

Membership Inference Attacks, Part II

April 15, 2025

① Exercise on means + variances.

from in-class exercises ≈ 2 weeks ago.

② Worst-case mechanism?

$$\mu \sim N(0, I_d)$$

Similar to exercise,
but leaves mechanism
uncertain about μ .

$$x_1, \dots, x_n \sim_{iid} N(\mu, I_d)$$

$$M(x) = \bar{x} + e \text{ where } \|e\|_2 \leq \alpha \sqrt{d} \quad (\text{RMSE } \alpha)$$

$$T(a, y) = \langle a - \mu, y - \mu \rangle$$

$$\text{OUT: } E(T|_{\text{OUT}}) = E \langle x - \mu, \bar{x} + e - \mu \rangle$$

$$= E \langle x - \mu, \bar{x} - \mu \rangle + E \langle x - \mu, e \rangle$$

↑ indep, mean 0

↑ mean 0, indep of e .

$$\text{Var}(T|_{\text{OUT}}) = d \text{Var}(x(1) - \mu(1)) \text{Var}(\bar{x}(1) - \mu(1)) + \text{Var}(N(0, \alpha^2 d))$$

use the fact

$$\text{that } \text{Var}(A+B) \leq 2(\text{Var}(A) + \text{Var}(B))$$

same as case of Gauss noise with $\rho \times d$.

if A, B indep with mean 0

$$\text{Then } \text{Var}(AB) = E(A^2 B^2)$$

$$= E(A^2) E(B^2)$$

$$= \text{Var}(A) \text{Var}(B)$$

$$\text{IN: } E(T|_{\text{IN}}) = E_i \langle x_i - \mu, M(x) - \mu \rangle$$

$$= \langle \bar{x} - \mu, \mu(x) - \mu \rangle$$

$$= \langle \bar{x} - \mu, \bar{x} - \mu \rangle + \langle \bar{x} - \mu, e \rangle$$

• Now condition on \bar{x} and e (where $\|e\|_2 \leq \alpha \sqrt{d}$)

• What is distrib of $x_i | \bar{X} = \bar{x}$?

Claim

From analyst's pov (that is, conditioned on x_1, \dots, x_n)

$$\mu \sim \frac{n}{n+1} \bar{x} + \frac{1}{n+1} z \quad \text{where } z \sim N(0, I_d)$$

~~indep of \bar{x}~~

Proof:

- $\text{Cov}(\mu, \bar{x}) = \begin{pmatrix} I & I \\ I & (1+\frac{1}{n})I \end{pmatrix}$
- Suppose that $\mu = a\bar{x} + bZ$ where $Z \sim N(0, I_d)$
 Then $\text{Cov}(\mu, \bar{x}) = \begin{pmatrix} a^2 \bar{x} + b^2 I & a \bar{x} \\ a \bar{x} & \bar{x} \end{pmatrix} \Rightarrow b^2 = \frac{1}{n+1}$
 $a^2 = \frac{1}{1+n} = \frac{n}{n+1}$

- So $\mu = \frac{n}{n+1} \bar{x} + bZ$.

- Furthermore, all the info about μ in x_1, \dots, x_n is contained in \bar{x} .
 (that is, \bar{x} is a "sufficient statistic" for μ)
 → To see why, note that

$$\begin{aligned}
 p(x_1, \dots, x_n | \mu) &\propto \prod_i \exp\left(-\frac{1}{2} \|x_i - \mu\|^2\right) \\
 &= \exp\left(-\frac{1}{2} \sum_i \|x_i - \mu\|^2\right) \\
 &= \exp\left(-\frac{1}{2} \sum_i \left(\|x_i - \bar{x}\|^2 + \|\bar{x} - \mu\|^2 + 2\langle x_i - \bar{x}, \bar{x} - \mu \rangle\right)\right) \\
 &= \exp\left(-\frac{1}{2} \left(n\|\bar{x} - \mu\|^2 + \sum_i \|x_i - \bar{x}\|^2 + 2\langle 0, \bar{x} - \mu \rangle\right)\right) \\
 &= \underbrace{\exp\left(-\frac{1}{2} n\|\bar{x} - \mu\|^2\right)}_{\propto P(\bar{x} | \mu)} \cdot \underbrace{\exp\left(-\frac{1}{2} \sum_i \|x_i - \bar{x}\|^2\right)}_{\propto P(x_1, \dots, x_n | \bar{x}, \mu)}
 \end{aligned}$$

QED

• We can now use the claim

$$\begin{aligned} \mathbb{E}(T | \mathcal{I}_N) &= \mathbb{E}\left(\|\bar{x} - \mu\|^2 + \langle \bar{x} - \mu, e \rangle\right) \\ &= d/n + \mathbb{E}\langle \bar{x} - \mu, e \rangle \end{aligned}$$

• Conditioning on \bar{x} and e , we get:

$$= \frac{d}{n} + \mathbb{E}_{\bar{x}, e} \mathbb{E}_{\bar{z}} \langle \bar{x} - \frac{n}{n+1}\bar{x} + \frac{1}{n+1}\bar{z}, e \rangle$$

$$= \frac{d}{n} + \mathbb{E}_{\bar{x}, e} \frac{1}{n+1} \langle \bar{x}, e \rangle$$

$$\geq \frac{d}{n} - \frac{1}{n+1} \underbrace{\|\bar{x}\| \cdot \|e\|}_{\approx \sqrt{d}} \underbrace{\leq \alpha \sqrt{d}}$$

$$\geq \frac{d}{n} (1 - \alpha)$$

which is less than $\frac{d}{n}$ for $\alpha \leq \frac{1}{2}$.