# Privacy in the Year 2027

Vitaly Shmatikov



www.shivonzilis.com/machineintelligence





# 2014

### Users' data

#### Services



#### **Threats**

- -Collection of sensitive personal data
- -Anonymization and re-identification
- -Inference attacks

5

-Side channels

### 2020



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?

# Typical Task: Classification





# Perceptron (1957, Cornell)

- 400 pixel camera
- Designed for image recognition
- "Knowledge" encoded as weights in potentiometers (variable resistors)
- Weights updated during learning performed by electric motors





## What Perceptron Can Learn



# Feedforward Neural Networks



### **Activation Functions**



# Universal Approximation Theorem

 Multilayer perceptron with a single hidden layer and linear output layer can approximate any continuous function on a compact subset of R<sup>n</sup> to within any desired degree of accuracy

- Under some assumptions about activation functions

# Why Deep Neural Networks

- Functions representable with a deep network can require exponential number of hidden units with a shallow (single hidden layer) network
- Piecewise linear networks (e.g., using ReLU) can represent functions that have a number of regions exponential in depth of network
  - Capture repeated, mirroring, symmetric patterns in the data
  - Often better generalization

# **Convolutional Neural Networks**



## Example: Face Recognition



## **Parameter Training**



Find parameters that minimize the classification error

# Loss Function (Cost Function, Objective Function)



# Loss Function

- Measures the "cost" of fitting a model to data
- Examples:
  - L<sup>2</sup> squared difference between model output and known correct output ("ground truth")
  - L<sup>I</sup> absolute difference between model output and known correct output
  - Cross-entropy between model output (interpreted as probability) and correct output

# Measuring Accuracy

- Training dataset: model is "trained" to fit this
- Validation dataset: model is repeatedly tested on this data during training
- Testing dataset: measure accuracy of a trained model

Testing accuracy vs. training accuracy



I) Feed-forward



# Parameter Training

2) Back-propagation



### "Batch" Gradient Descent



# Stochastic Gradient Descent

- Need to compute sum of n terms, n is large
- Sample instead of computing the full sum
  - Example: train on 100 photos to compute an estimate, then repeat and update the estimates

Repeat until an approximate minimum is reached:

- Pick a random sample of training examples
- For i in 1, 2 ... n

$$w:=w-\eta 
abla Q_i(w)$$

# Gradient Descent with Backpropagation



- Initialize weights w<sub>0</sub>
- Repeatedly apply gradient descent:

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \gamma_n \quad \nabla E(\mathbf{w}_n)$$

 Stop when validation error is "small"

# Stochastic Gradient Descent with Backpropagation



- Initialize weights w<sub>0</sub>
- Randomly shuffle dataset
- For each batch i calculate gradient descent using backpropagation and apply  $w_{n+1} = w_n - \gamma_n \bigvee E_i(w_n)$
- Stop when validation error is "small"

# Parameter Training using SGD



# Parameter Training using SGD

Parameter update



Repeat for new batches of training data

# **Privacy**?

### Sensitive data Machine learning

Medical images Clinical records Text documents Personal photos Retail purchases

171



### **Model Inversion**



#### Fredrikson et al.

# Model Inversion in Action



slide 31

#### <u>Hitaj et al.</u>

# Deep Models under the GAN



# **Does Inference Breach "Privacy"?**



# **Recommended Reading**

# Frank McSherry. "Statistical inference considered harmful"



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

# Privacy in Statistical Databases

Individuals

Researchers



Large collections of personal information

- census data
- national security
- medical/public health
- social networks
- recommendation systems
- education



Utility: release aggregate statistics Privacy: ??? (intuition: individual information stays "hidden")


"Relax – it can only see metadata."



Slide credit: Adam Smith

#### **Remove Obvious Identifiers?**



# Membership Inference Attacks

• <u>Exact</u> high-dimensional summaries allow an attacker to test membership in a data set

[Homer et al. 2008]

- Caused US NIH to change data sharing practices for genomic data
- <u>Distorted</u> high-dimensional summaries allow an attacker to test membership in a data set

[Dwork et al. FOCS 2015]

#### Homer at al. Attack



# Machine Learning as a Service



# **Exploiting Trained Models**







# 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
- Synthetic: exploit target model Sample from inputs classified by the target model with high confidence input space target's training inputs

#### Important Point



 It is <u>not</u> the case that attack model simply learns to say that all inputs classified by the target model with high confidence belong to its training dataset

# Synthesizing Shadow Training Data



## Membership Inference Attack



Was this image part of the training set?





slide 50

## Next Step: Reconstruction





## Attack Success vs. Test-Train Gap



**Privacy:** Does the model leak information about data in the training set? <u>Learning</u>: Does the model generalize to data outside the training set?



data universe

# Generalizability Is Not Privacy

- Deep neural networks have huge memorization capacity
- A well-generalized model can still leak information about its training dataset
  - Good test performance on the primary task does not preclude good performance on another task (e.g., membership inference or reconstruction)



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 <u>and</u> accuracy



## "Classical" Intuition for Privacy

- Dalenius (1977): "If the release of statistics S makes it possible to determine the value [of private information] more accurately than is possible without access to S, a disclosure has taken place"
  - Privacy means that anything that can be learned about a respondent from the statistical database can be learned without access to the database

## Problems with Classical Intuition

- Popular interpretation: prior and posterior views about an individual shouldn't change "too much"
  - What if my (incorrect) prior is that every student in this room has three arms?
- How much is "too much?"
  - Can't achieve cryptographically small levels of disclosure and keep the data useful
  - Users <u>are</u> supposed to learn unpredictable things about the data

#### **Differential Privacy: Intuition**



If you change or remove one person's data, distribution of outputs should not change much

#### **Differential Privacy**



# **Deployed Differential Privacy**



#### Laplace Mechanism



- Intuition: f(x) can be released accurately when f is insensitive to individual entries  $x_1, \ldots x_n$
- Global sensitivity  $GS_f = \max_{neighbors x,x'} ||f(x) f(x')||_1$ 
  - Example:  $GS_{average} = I/n$  for sets of bits
- Theorem:  $f(x) + Lap(GS_f/\epsilon)$  is  $\epsilon$ -indistinguishable

- Noise generated from Laplace distribution

Lipschitz

constant of f

# **Achieving Privacy**

<u>Theorem</u>

If  $A(x) = f(x) + Lap\left(\frac{GS_f}{\varepsilon}\right)$  then A is  $\varepsilon$ -indistinguishable.

Laplace distribution  $Lap(\lambda)$  has density  $h(y) \propto e^{-\frac{\|y\|_1}{\lambda}}$ 

 $h(y+\delta) \bigwedge h(y)$ 

Sliding property of  $\operatorname{Lap}\left(\frac{\operatorname{GS}_{f}}{\varepsilon}\right)$ :  $\frac{h(y)}{h(y+\delta)} \leq e^{\varepsilon \cdot \frac{\|\delta\|}{\operatorname{GS}_{f}}}$  for all  $y, \delta$  *Proof idea:* A(x): blue curve A(x'): red curve  $\delta = f(x) - f(x') \leq \operatorname{GS}_{f}$ 



Impossible

- Suppose you know that I smoke
- Clinical study: "smoking and cancer correlated"
- You learn something about me
  - ... whether or not my data were used

### What Differential Privacy Means

# You learn (almost) the same things about me whether or not my data are used

No matter what you know ahead of time



## Privacy Concerns

- Training data is sensitive
  - speech, photo images, written documents
- Users have no control over the learning objective
- Using trained networks requires users to share their private data with service providers

# Possible Consequences

- Users' data might be used in wrong context
  - Compromises and data breaches
  - Inference of sensitive information
  - Training of intrusive models
- Holders of sensitive data cannot benefit from large-scale deep learning because they may not share or pool their datasets for training
  - Biomedical researchers?
  - Social scientists?




# Distributed Selective SGD (DSSGD)

#### Selective SGD



#### Share with others

#### Selective SGD



Share with others









## **Distributed Selective SGD**

- Local training, global convergence
- High training stochasticity
- Less overfitting

**Evaluation Datasets** 

#### **MNIST**







Task: Find the digit in the image — 10 class classification



# **Privacy Properties**

- Participants' datasets remain private
- Full control over parameter selection
- Known learning objective
- Resulting model available to all parties

### Indirect Information Leakage through gradient sharing



## Prevent Indirect Leakage

• Differentially private parameter selection and gradient sharing



# Sparse Vector Technique

• Select a small fraction of (perturbed) gradients that are above a given (perturbed) threshold



C. Dwork and A. Roth. The algorithmic foundations of differential privacy. *Theoretical Computer Science*, 9(3-4):211–407, 2013.



#### Federated Learning: Collaborative Machine Learning without Centralized Training Data

Thursday, April 06, 2017



