction

VS

the **accurate** model that bans the interaction

$F^*(x)$ : Target function

$F(x)$ : Highly accurate model

$R_{ij}(x)$ : Highly accurate but ban interaction between  $x_i$  and  $x_j$

$$\text{stRMSE}(F(\mathbf{x})) = \frac{\text{RMSE}(F(\mathbf{x}))}{\text{StD}(F^*(\mathbf{x}))},$$

Detecting statistical interactions with additive groves of trees

$$I_{ij}(F(\mathbf{x})) = \text{stRMSE}(F(\mathbf{x})) - \text{stRMSE}(R_{ij}(\mathbf{x}))$$

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization  
Machines

Hybrid models with  
neuron networks

Summary



# Hessian - second derivative

Post-hoc statistic for interaction

The prediction has non-zero hessian over the interaction variables.

$$\frac{\partial F(\mathbf{x})}{\partial x_i \partial x_j} \neq 0$$

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

# Hessian - second derivative

Post-hoc statistic for interaction

The prediction has non-zero hessian over the interaction variables.

$$\frac{\partial F(\mathbf{x})}{\partial x_i \partial x_j} \neq 0$$

## Learning Global Pairwise Interactions with Bayesian Neural Networks

$$\text{EAH}_g^{i,j}(\mathbf{W}) = \mathbb{E}_{p(\mathbf{x})} \left[ \left| \frac{\partial^2 g^{\mathbf{W}}(\mathbf{x})}{\partial x_i \partial x_j} \right| \right]$$
$$\text{AEH}_g^{i,j}(\mathbf{W}) = \left| \mathbb{E}_{p(\mathbf{x})} \left[ \frac{\partial^2 g^{\mathbf{W}}(\mathbf{x})}{\partial x_i \partial x_j} \right] \right|$$

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic  
Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization  
Machines

Hybrid models with  
neuron networks

Summary

# Table of Contents

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

## 1 Definitions of feature interaction

- Post-hoc statistic
- Model-specific

## 2 Models

- Tree-based ensembles
- GA<sup>2</sup>M
- Neuron Networks
- Factorization Machines
- Hybrid models with neuron networks

# Model-specific depiction of interaction

## Recent Works on Feature Interaction

Xingzhi Sun

## Outline

## Definitions of feature interaction

Post-hoc statistic

## Model-specific

## Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

## Summary

- Explicit: parameters, such as weights of interaction terms.
- Implicit: Neuron networks, embedding vectors, etc.

# Table of Contents

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

## 1 Definitions of feature interaction

- Post-hoc statistic
- Model-specific

## 2 Models

- Tree-based ensembles
- GA<sup>2</sup>M
- Neuron Networks
- Factorization Machines
- Hybrid models with neuron networks

# Predictive learning via rule ensembles

## Tree-based ensembles

- Main effects: Linear
- Interaction: Rules derived from decision trees

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

**Tree-based ensembles**

$GA^2M$

Neuron Networks

Factorization

Machines

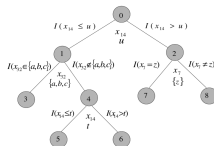
Hybrid models with  
neuron networks

Summary

# Predictive learning via rule ensembles

## Tree-based ensembles

- Main effects: Linear
- Interaction: Rules derived from decision trees



$$r_1(\mathbf{x}) = I(x_{14} \leq u),$$

$$r_4(\mathbf{x}) = I(x_{14} \leq u) \cdot I(x_{32} \notin \{a, b, c\}),$$

$$r_6(\mathbf{x}) = I(x_{14} < u) \cdot I(x_{32} \notin \{a, b, c\}),$$

$$r_7(\mathbf{x}) = I(x_{14} > u) \cdot I(x_7 = z).$$

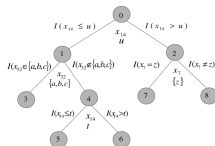
$$r_m(\mathbf{x}) = \prod_{s_{jm} \neq s_j} I(x_j \in s_{jm})$$

Figure: A decision tree and its corresponding rule term

# Predictive learning via rule ensembles

## Tree-based ensembles

- Main effects: Linear
- Interaction: Rules derived from decision trees



$$r_1(\mathbf{x}) = I(x_{14} \leq u),$$

$$r_4(\mathbf{x}) = I(x_{14} \leq u) \cdot I(x_{32} \notin \{a, b, c\}),$$

$$r_6(\mathbf{x}) = I(t < x_{14} \leq u) \cdot I(x_{32} \notin \{a, b, c\}),$$

$$r_7(\mathbf{x}) = I(x_{14} > u) \cdot I(x_7 = z).$$

$$r_m(\mathbf{x}) = \prod_{s_{jm} \neq s_j} I(x_j \in s_{jm})$$

Figure: A decision tree and its corresponding rule term

RuleFit (Friedman and Popescu, 2008)

$$F(\mathbf{x}) = \hat{a}_0 + \sum_{k=1}^K \hat{a}_k r_k(\mathbf{x}) + \sum_{j=1}^n \hat{b}_j l_j(x_j)$$

Learn the model with regularized regression.

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

$GA^2M$

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary



# Predictive learning via rule ensembles

## Tree-based ensembles

### Recent Works on Feature Interaction

Xingzhi Sun

### Outline

### Definitions of feature interaction

Post-hoc statistic

Model-specific

### Models

**Tree-based ensembles**

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

### Summary

## RuleFit (Friedman and Popescu, 2008)

$$F(\mathbf{x}) = \hat{a}_0 + \sum_{k=1}^K \hat{a}_k r_k(\mathbf{x}) + \sum_{j=1}^n \hat{b}_j l_j(x_j)$$

# Predictive learning via rule ensembles

## Tree-based ensembles

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

### RuleFit (Friedman and Popescu, 2008)

$$F(\mathbf{x}) = \hat{a}_0 + \sum_{k=1}^K \hat{a}_k r_k(\mathbf{x}) + \sum_{j=1}^n \hat{b}_j l_j(x_j)$$

### statistic for interaction

$$H_{jk}^2 = \sum_{i=1}^N \left[ \hat{F}_{jk}(x_{ij}, x_{ik}) - \hat{F}_j(x_{ij}) - \hat{F}_k(x_{ik}) \right]^2 / \sum_{i=1}^N \hat{F}_{jk}^2(x_{ij}, x_{ik})$$

Detect interaction with significant  $H_{jk}$ , whose distribution is obtained by bootstrapping.

# Detecting statistical interactions with additive groves of trees

## Tree-based ensembels

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic  
Model-specific

Models

Tree-based ensembels

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

the **accurate** model that contains the interaction

vs

the **accurate** model that bans the interaction

statistic for interaction

$$I_{ij}(F(\mathbf{x})) = \text{stRMSE}(F(\mathbf{x})) - \text{stRMSE}(R_{ij}(\mathbf{x}))$$

# Detecting statistical interactions with additive groves of trees

## Tree-based ensembels

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic  
Model-specific

Models

Tree-based ensembels

GA<sup>2</sup>M

Neuron Networks

Factorization  
Machines

Hybrid models with  
neuron networks

Summary

the **accurate** model that contains the interaction

vs

the **accurate** model that bans the interaction

statistic for interaction

$$I_{ij}(F(\mathbf{x})) = \text{stRMSE}(F(\mathbf{x})) - \text{stRMSE}(R_{ij}(\mathbf{x}))$$

- Obtain a highly accurate model through a tree ensemble that is later bagged:  $F_0(x) = \sum_{i=1}^K T_i(x)$
- Restrict interaction by forbidding one of the interacting variables when growing a tree.

# Detecting statistical interactions with additive groves of trees

## Tree-based ensembels

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic  
Model-specific

Models

Tree-based ensembels

$GA^2M$

Neuron Networks

Factorization  
Machines

Hybrid models with  
neuron networks

Summary

the **accurate** model that contains the interaction

vs

the **accurate** model that bans the interaction

statistic for interaction

$$I_{ij}(F(\mathbf{x})) = \text{stRMSE}(F(\mathbf{x})) - \text{stRMSE}(R_{ij}(\mathbf{x}))$$

- Obtain a highly accurate model through a tree ensemble that is later bagged:  $F_0(x) = \sum_{i=1}^K T_i(x)$
- Restrict interaction by forbidding one of the interacting variables when growing a tree.

Beats the Rulefit statistic in avoiding spurious interaction at sparse regions.

# Table of Contents

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

## 1 Definitions of feature interaction

- Post-hoc statistic
- Model-specific

## 2 Models

- Tree-based ensembles
- GA<sup>2</sup>M
- Neuron Networks
- Factorization Machines
- Hybrid models with neuron networks

# Accurate Intelligible Models with Pairwise Interactions

## GA<sup>2</sup>M

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

**GA<sup>2</sup>M**

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

- Main effects: GAM.
- Interaction: 2D bin functions  $f_{ij}$  on the residual of GAM.

# Accurate Intelligent Models with Pairwise Interactions

## GA<sup>2</sup>M

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization  
Machines

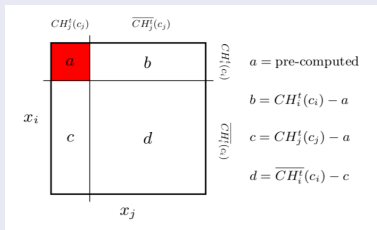
Hybrid models with  
neuron networks

Summary

- Main effects: GAM.
- Interaction: 2D bin functions  $f_{ij}$  on the residual of GAM.

## FAST

Speed up the calculation of bin averages: pre-caluculate a CDF lookup table for reusing.





# Table of Contents

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

## 1 Definitions of feature interaction

- Post-hoc statistic
- Model-specific

## 2 Models

- Tree-based ensembles
- GA<sup>2</sup>M
- Neuron Networks
- Factorization Machines
- Hybrid models with neuron networks

# Learning Global Pairwise Interactions with Bayesian Neural Networks

## Neuron Networks

### Recent Works on Feature Interaction

Xingzhi Sun

### Outline

### Definitions of feature interaction

Post-hoc statistic

Model-specific

### Models

Tree-based ensembles

GA<sup>2</sup>M

**Neuron Networks**

Factorization

Machines

Hybrid models with  
neuron networks

### Summary

## Definition of feature interaction for smooth models

The prediction has non-zero hessian over the interaction variables.

$$\frac{\partial F(\mathbf{x})}{\partial x_i \partial x_j} \neq 0$$

# Learning Global Pairwise Interactions with Bayesian Neural Networks

Neuron Networks

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization  
Machines

Hybrid models with  
neuron networks

Summary

## Definition of feature interaction for smooth models

The prediction has non-zero hessian over the interaction variables.

$$\frac{\partial F(\mathbf{x})}{\partial x_i \partial x_j} \neq 0$$

In practice, we take the expectation of the hessian over the distribution of  $\mathbf{x}$

# Learning Global Pairwise Interactions with Bayesian Neural Networks

## Neuron Networks

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

### Two ways of taking expectation

$$\text{EAH}_g^{i,j}(\mathbf{W}) = \mathbb{E}_{p(\mathbf{x})} \left[ \left| \frac{\partial^2 g^{\mathbf{W}}(\mathbf{x})}{\partial x_i \partial x_j} \right| \right] \quad \text{lowest FNR, highest FPR.}$$

$$\text{AEH}_g^{i,j}(\mathbf{W}) = \left| \mathbb{E}_{p(\mathbf{x})} \left[ \frac{\partial^2 g^{\mathbf{W}}(\mathbf{x})}{\partial x_i \partial x_j} \right] \right| \quad \text{lowest FPR, highest FNR.}$$

- EAH avoids  $(+, -)$  noise to cancel and could capture spurious interactions.
- AEH could have true interactions cancel out and fail to capture true interactions.

# Learning Global Pairwise Interactions with Bayesian Neural Networks

Neuron Networks

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

strike a balance: Group Expected Hessian (GEH)

$$M\text{-GEH}_g^{i,j}(\mathbf{W}) = \sum_{m=1}^M \frac{|A_m|}{\sum_{k=1}^M |A_k|} \left| \mathbb{E}_{p(\mathbf{x}|\mathbf{x} \in A_m)} \left[ \frac{\partial^2 g^{\mathbf{W}}(\mathbf{x})}{\partial x_i \partial x_j} \right] \right|$$

Partition the datapoints into  $M$  clusters, and expect the interaction is similar within each cluster, where only the noise is canceled out.

# Learning Global Pairwise Interactions with Bayesian Neural Networks

## Neuron Networks

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

strike a balance: Group Expected Hessian (GEH)

$$M\text{-GEH}_g^{i,j}(\mathbf{W}) = \sum_{m=1}^M \frac{|A_m|}{\sum_{k=1}^M |A_k|} \left| \mathbb{E}_{p(\mathbf{x}|\mathbf{x} \in A_m)} \left[ \frac{\partial^2 g^{\mathbf{W}}(\mathbf{x})}{\partial x_i \partial x_j} \right] \right|$$

Partition the datapoints into  $M$  clusters, and expect the interaction is similar within each cluster, where only the noise is canceled out.

Bayesian NN allows for the distribution of the  $M$ -GEH statistic and thus mean, std, confidence intervals, etc.

# Detecting Statistical Interactions from Neural Network Weights

## Neuron Networks

### Recent Works on Feature Interaction

Xingzhi Sun

### Outline

### Definitions of feature interaction

Post-hoc statistic

Model-specific

### Models

Tree-based ensembles

$GA^2M$

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

### Summary

How does a neuron network capture an interaction?

- Features share units of the first hidden layer.
- The shared units are passed to the output through descendent edges.

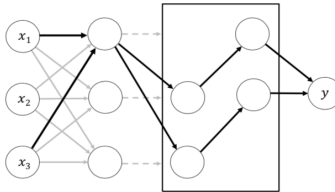


Figure: Interaction in an NN

# Detecting Statistical Interactions from Neural Network Weights

Neuron Networks

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

Write in the form of NN weights:

NID (Tsang et al., 2018)

$$\mathbf{z}^{(\ell)} = |\mathbf{w}^y|^\top \left| \mathbf{W}^{(L)} \right| \cdot \left| \mathbf{W}^{(L-1)} \right| \dots \left| \mathbf{W}^{(\ell+1)} \right|$$
$$\omega_i(\mathcal{I}) = z_i^{(1)} \mu \left( \left| \mathbf{W}_{i,\mathcal{I}}^{(1)} \right| \right)$$



# Detecting Statistical Interactions from Neural Network Weights

Neuron Networks

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

$GA^2M$

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

Write in the form of NN weights:

NID (Tsang et al., 2018)

$$\mathbf{z}^{(\ell)} = |\mathbf{w}^y|^\top \left| \mathbf{W}^{(L)} \right| \cdot \left| \mathbf{W}^{(L-1)} \right| \dots \left| \mathbf{W}^{(\ell+1)} \right|$$
$$\omega_i(\mathcal{I}) = z_i^{(1)} \mu \left( \left| \mathbf{W}_{i,\mathcal{I}}^{(1)} \right| \right)$$

The algorithm is fast because only the features with top weights is considered in each iteration.

# Neural Interaction Transparency: Disentangling Learned Interactions for Improved Interpretability

## Neuron Networks

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

- The first hidden layer captures first-order interactions.
- The following layers capture higher-order interactions.
- But they entangle and could contain spurious interaction:  
 $x_1x_2 + x_3x_4 \rightarrow x_1, x_2, x_3, x_4$  instead of  $x_1, x_2$  and  $x_3, x_4$

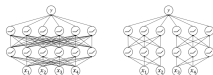


Figure: Entangled vs disentangled

# Neural Interaction Transparency: Disentangling Learned Interactions for Improved Interpretability

## Neuron Networks

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization  
Machines

Hybrid models with  
neuron networks

Summary

- The first hidden layer captures first-order interactions.
- The following layers capture higher-order interactions.
- But they entangle and could contain spurious interaction:  
 $x_1x_2 + x_3x_4 \rightarrow x_1, x_2, x_3, x_4$  instead of  $x_1, x_2$  and  $x_3, x_4$

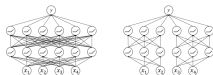
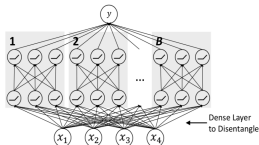


Figure: Entangled vs disentangled

Solution: Add a penalty on the weight matrix to control maximum allowed order of interaction.



# Feature Interaction Interpretability: A Case for Explaining Ad-Recommendation Systems via Neural Interaction Detection Neuron Networks

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

$GA^2M$

**Neuron Networks**

Factorization

Machines

Hybrid models with  
neuron networks

Summary

Given a black-box model, how to interpret global interaction?

- detect local interaction
- count the local interactions.
- interactions with many occurrences are interpreted as global.

# Feature Interaction Interpretability: A Case for Explaining Ad-Recommendation Systems via Neural Interaction Detection Neuron Networks

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

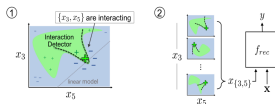
Factorization  
Machines

Hybrid models with  
neuron networks

Summary

Given a black-box model, how to interpret global interaction?

- detect local interaction
- count the local interactions.
- interactions with many occurrences are interpreted as global.



GLIDER (Tsang et al., 2020)

- 1 For each data instance  $x$ , perturb to get a local dataset, and predict on that dataset.
- 2 on the local dataset, detect interaction with NID.
- 3 count and rank occurrences of each interaction.

# Table of Contents

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization  
Machines

Hybrid models with  
neuron networks

Summary

## 1 Definitions of feature interaction

- Post-hoc statistic
- Model-specific

## 2 Models

- Tree-based ensembles
- GA<sup>2</sup>M
- Neuron Networks
- Factorization Machines
- Hybrid models with neuron networks

# Factorization machines

## Factorization Machines

Background: Highly-sparse categorical data with one-hot coding. e.g. Shopping history, movie reviews.

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

$GA^2M$

Neuron Networks

**Factorization  
Machines**

Hybrid models with  
neuron networks

Summary

# Factorization machines

## Factorization Machines

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization  
Machines

Hybrid models with  
neuron networks

Summary

Background: Highly-sparse categorical data with one-hot coding. e.g. Shopping history, movie reviews.

### FM (Rendle, 2010)

$$\hat{y}(\mathbf{x}) := w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j$$

learn the parameters with gradient descent

Each feature  $x_i$  corresponds to an *embedding vector*  $\mathbf{v}_i$ .

Interaction: the inner products of the vectors  $\langle \mathbf{v}_i, \mathbf{v}_j \rangle$ .



# Factorization machines

## Factorization Machines

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

$GA^2M$

Neuron Networks

Factorization  
Machines

Hybrid models with  
neuron networks

Summary

Background: Highly-sparse categorical data with one-hot coding. e.g. Shopping history, movie reviews.

### FM (Rendle, 2010)

$$\hat{y}(\mathbf{x}) := w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j$$

learn the parameters with gradient descent

Each feature  $x_i$  corresponds to an *embedding vector*  $\mathbf{v}_i$ .

Interaction: the inner products of the vectors  $\langle \mathbf{v}_i, \mathbf{v}_j \rangle$ .

Advantages over regression (polynomial kernel SVM):

- Generalizes to instances that do not appear in the training set.
- linear complexity.

# Variants of FM

## Factorization Machines

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization  
Machines

Hybrid models with  
neuron networks

Summary

### Field-aware Factorization Machines for CTR Prediction

- 1 Categorize features into fields: Clothes, Food, Electronics,...
- 2 For each feature, instead of one embedding vector, learn a vector for each field.
- 3 When calculating interaction, use the vectors matching each other's field to take inner product.

$$\phi_{\text{FFM}}(\mathbf{w}, \mathbf{x}) = \sum_{i=1}^n \sum_{j=i+1}^n (\mathbf{w}_{i,f_2} \cdot \mathbf{w}_j, f_1) x_i x_j$$

### Attentional Factorization Machines

Instead of taking the inner product of the embedding vectors, take weighed outer product with weights learned by an Attention Neuron Network.

# Table of Contents

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

## 1 Definitions of feature interaction

- Post-hoc statistic
- Model-specific

## 2 Models

- Tree-based ensembles
- GA<sup>2</sup>M
- Neuron Networks
- Factorization Machines
- Hybrid models with neuron networks

# Wide & Deep Learning for Recommender Systems

Hybrid models with neuron networks

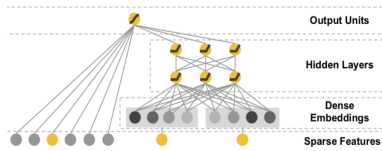


Figure: Wide & Deep

A **jointly-trained** hybrid model:

- Linear crossing model: *manual* low-order interactions.
- Deep model: *automatic* high-order interactions.

# Deep & Cross Network for Ad Click Predictions

## Hybrid models with neuron networks

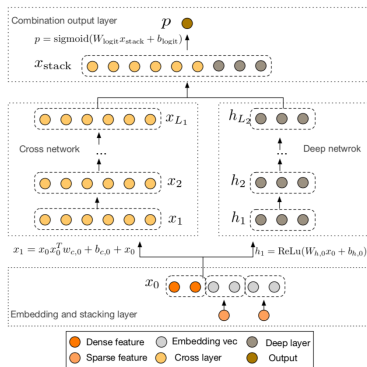


Figure: Deep & Cross

A **jointly-trained** hybrid model:

- Crossing network: *automatic* low-order interactions.
- Deep network: *automatic* high-order interactions.

# DeepFM: A Factorization-Machine based Neural Network for CTR Prediction

Hybrid models with neuron networks

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization  
Machines

Hybrid models with  
neuron networks

Summary

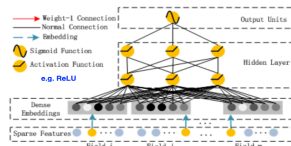
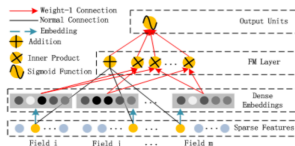


Figure: DeepFM

A **jointly-trained** hybrid model:

- FM: *automatic* low-order interactions.
- Deep model: *automatic* high-order interactions.

# xDeepFM: Combining Explicit and Implicit Feature Interactions for Recommender Systems

Hybrid models with neuron networks

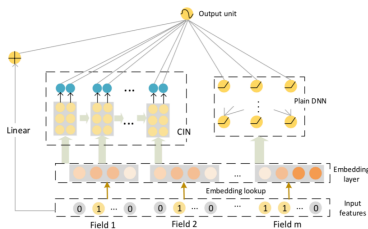


Figure: xDeepFM

A **jointly-trained** hybrid model:

- Linear model: main effects.
- Compressed Interaction Network: *automatic* low-order interactions: Interaction of each order goes to the output.
- Deep neuron network: *automatic* high-order interactions.

# Deep Interest Network for Click-Through Rate Prediction

## Attention Neuron Networks

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

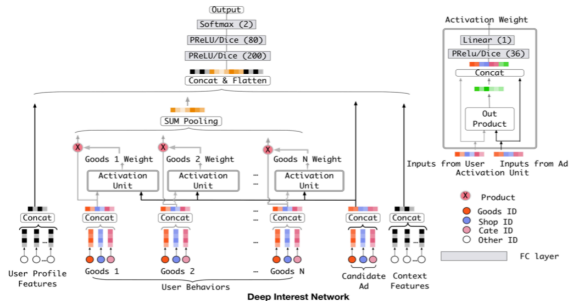


Figure: Deep Interest Network

The network comprises Activation Units, an attention network that leverages the user preference history and identify the true interest from diverse interests.



# AutoInt: Automatic feature interaction learning via self-attentive neural networks

Attention Neuron Networks

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization

Machines

Hybrid models with  
neuron networks

Summary

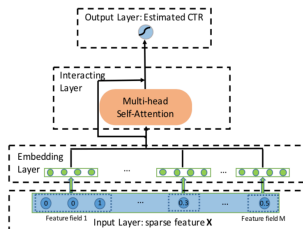


Figure: AutoInt

Each *multi-head self-attentive* layer captures interaction by learning interaction weights.  
The model is a black-box with all the interactions entangled.

# Summary

## Models for capturing feature interactions

Recent Works  
on Feature  
Interaction

Xingzhi Sun

Outline

Definitions of  
feature  
interaction

Post-hoc statistic

Model-specific

Models

Tree-based ensembles

GA<sup>2</sup>M

Neuron Networks

Factorization  
Machines

Hybrid models with  
neuron networks

Summary

### Models for capturing feature interactions:

- post-hoc
  - Tree Ensembles: non-additiveness, expressiveness
  - Neuron Networks: hessian, weights
- ad-hoc
  - Regularized NN: additive terms
  - GA2M: simple, fast
  - FM: sparse data, generalization, fast
  - Non-Deep and Deep hybrids: low-order and high-order