Machine Learning and Privacy

Vitaly Shmatikov

Machine Intelligence LANDSCAPE

ARTIFICIAL INTELLIGENCE



DEEP LEARNING



MACHINE LEARNING



NLP **PLATFORMS**



PREDICTIVE APIS



IMAGE RECOGNITION

PERSONAL

ASSISTANT

clarifai MADBITS DNNresearch | DEXTRO VISENZE I lookflow

SPEECH RECOGNITION

CORE TECHNOLOGIES

popup archive NUANCE

RETHINKING ENTERPRISE

SALES



SECURITY / AUTHENTICATION



FRAUD DETECTION



HR / RECRUITING



MARKETING

INTELLIGENCE TOOLS

MADATAD Q Palantir # FirstRain

> transcriptic ** ZEPHYR

bing TUTE

ADTECH





AYASDI kaggle TACHYUS biota

. Flutura

Tech 2015: Deep Learning And Machine Intelligence Will Eat The World

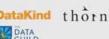
FINANCE





PHILANTHROPIES





Diligence Engine

AUTOMOTIVE



DIAGNOSTICS





RETAIL

MEDICAL

▼Parzival

@grand round table



RETHINKING HUMANS / HCI

AUGMENTED REALITY





COMPUTING

Gesture lak

△ Prismatic al



ROBOTICS



EMOTIONAL RECOGNITION

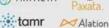


SUPPORTING TECHNOLOGIES

HARDWARE











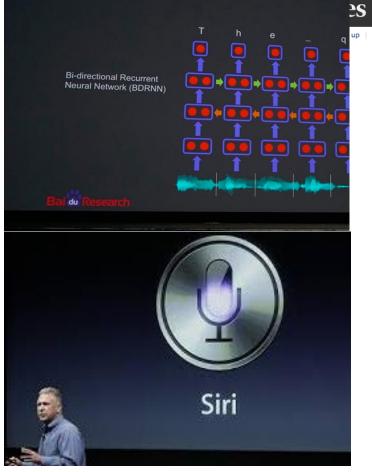
Baidu Deep Spen-

Sign up



Robert Hof Contributor TECH 12/18/2014 @ 9:00AM | 48,377 views

Baidu Announces Breakthrough In Speech Recognition, Claiming To Top Google And Apple







New Posts +1

Michael Thomsen Contributor

TECH 2/19/2015 @ 1:06PM 4,996 views

Microsoft's Deep Learning Project Outperforms Humans In Image Recognition



GT: horse cart 1: horse cart 2: minibus 3: oxcart 4: stretcher 5: half track



GT: coucal
1: coucal
2: indigo bunting
3: lorikeet
4: walking stick
5: custard apple



GT: birdhouse
1: birdhouse
2: sliding door
3: window screen
4: mailbox
5: pot



GT: komondor 1: komondor 2: patio 3: llama 4: mobile home



2: garbage truck

4: trailer truck

5: go-kart

GT: yellow lady's slipper 1: yellow lady's slipper 2: slug 3: hen-of-the-woods 4: stinkhorn 5: coral fungus



1: drumstick 2: candle 3: wooden spoon 4: spatula



1: Band Aid 2: ruler 3: rubber eraser 4: pencil box

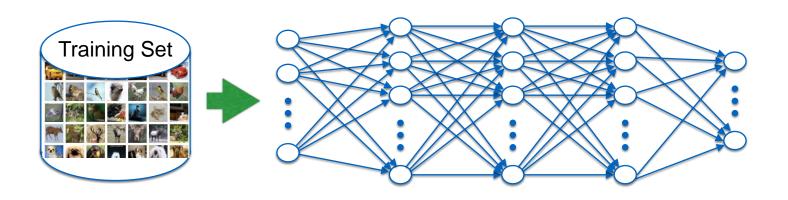


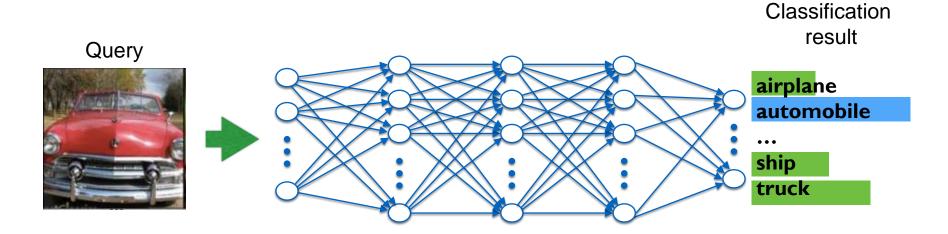
GT: spotlight 1: grand piano 2: folding chair



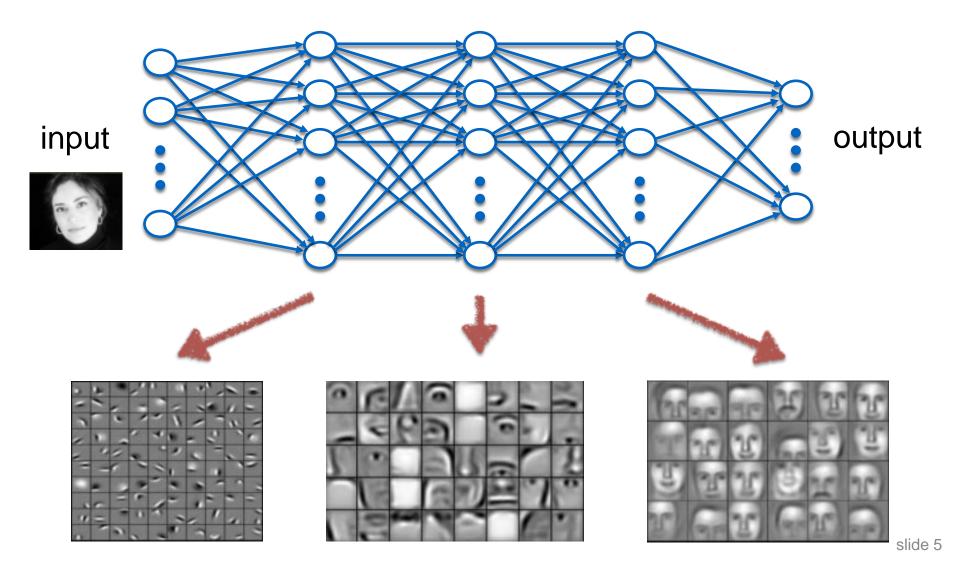
GT: spotlight 1: acoustic guitar 2: stage

Typical Task: Classification

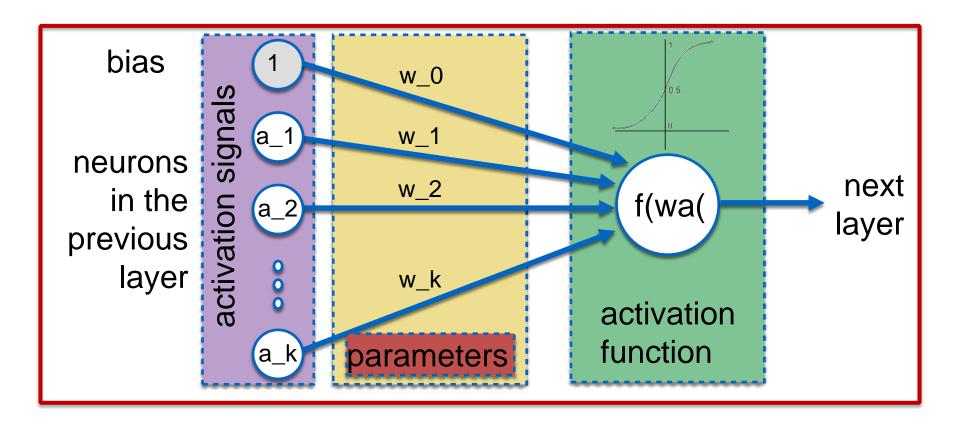




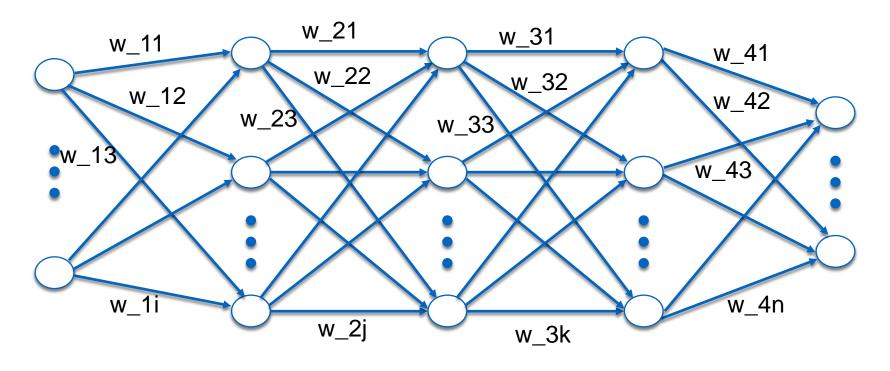
Deep Neural Networks



Deep Neural Networks

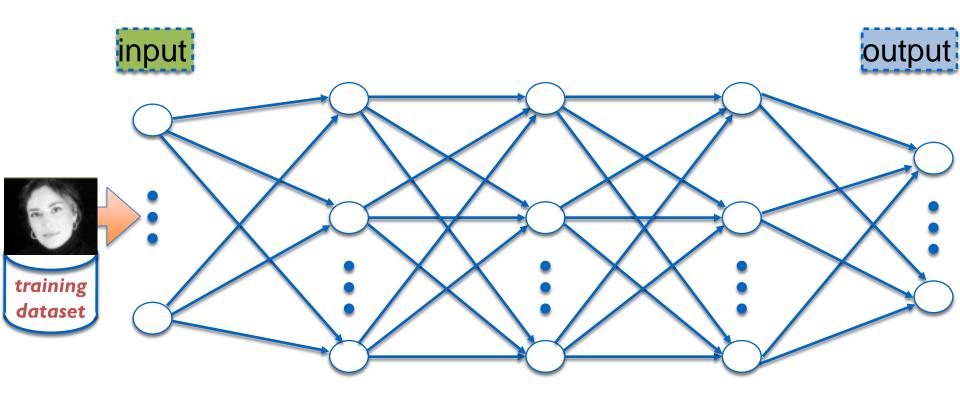


Learn parameters using Stochastic Gradient Descent (SGD)

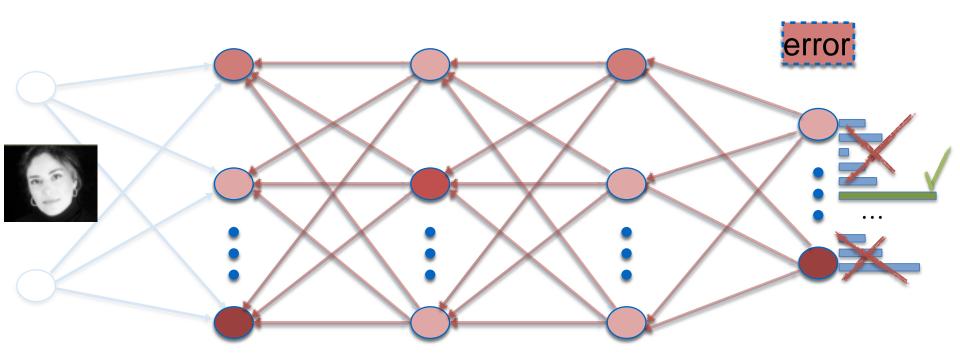


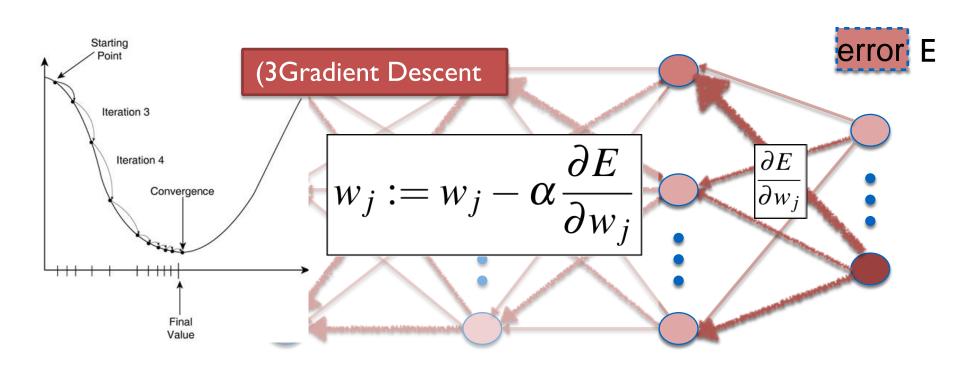
Find parameters that minimize the classification error

(IFeed-Forward

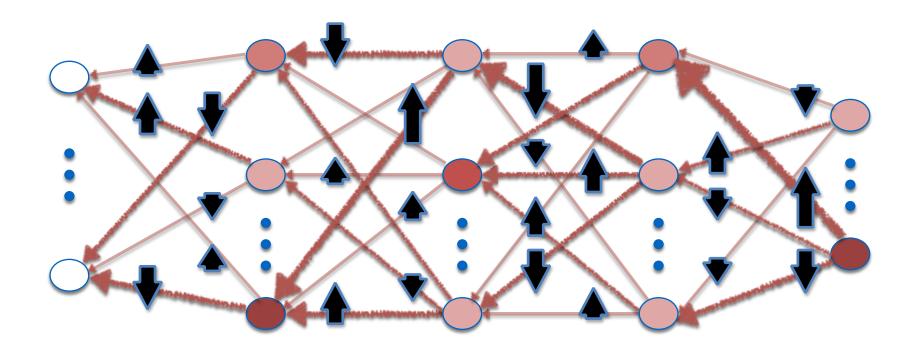


(2Back-propagation





Parameter Update



Repeat for new batches of training data

2014

Users' data Services **Threats** -Collection of sensitive personal data -Anonymization and re-identification -Inference attacks -Side channels

2018

Users' data

Machine learning

Services

Output

Description

Services

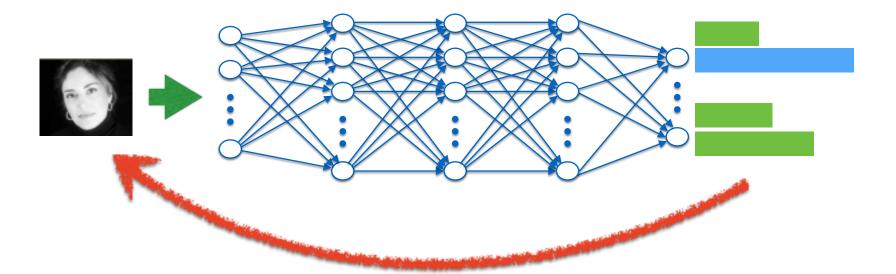
Do trained models leak sensitive data?

Is it possible to train a "good" model while respecting privacy of training data?

Is it possible to keep the model itself private?

Model Inversion

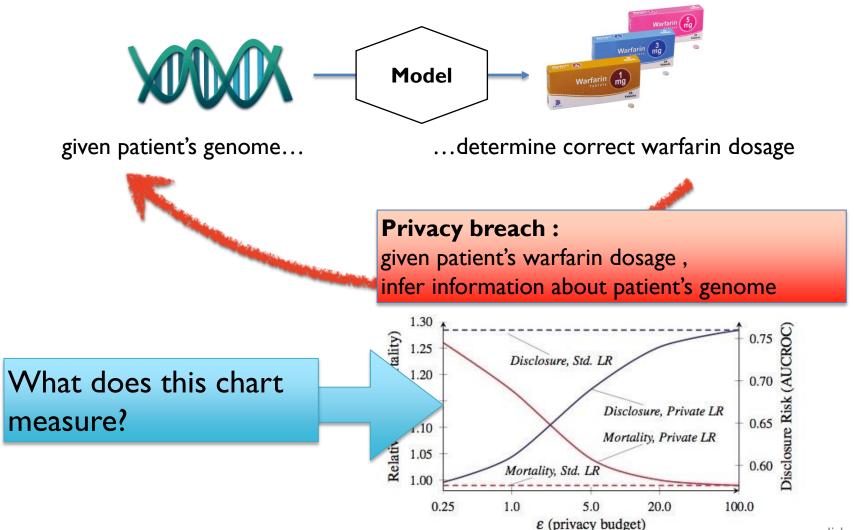
Fredrikson et al.



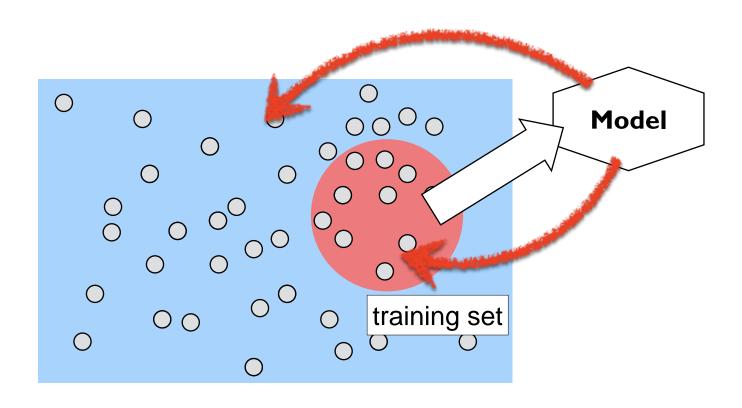
Given an output of a machine learning model, infer something about the input

"unexpected attributes"

Model Inversion in Action



Does Inference Breach Privacy?



Recommended Reading

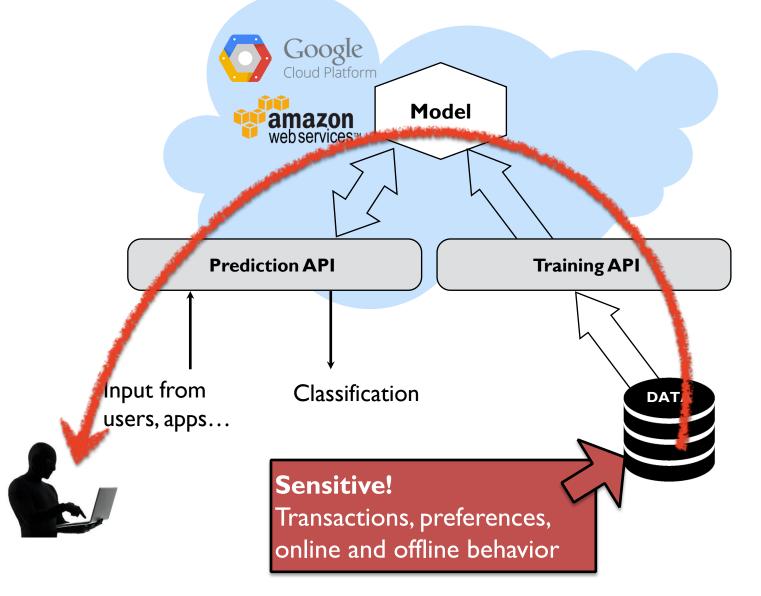
Frank McSherry.

"Statistical inference considered harmful"

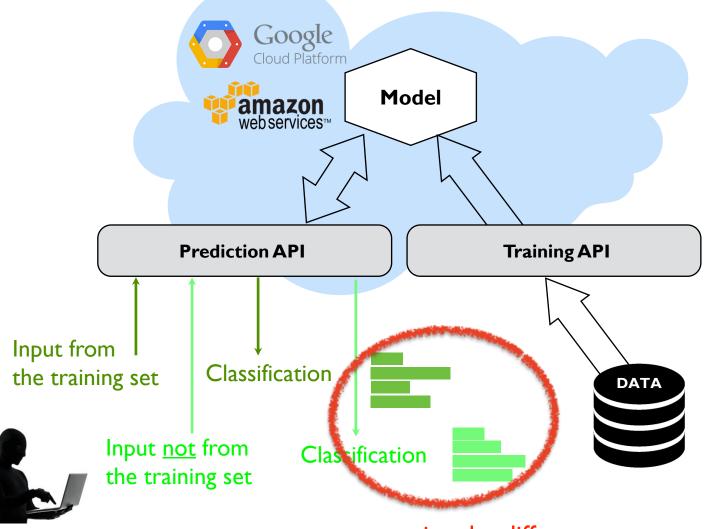


https://github.com/frankmcsherry/blog/blob/master/posts/2016-06-14.md

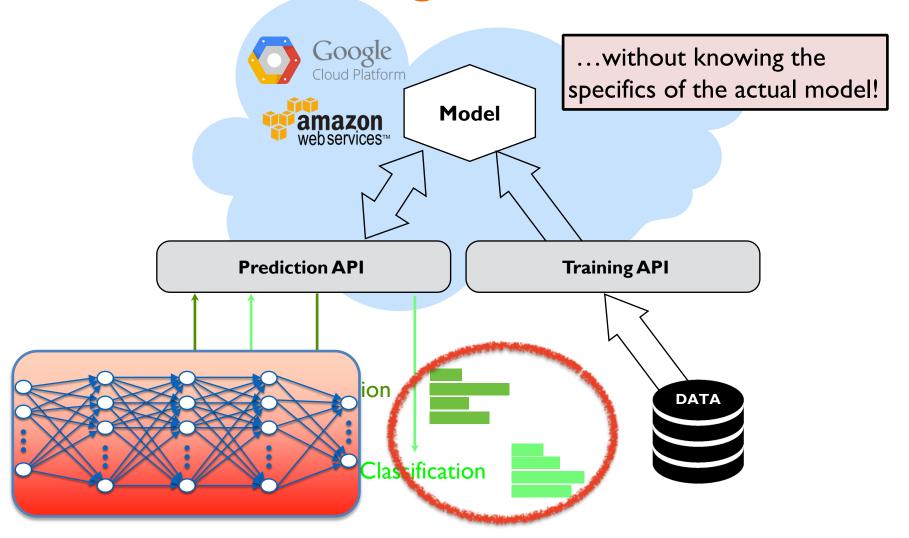
Machine Learning as a Service



Exploiting Trained Models



ML Against ML

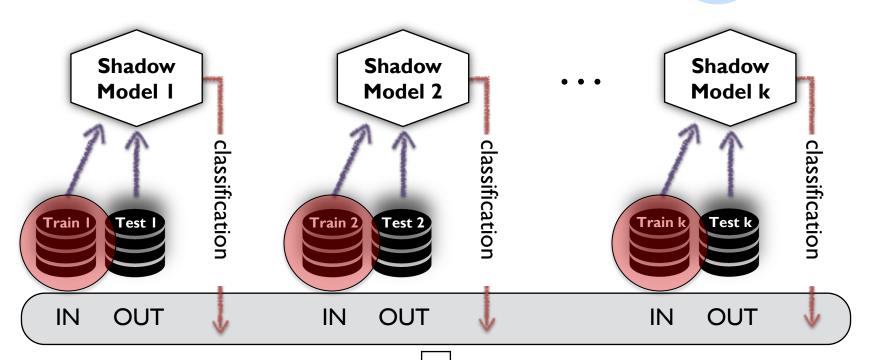


Train a model to...

recognize the difference

Training Attack Model using Shadow Models





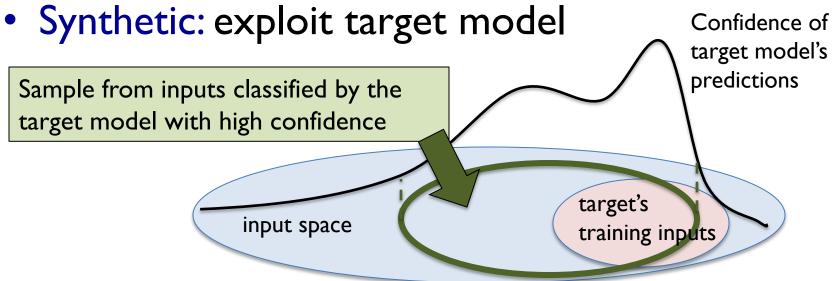


Train the attack model

to predict if an input was a member of the training set (in (or a non-member (out(

Training Data for Shadow Models

- Real: must be similar to training data of the target model (drawn from same distribution)
- Synthetic: sample feature values from (known) marginal distributions



Synthesizing Shadow Training Data

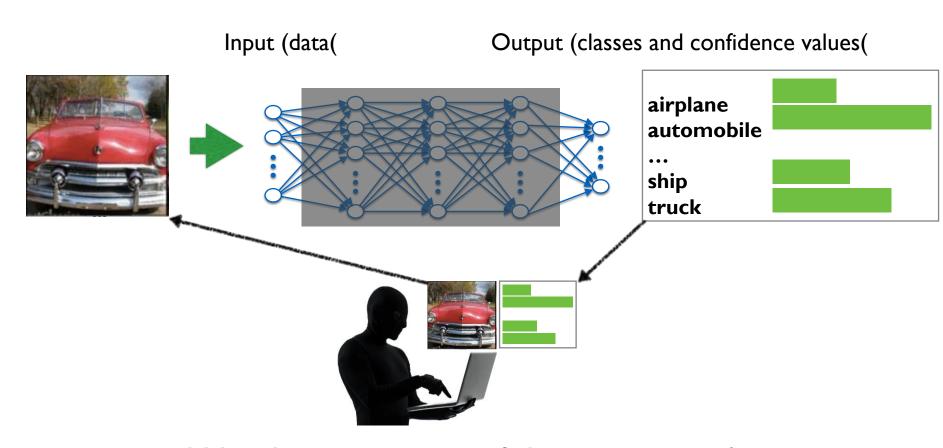
Hill-climb the space of possible inputs to find those classified by the target model with high confidence

Sample from these inputs to synthesize the training dataset for shadow models

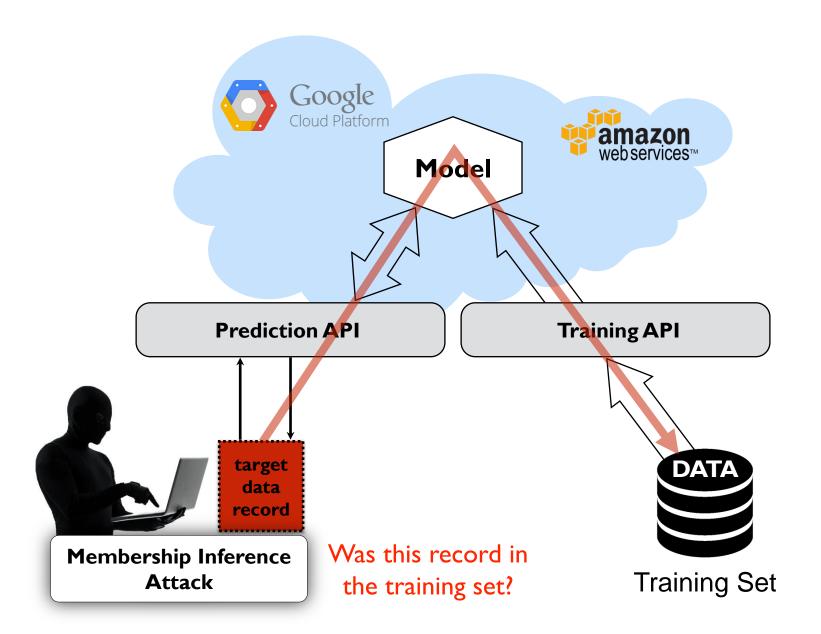
If many candidate inputs rejected by the target model, re-randomize some features and try again

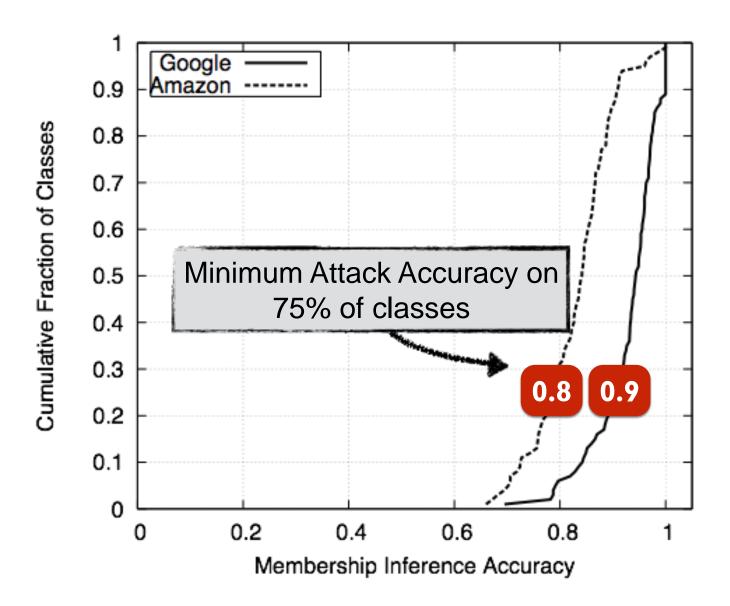
```
Algorithm 1 Data synthesis using the target model
 1: procedure SYNTHESIZE(class : c)
          \mathbf{x} \leftarrow \text{RANDRECORD}()
                                                 ▶ initilize a record randomly
          y_c^* \leftarrow 0
          j \leftarrow 0
           k \leftarrow k_{max}
          for iteration = 1 \cdots iter_{max} do
                \mathbf{y} \leftarrow f_{\mathsf{target}}(\mathbf{x})
                                                      ▶ query the target model
                if y_c \geq y_c^* then
                                                             ▶ accept the record
                     if y_c > \text{conf}_{min} and c = \arg \max(\mathbf{y}) then
                          if rand() < y_c then
                                                                          ▶ sample
                               return x
                                                                 > synthetic data
                          end if
                     end if
                     \mathbf{x}^* \leftarrow \mathbf{x}
                     y_c^* \leftarrow y_c
                     j \leftarrow 0
                else
17:
                     j \leftarrow j + 1
 18:
                     if j > rej_{max} then \triangleright many consecutive rejects
 19:
                          k \leftarrow \max(k_{min}, \lceil k/2 \rceil)
                          j \leftarrow 0
                     end if
22:
                end if
23:
                \mathbf{x} \leftarrow \text{RANDRECORD}(\mathbf{x}^*, k) \triangleright randomize \ k \ features
24:
           end for
25:
           return |
                                                           ▶ failed to synthesize
26:
27: end procedure
```

Membership Inference Attack



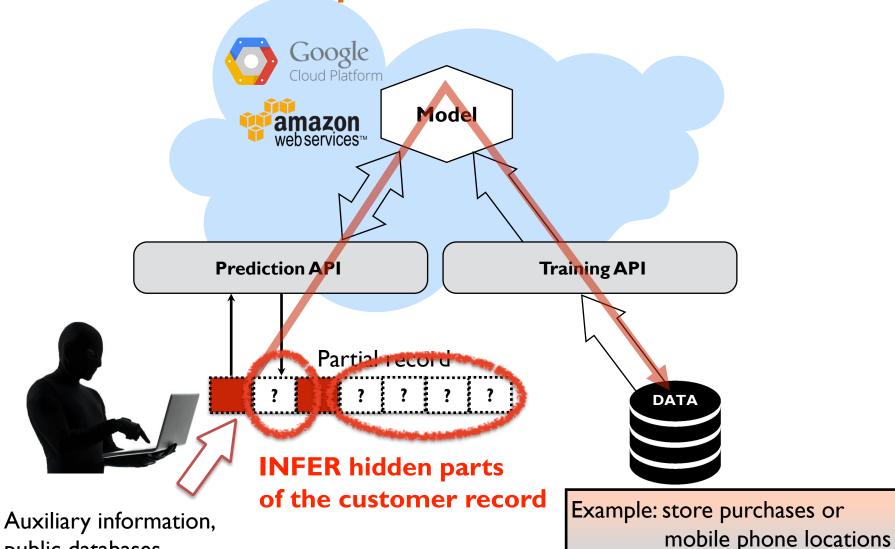
Was this image part of the training set?





Purchase Dataset — Classify Customers

Next Step: Reconstruction

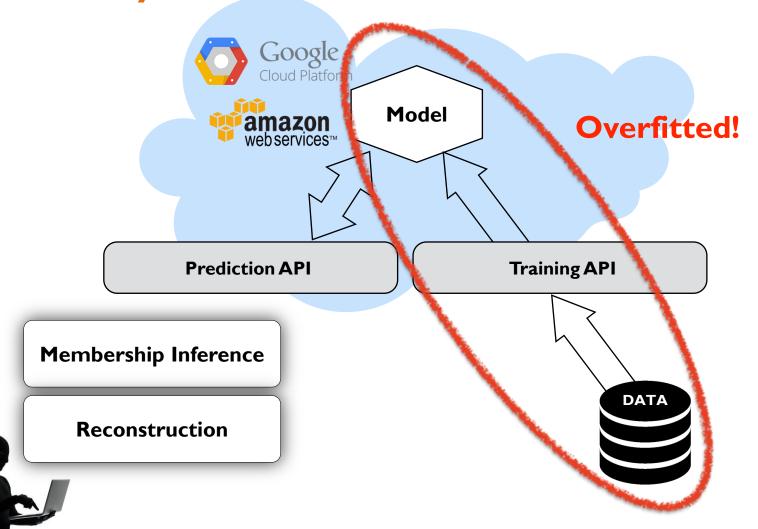


public databases,

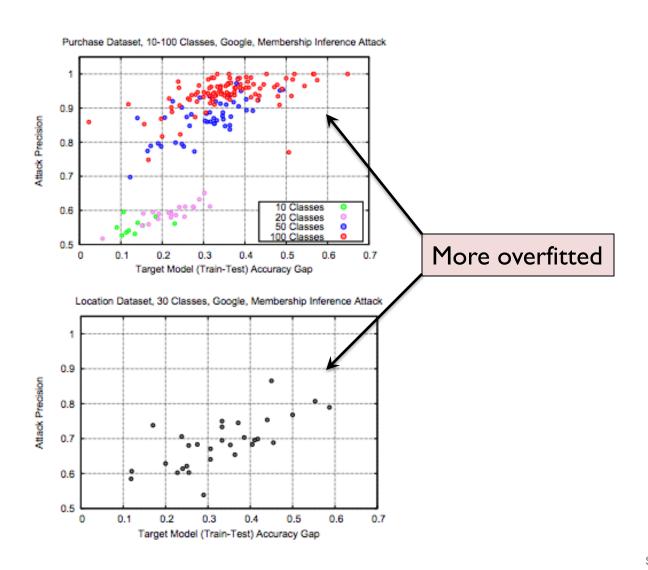
accidentally revealed data

slide 27

Why Do These Attacks Work?



Attack Success vs. Test-Train Gap

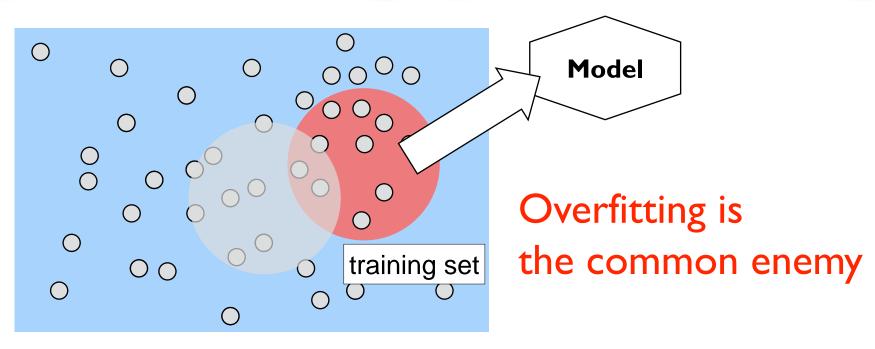


Privacy:

Does the model leak information about data in the training set?

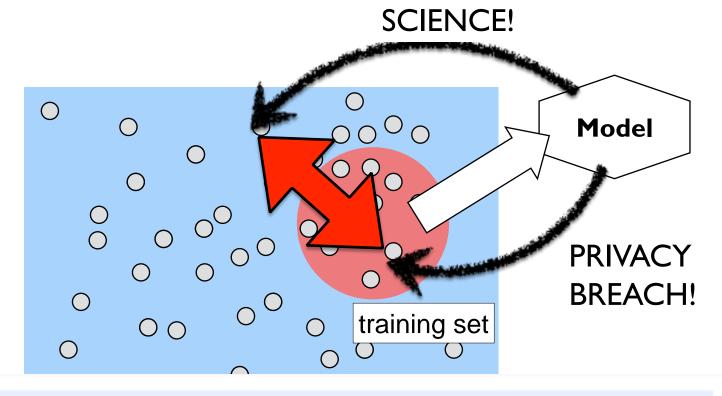
Learning:

Does the model generalize to data outside the training set?



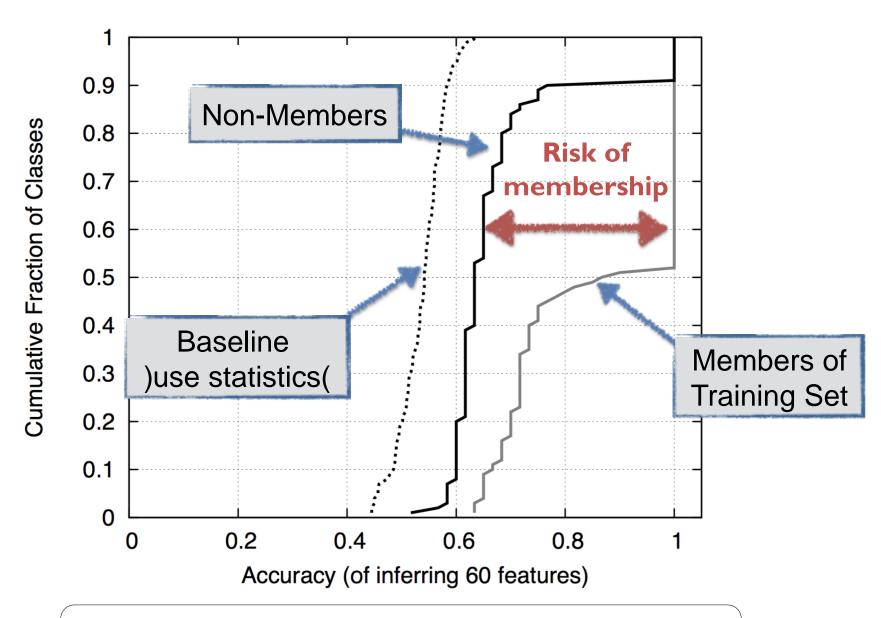
data universe

Does Inference Breach "Privacy?"



Privacy breach = risk of membership:

Gap between what can be inferred from the model about a member of the training set and an arbitrary input from the population



Future

- Modern machine learning is both a threat and an opportunity for data privacy
- For once, privacy and utility are not in conflict: overfitting is the common enemy
 - Overfitted models leak training data
 - Overfitted models lack predictive power
- Need generalizability and accuracy

