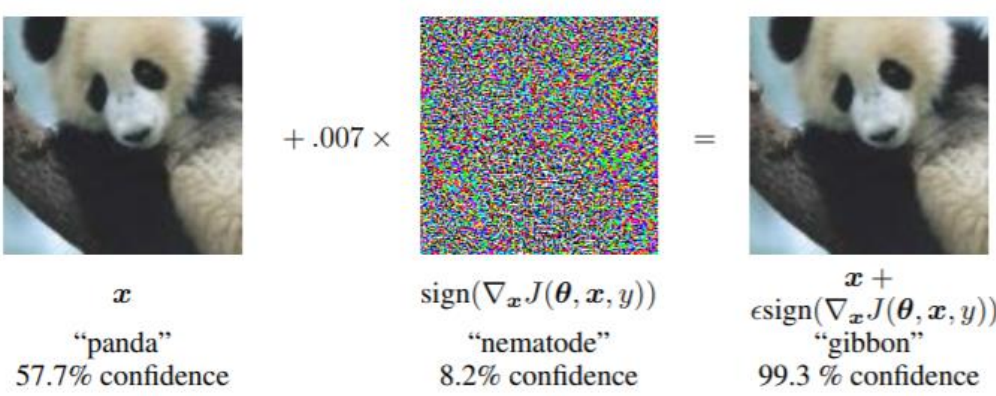# Adversarial attack to Semantic Parser

Weiliang Tang, Shilin He (TA) , Michael Lyu (Prof.)

NLP

## Introduction of New Adversarial Task for Semantic Parser

**Adversarial attack to image classification model**

$x$
"panda"
57.7% confidence

$+ .007 \times$

$sign(\nabla_x J(\theta, x, y))$
"nematode"
8.2% confidence

$=$

$x + \epsilon sign(\nabla_x J(\theta, x, y))$
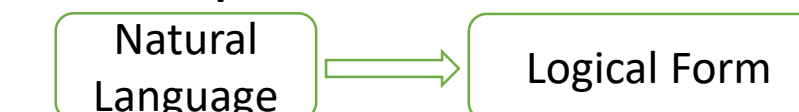"gibbon"
99.3 % confidence

**Adversarial attack to text classification model**

South Africa's historic Soweto township marks its 100th birthday on Tuesday in a mood of optimism.
57% **World**

South Africa's historic Soweto township marks its 100th birthday on Tuesday in a moo**P** of optimism.
95% **Sci/Tech**

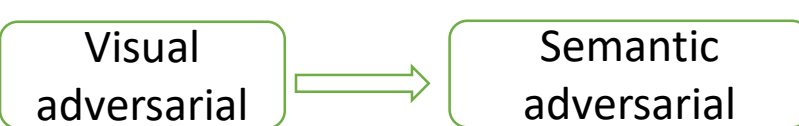**New Challenge to Attack Semantic Parser:**
- The input is short, change of input is clearly distinguishable told visually
- The input space is discrete,

**Semantic parser:**

Natural Language → Logical Form

**A new definition of adversarial example: x\*** where
- Semantic(x) = Semantic(x\*)
- Semantic(Model(x)) ≠Semantic(Model(x\*))

Visual adversarial → Semantic adversarial

Long input
Continuous input

Short discrete input
Semantic task

what is the name of the **loser** when the winner was new england patriots , …?

NL2SQL Model

SELECT **loser** WHERE winner cond_op new england patriots …

what is the name of the **losers** when the winner was new england patriots , …?

Same NL2SQL Model

SELECT **winner** WHERE winner cond_op new england patriots …
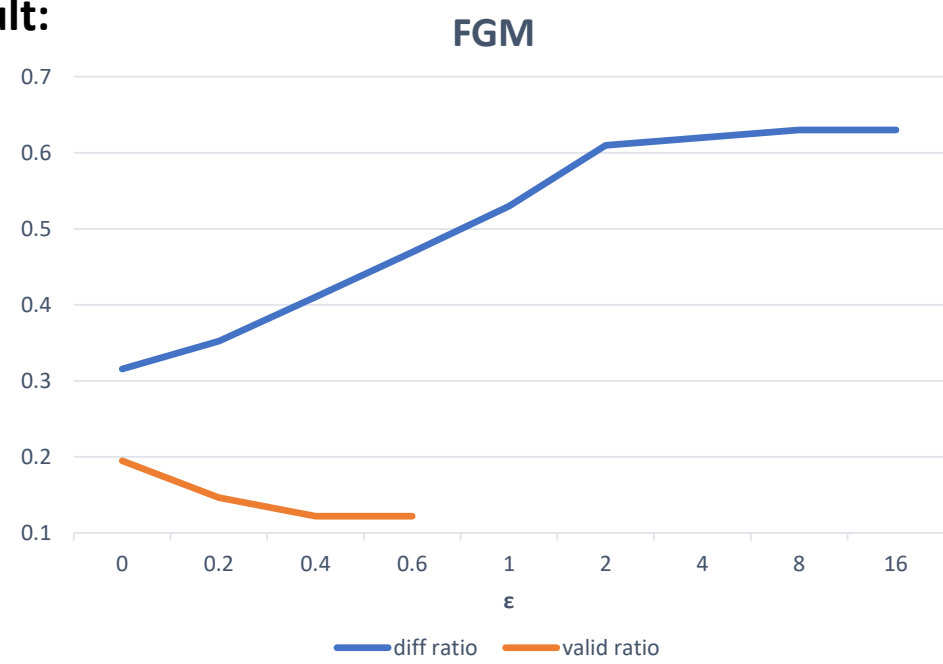
## Generating Adversarial Examples

**Measurement:**
**Correct ratio**: correct predictions/ input data
**Differ ratio**: diff predictions/ perturbed input data.
**Valid ratio**: predictions which keep the semantic meaning unchanged / different outputs.

**Basic Method: Fast Gradient Method**
**Algorithm:**

```
FGM
1   // grad_data = (input_len × embedding_size)
2   for i = 0 to length[grad_data] − 1
3       word_grad[i] = ‖grad_data[i]‖
4       target_word = arg max(word_grad)
5   perturbed_word =   arg min    ‖word[idx] + ε · grad_data[idx] − w‖
                      w∈embed_space
```

**Experiment Result:**



FGM

- The larger the ε is, the higher diff ratio and lower valid ratio it will be when ε is relatively small.
- Some pattern is shown in the successful perturbed examples:
1. Among all the successful example, 36% is done by changing a word in single form to plural form.

What is the air **force** cross when … => SELECT **airforcecross** WHERE…
What is the air **forces** cross when … => SELECT **navyforcecross** WHERE…

what **gender** is quentin ? => SELECT **gender** WHERE name = quentin
what **genders** is quentin ? => SELECT **status** WHERE name = quentin
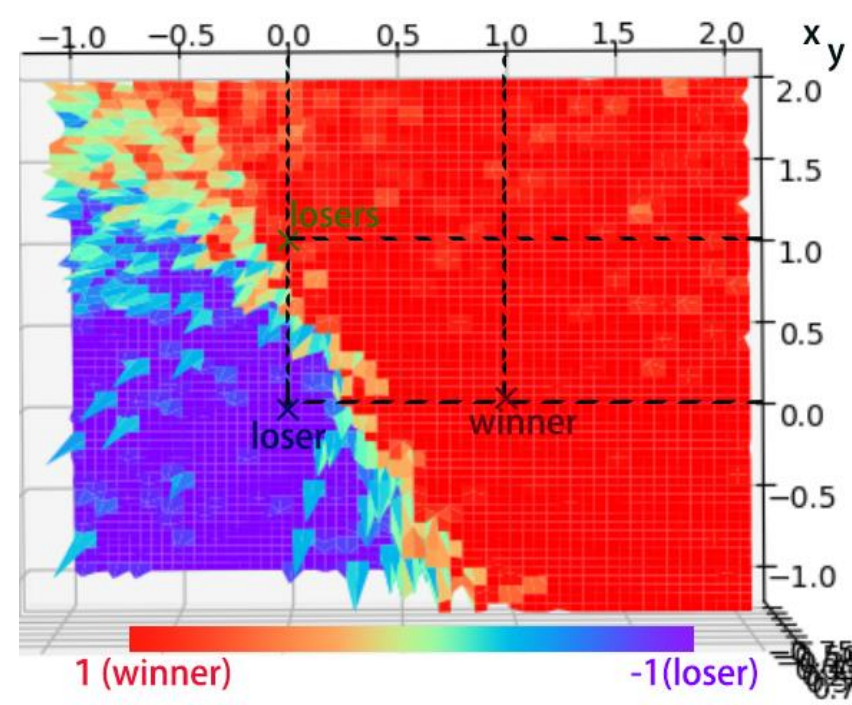
2. Substitute the word with its synonym

How many types of organization … => SELECT MAX **types** WHERE…
How many kinds of organization … => SELECT MAX **organization** WHERE…

- **Drawback:**
The choice of word neglects the semantic environment around it, one word can be perturbed only into another fixed word under on circumstances

### Fast Gradient Method

**Reason: Under-fitting problem in NL2SQL task**

The distribution of z (see below) on a plane in the high dimensional space
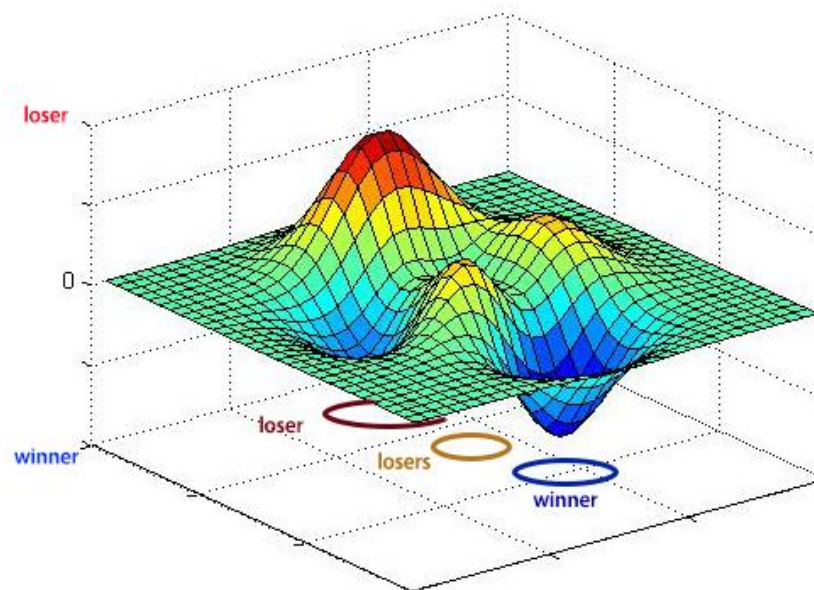


1 (winner)      -1(loser)

$$z = \frac{(p(sel\_op='loser') - p(sel\_op='winner'))}{(p(sel\_op='loser') + p(sel\_op='winner'))}$$

$$p(sel\_op) = Model(\vec{v}('losers') + x \cdot (\vec{v}('losers') - \vec{v}('loser')) + y \cdot (\vec{v}('winner') - \vec{v}('loser')); \theta)$$
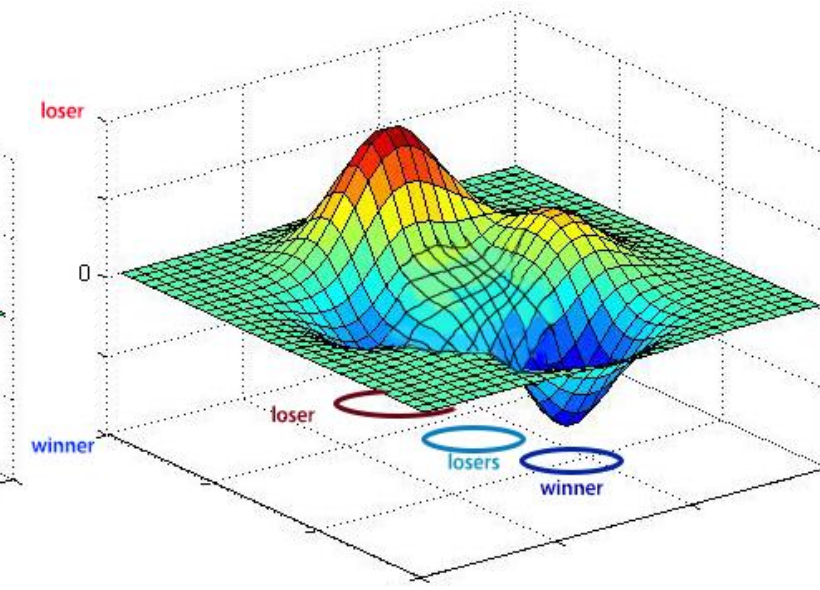
The plane where "loser", "losers" and "winner" lie in

when $x = 1, y = 0, p(sel\_op) = Model(\vec{v}('losers'))$
when $x = 0, y = 1, p(sel\_op) = Model(\vec{v}('winner'))$
when $x = y = 0, p(sel\_op) = Model(\vec{v}('loser'))$

**What's supposed to be**          **What it looks like actually**



**Under fitting problem:** some words are crowded in a small area, the word untrained is easily been misguided by the trained words around it
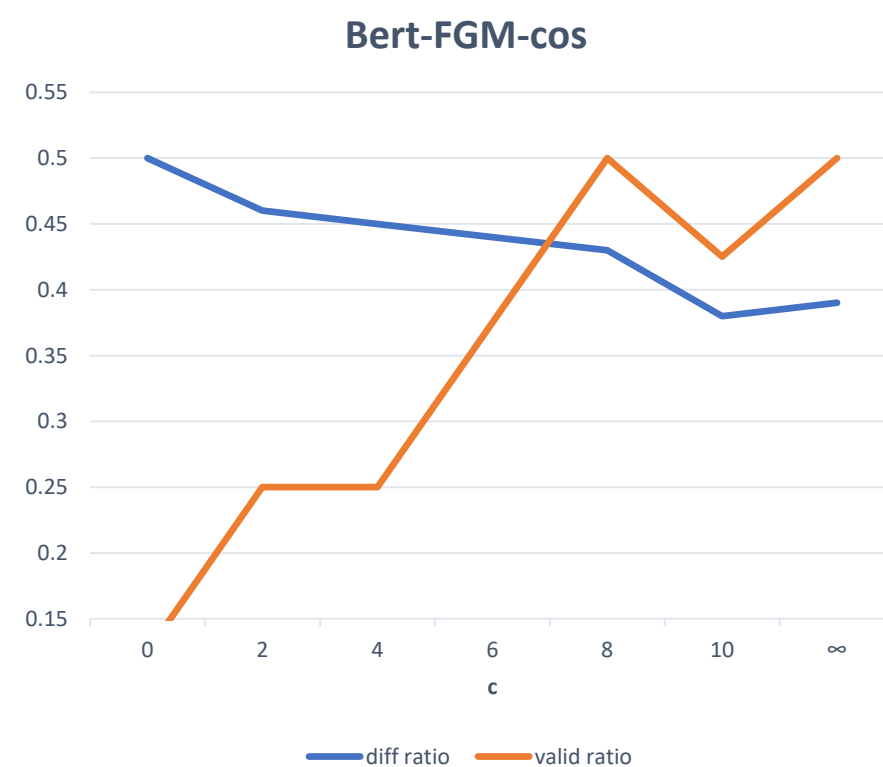
- **New adversarial feature for NL2SQL model:** The header of SQL usually are of the same type and sometimes very close to each other, the header can be vulnerable under adversarial attack

---

```
BERT-FGM
1   for i = 0 to 3
2       // grad_data = (input_len × embedding_size)
3       for i = 0 to length[grad_data] − 1
4           word_grad[i] = ‖grad_data[i]‖
5       target_word_list = n_arg max(word_grad)
6       for i = 0 to length[idx_list] − 1
7           target_word = target_word_list[i]
8           bert_list = Bert(sen, target_word, 10)
9           word_list =   arg max   c · bert_prob[w] + cos_simi(ε · grad_data[idx], w − target_word)
                        w∈bert_list
10      perturbed_word =   arg max   c · bert_prob[w] + cos_simi(ε · grad_data[idx], w − target_word)
                         w∈word_list
11      word ⇒ perturbed_word
```

**Algorithm:**

**Experiment result:**

- A trade off between diff ratio and valid ratio
  - The smaller the c is, the more dominant the cosine_similarity will be, the word is more likely to follow the gradient straightly, the higher diff ratio is.
  - The bigger the c is , the more dominant the bert_probability will be, the word is more likely to make sense, the higher valid ratio is, but it may not follow the gradient too much.
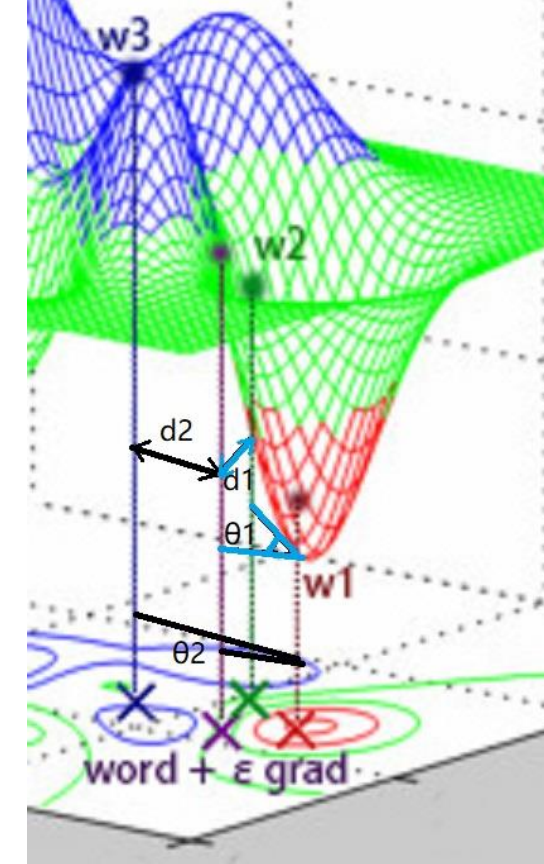- This method successfully elaborate the valid ratio compared to previous simple FGM method

- A more variety of forms of successful example is shown
- A more semantic consistency is shown after substitution

### Improvement Using Bert

**Cosine similarity is a more reasonable choice：**
- d1 < d2
- loss(d2) > loss(d1)
- cos_similarity describe the degree of following the gradient better since Bert ensures the small distance already



- Unreasonable result occurs if using norm distance



Bert-FGM-norm



Bert-FGM-cos

| Original sentence | Perturbed sentence |
|---|---|
| what is height , **when** rank is less than 20… | What is height, **where** rank is less than 20… |
| **When** total goals have a fa cup apps larger than 1 ,…,what is the total number? | **if** total goals have a fa cup apps larger than 1 ,…,what is the total number? |
| what is the smallest period -lrb- days -rrb- to have a planetary mass **of** 1, and … | what is the smallest period -lrb- days -rrb- to have a planetary mass **at** 1, and … |

**Weiliang Tang**
Email:1155107670@link.cuhk.edu.hk