Babol Noshirvani
University of Technology

Privacy Preserving Machine Learning

Mohammad Hoseinpour
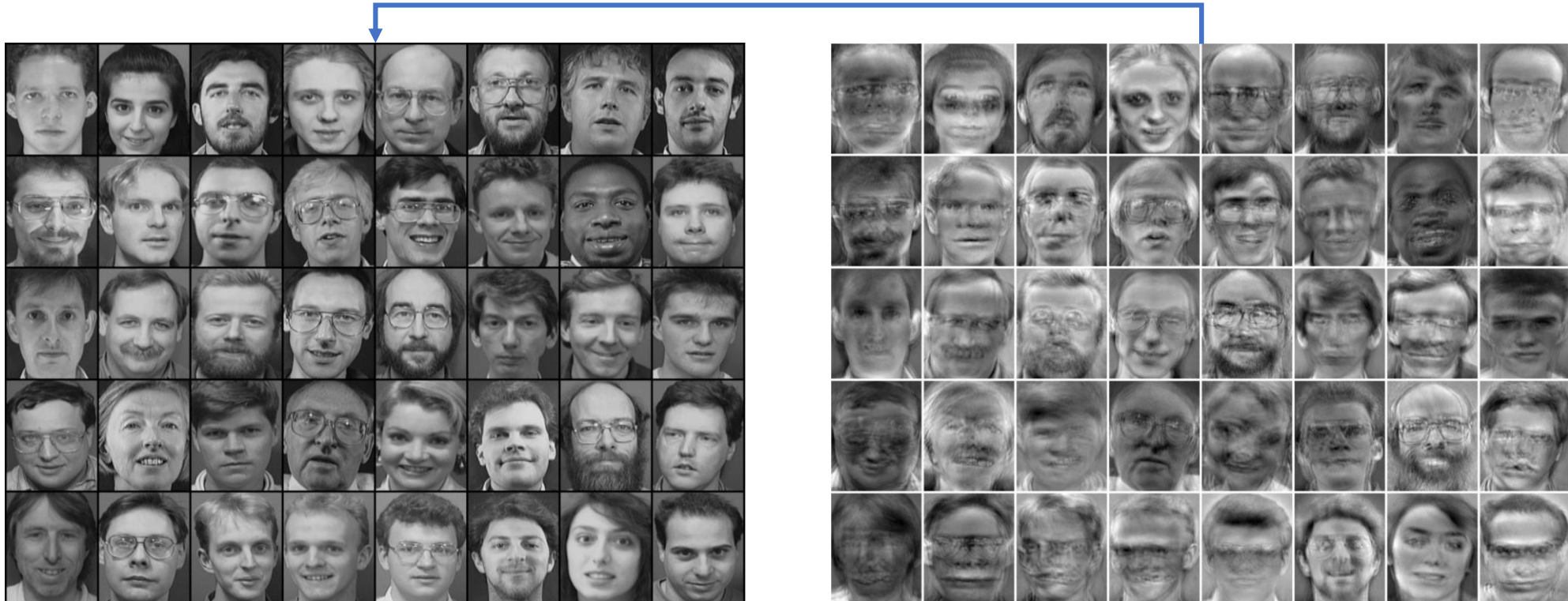
Supervisor : Prof. Ali Aghagolzadeh

September, 2022

# Motivation

- Machine Learning (ML) models can **memorize** training datasets

- Training ML models over **private datasets** can **violate** the **privacy of individuals**

- **Training data extraction attacks:**
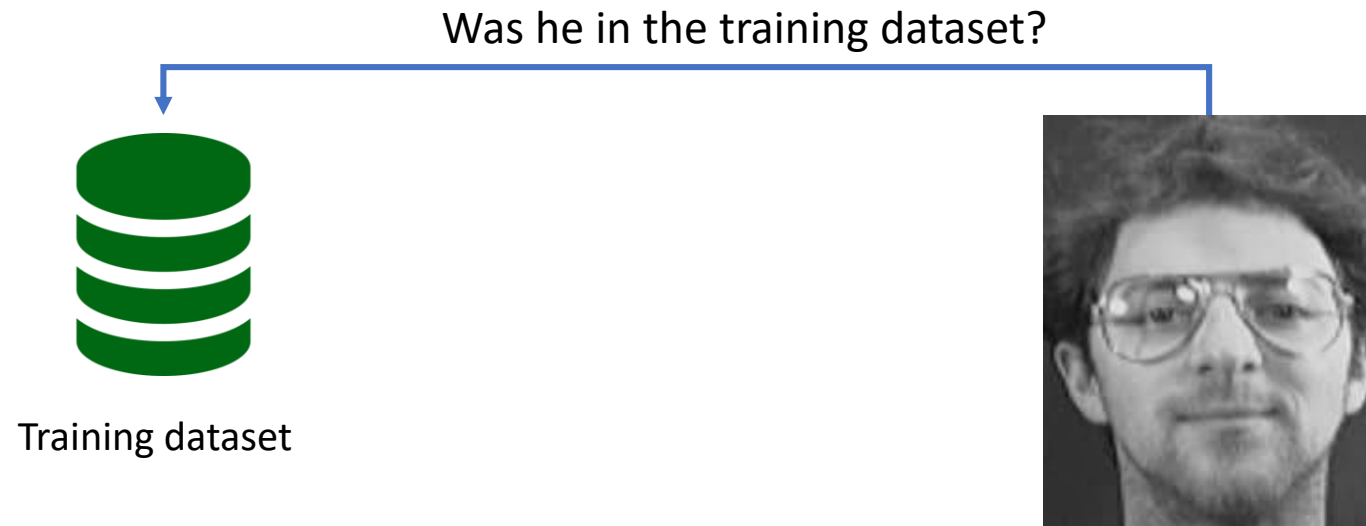  - Fredrikson et al. (2015), "Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures"


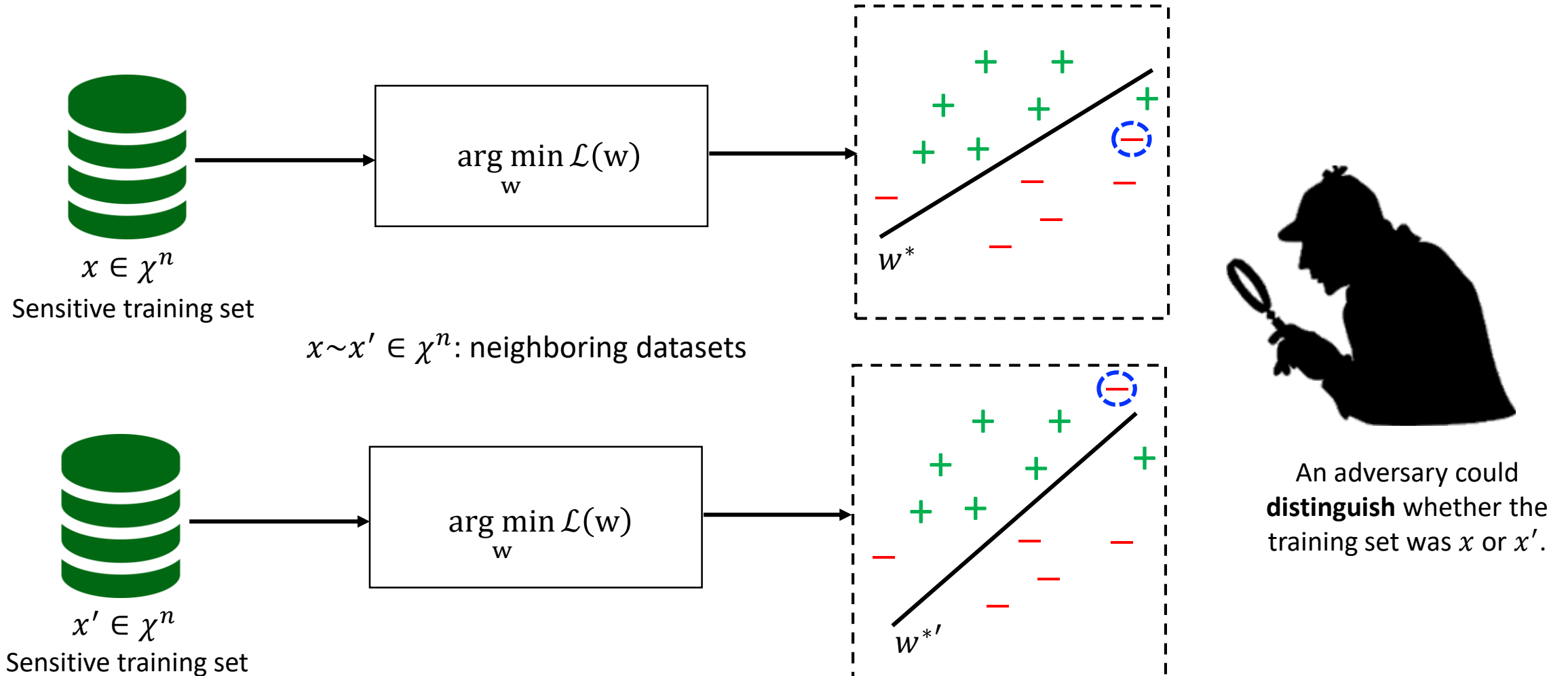
An example of model inversion attack

# Motivation

- Machine Learning (ML) models can **memorize** training datasets

- Training ML models over **private datasets** can **violate** the **privacy of individuals**

- **Membership inference attacks:**
  - Shokri et al. (2016), "Membership Inference Attacks Against Machine Learning"

Was he in the training dataset?



Training dataset

# Non-Private Logistic Regression

- The **decision boundary** of the classifier is **sensitive to the individual data points** in the training set.



$x \in \chi^n$

Sensitive training set

$x \sim x' \in \chi^n$: neighboring datasets

$x' \in \chi^n$

Sensitive training set

$\arg\min_{w} \mathcal{L}(w)$

$w^*$

$w^{*'}$

An adversary could **distinguish** whether the training set was $x$ or $x'$.
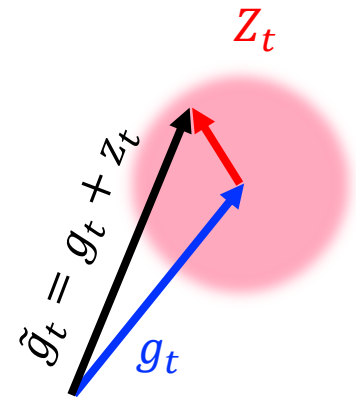
# Private Logistic Regression

- We apply **Gaussian mechanism** for privatizing the **updating rule** of the gradient descent

Empirical Risk: $$\mathcal{L}(\mathrm{w}) = \frac{1}{n}\sum_{i=1}^{n}\ell\big(\mathrm{w},(x_i,y_i)\big) + \lambda R(\mathrm{w})$$

---
**Algorithm 1** Noisy Projected Gradient Descent $(\mathcal{L},\mathcal{C},\eta,\sigma)$
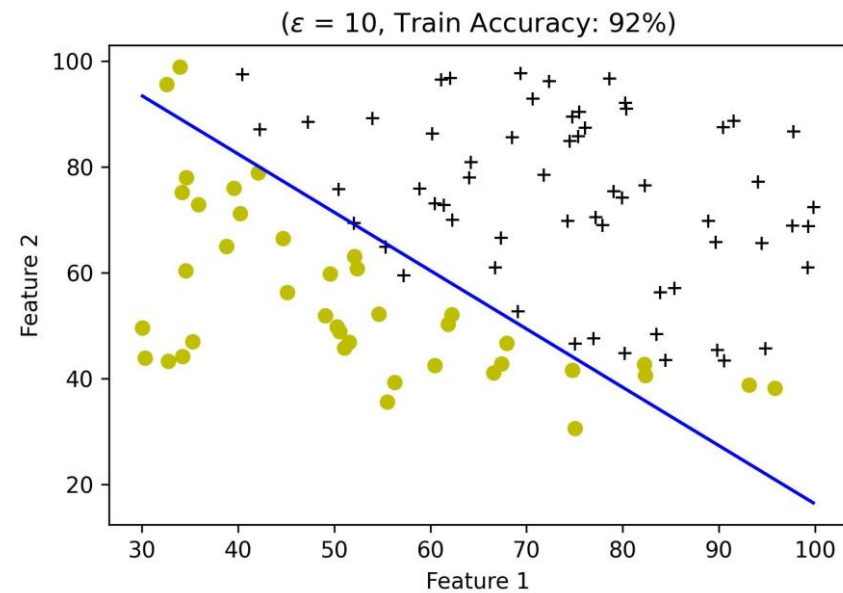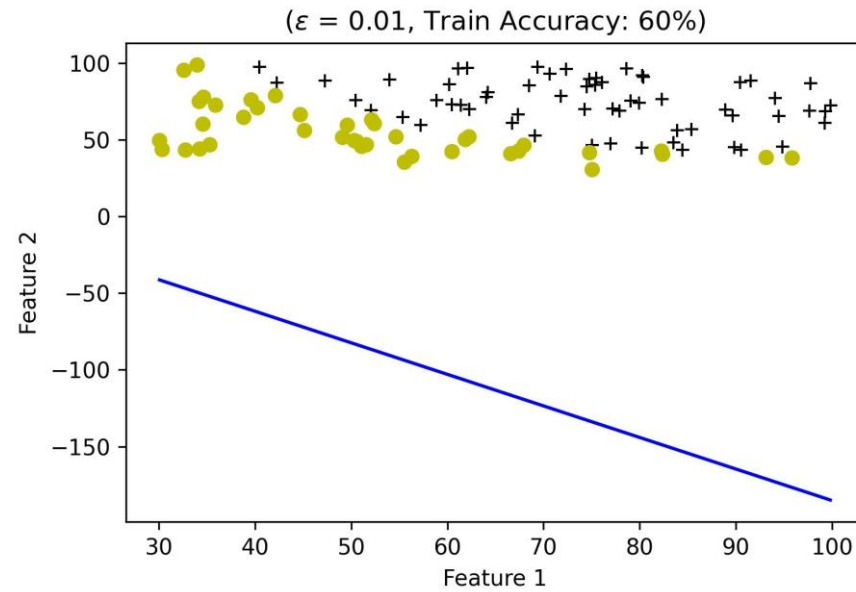
---
**Inputs**: Set $\mathcal{C} \subseteq \mathbb{R}^d$, noise parameter $\sigma$, learning rate $\eta$, loss function $\mathcal{L}(\mathrm{w})$.

1: $\mathrm{w}_0 \longleftarrow$ arbitrary point in $\mathcal{C}$;
2: **for** $t = 1, 2, ..., T$ **do**
3: $\quad g_t \longleftarrow \nabla\mathcal{L}(\mathrm{w}_{t-1})$;
4: $\quad \tilde{g}_t \longleftarrow g_t + \mathcal{N}(0,\sigma^2 I_d)$;
5: $\quad u_t \longleftarrow \mathrm{w}_{t-1} - \eta\tilde{g}_t$;
6: $\quad \mathrm{w}_t \longleftarrow \Pi_{\mathcal{C}}(u_t)$;
7: **end for**
8: **return** $\mathrm{w}_T$;

---



Perturbed gradient vector due to the additive Gaussian noise

# Results

# Thank You