

- This block contains a dense handwritten note on differential privacy, covering various topics such as attacks, mechanisms, and specific protocols like BGW and PATEGAN.

1. Attacks

 - 1.1 Homogeneity attack:** A QI-block is ϵ -diverse if it contains at least ϵ well-represented values for the sensitive attribute S . A table is ϵ -diverse if every QI-block is ϵ -diverse.
 - 1.2 Background knowledge attack:** If R is a released table and Q is a query, R is said to satisfy ϵ -anonymity if and only if for each Q , Q appears with at least ϵ occurrences in R .
 - 1.3 RT (Row Transparency):** Approaches to make data ϵ -anonymous: 1. Generalization 2. Suppression 3. Find optimal ϵ -anonymized table with generalization and suppression is NP-Hard.
 - 1.4 Attacks against ϵ -anonymity:** 1. **Homogeneity attack:** A QI-block is ϵ -diverse if it contains at least ϵ well-represented values for the sensitive attribute S . A table is ϵ -diverse if every QI-block is ϵ -diverse.
 - 1.5 ℓ -diversity:** A QI-block is ℓ -diverse if it contains at least ℓ well-represented values for the sensitive attribute S . A table is ℓ -diverse if every QI-block is ℓ -diverse.
 - 1.6 Types:** 1. Probabilistic: Frequency of most frequent value in equivalence class is bounded by $\frac{1}{\epsilon}$. 2. Entropy: Entropy of the distribution of sensitive values in an equivalence class is at least $\log(\frac{1}{\epsilon}) \cdot H(X) = -\frac{1}{\epsilon} \sum p_i \log(p_i)$, where p_i is freq. of sensitive values in each class. 3. Recursive ((ϵ, δ)): $\ell < C(\epsilon + \delta + \dots + \epsilon m)$, ℓ - freq. of the 1st most freq. idea: the most freq. value v does not appear too frequently. 4. Checking: straightforward, creating NP-complete.
 - 1.7 Attacks:** 1. Skewness: distribution of sensitive attribute is invariance (for many records), every class E must have at least m rows. 2. The values for sensitive attributes S must all be different. 3. set of distinct sensitive values in each eq. class must be the same. methods: generalize quasi-identifiers - counterfactual tuples. others: same as the one ℓ -closeness idea: distribution of sensitive attribute values within each eq. class should be "close" to the distribution of sensit. attr. in the entire table. methods: 1. for numeric attributes: normalize the distance by range [min, max]. 2. for categorical: define distance based on a hierarchy. attacks: background knowledge.

2. Semantic Similarity

 - 2.1 Probabilistic Inference:** ϵ -invariance (for many records), every class E must have at least m rows. 2. The values for sensitive attributes S must all be different. 3. set of distinct sensitive values in each eq. class must be the same. methods: generalize quasi-identifiers - counterfactual tuples. others: same as the one ℓ -closeness idea: distribution of sensitive attribute values within each eq. class should be "close" to the distribution of sensit. attr. in the entire table. methods: 1. for numeric attributes: normalize the distance by range [min, max]. 2. for categorical: define distance based on a hierarchy. attacks: background knowledge.

3. Settings

 - 3.1 Players:** P_1, P_2, \dots, P_n wants to compute $f(x_1, x_2) = (y_1, y_2, \dots, y_n)$. They can communicate. Assumption: ideal world: T - completely trusted, incorruptible party. Each P_i sends x_i to T - T computes $f(x_1, x_2)$ and return result. In real world: there is no T , but there is a protocol π , which is secure if it emulates the ideal execution. So we simulate adversary \mathcal{A} : $\mathcal{A}(M, E) \rightarrow \mathcal{I} \stackrel{\text{def}}{=} \mathcal{P}(L, E) \rightarrow \mathcal{I} \subseteq \mathcal{I}$ is ideal world. Types of adversaries:
 - 1. Semi-honest (follows π , but nothing)
 - 2. Malicious (deviate from π)**3.2 BGW (Gentry, Halevi, Sahai):** Garbled circuits (GC) protocol. Idea: lookup table. 2. Oblivious transfer: R generates S , P and takes random key p_k . R sends $\{p_k, p_k\} - \rightarrow R$ to S . For $\delta = 1$: S returns $\{b_{00}, b_{01}, b_{10}, b_{11}\}$.
• R decrypts b_{ij} using S :
 - A sends to B : 1. circuit and the its key. 2. via OT ask keys for its input. $\text{Enc}_{K_A}(b_{ij})$.
 B decrytp and sends output to A . Security: semi-honest model (OT secure in semi-honest) / loss security if we evaluate a circuit more than once with diff. outputs.
 - 3.3 BGW:** Idea: arithmetic functions with more than two inputs \Rightarrow arithmetic circuits. We can represent any polynomial with a combination of addition and multiplication gates. We operate on \mathbb{F}_q . Idea: secret S , that we split into n shares - only $t+1$ is enough to reconstruct S . Shamir's secret sharing: n players, threshold $t+1$, n selects finite field \mathbb{F}_q , $q = n^t$. Secret is $S \in \mathbb{F}_q$. Select n distinct elements $P_1, P_2, \dots, P_n \in \mathbb{F}_q$. Sample random polynomial p with degree t and $p(0) = S$ - compute the share of player $i: s_i = p(P_i)$ reconstruct Lagrange polynomial interpolation with any $t+1$, and evaluate $p(0) = S$. With multiplication degree of polynomial $pc \leq 2t$. $\leq 2t$ degree reduction step: $p(c) = \sum_{i=0}^{2t} a_i c^i$; each player holds one of c_i : each player generates and distributes shares c_j of their c_i . Each computes $c_j \in \mathbb{F}_q$. c_j is a share of c_i in a secret sharing with that wire. Semi-honest model: $2t+1$ users; adversary can corrupt up to $t+1$ parties. Zero-knowledge proofs: GMVW compiler: run π and prove in zero-knowledge that every msg. is the result of running π honestly. Can transform any SAC from semi-honest to malicious model.
 - 3.4 Federated SGD:** Clients C_1, \dots, C_n , input $\{(x_i^1, y_i^1)\}$. Server chooses loss L , initializes weights w . Chooses J for iteration $i \in [0, \dots, J]$. Server samples a random fraction of clients, C_J , sends them weights w . Each client C_j : compute loss for sample $L_j = L(\{x_i^j, y_i^j\})$. Compute gradient vector $g_j = \nabla_{w_j} L_j$. Send g_j to server. Server compute avg gradient $g_t = \frac{1}{|C|} \sum_j g_j$. Server update $w \leftarrow w - \eta g_t$.
 - 3.5 Attacks:** Adversary models: external server, clients. 1. Honest-but-curious server-target membership inference (server records an model update (gradient) from all C at epoch, after training membership for each C). 2. Honest-but-curious server gradient inversion (adversary steals training data from C , updates dummy inputs and labels to min distance between gradients or add more points so client grad. minimizes training). 3. Malicious user or server: gradient ascent (single datapoint per dataset, report inverse gradient, so next iter. someone SGD will abruptly reduce grad. of loss).
 - 3.6 Privacy FL with central DP (ex. Gboard):** protect against honest-but-curious server. algorithm (combination of DPSGD and federated SGD, user clip grad and server avg grads and applies Gaussian noise)
 - 3.7 FL with local DP:** each user adds DP to their grad. threat model: malicious server, user drawbacks: need very large # users and amount of noise. 3. FL with secure aggregation: malicious server, so use secure multi-party computation to agg grad. drawbacks: no DP (memory overhead); FL with distributed DP are secure aggregation + hybrid privacy preserving FL. 4. Local DP: each iteration: one epoch of DPSGD (central DP not local) on each client. gradients shared with threshold homomorphic encryption. Assume t honest clients. Accuracy is high, but secure agg allows us to use less noise. 2. Distributed discrete gaussian for FL with secure aggregation: Secure agg is a black box - plug in any SMC scheme, including ones for malicious server. Client distributes data to reduce communication overhead (fewer bits in distributed-reduced accuracy), reduced communication; client apply discrete Gaussian noise, achieves ϵ -concentrated DP.
 - 3.8 Synthetic data generation:** one column: distribution of data, histograms. 2. # of unique values; add Lap. noise, parallel composition, sample points. multi-column: n-way marginal distribution, ok for small data; for big, low coverage (noise overwhelms the signal).
 - 3.9 Priv-Bayes:** (mitigate low data counts) represent correlations in Bayesian Network or Markov random fields (better acc) derive low dim. table; generate DP version of it, sample.
 - 3.10 DP Bayesian Network:** input d attributes A_1, A_2, \dots, A_d ; initialize empty Bayesian network N . Consider all possible (k, l) attribute pairs A_i, A_j and evaluate mutual info $I(A_i, A_j | A_1, \dots, A_{i-1})$. Add DP noise into mutual info. Add selected $A_1, \dots, A_k \rightarrow A_l$ into N . Repeat with A_2, \dots, A_N .
 - 3.11 PATEGAN (GAN-like adversarial networks):** random noise \rightarrow generator \rightarrow discriminator \rightarrow noise. Add PATE into GANs.
 - 3.12 Implementation challenges:** both implementation determine value ranges for columns from raw data to break ID. Win approach of MIST(MIST): measure 4, 2, 3-way marginals using gaussian mechanism.
 - 3.13 Synthetic data that has those marginals, using graphical model. 3. generate data samples and measure performance (lower score better). Future? Other data type, image/time data; with rare events.**
 - 3.14 Privacy design strategies:** 1. minimize processing of personal data & separate (isolate data) & abstract (how data was processed). 2. hide (protect personal data) & inform (subject about processing their data) & control (subject can control data) & enforce (policy commitment) & demonstrate (document, report results).
 - 3.15 DevPrivOps:** 1. Plan (select privacy design strategies, threat modeling, apply LINDDUN framework, risk analysis) 2. Code (iterate lib that align with design objectives, document, TILT(TIRA)) 3. Log (monitoring) 4. Build (flag outdated versions, create distinct versions for diff. groups) 4. Test 5. Release 6. Deploy (use: foreign jurisdictions, use canary releases and A/B test to mitigate risk of unauthorized data processing) 7. Operate (penetration testing, observation and cross-validation against prev. assumption) & monitor (logging, auto-report).
 - 3.16 TILT:** Transparency info language and toolkit - standardized machine-readable format for transparency info in accordance with EU data protection law - data subjects get into they need (choice of interpretation) & data controller collects and process info efficiently. Ex: GDPR obligation \rightarrow categorization \rightarrow detailed analysis \rightarrow formal language def \rightarrow machine-readable format \rightarrow uniformly represented transparency info.
 - 3.17 TIRA:** extends OpenAPI specification to allow developers to annotate transparency info.

4. Mechanisms

 - 4.1 Settings:** database D - query S with the response R_s - mechanism $M: S \rightarrow R_s$. Without leaking sensitive info \rightarrow without leaking sensitive info.
 - 4.2 Naive mechanism:** always reports correct and try to limit query types. **1. Differencing attack:** (aggregating over large groups) \rightarrow add noise \rightarrow independent noise \rightarrow attack with repeated queries (median filter) \rightarrow Dinur-Nissim reconstruction attack. Such mechanisms are called **blatantly non-private**, adding η noise is blatantly non-private.
 - 4.3 Differential privacy:** $D: A \rightarrow D \rightarrow T$. Adversary gets access to T and details of A . Adversary should learn nothing new about an individual after seeing A and T , regardless any other background knowledge. A mechanism $M: X \rightarrow Y$ is ϵ -diff. private (ϵ -DP) if for any two neighboring D_1 and D_2 (the same except one row): $\forall T \subseteq Y. P[\exists (M(D_1), T)] \leq e^\epsilon P[\exists (M(D_2), T)]$. Smaller $\epsilon \Rightarrow$ more privacy; for small $\epsilon: e^{-\epsilon} \approx 1 + \epsilon$. Properties: 1. Safety against post-processing 2. Sequential composition: $e \cdot e = e$.
 - 4.4 Parallel composition:** $(\lambda x. x)(\lambda y. y)$ \rightarrow $\lambda z. z$ \rightarrow $\lambda y. y$ \rightarrow $\lambda x. x$. Group privacy: $P[\exists (M(D_1), T)] \leq e^\epsilon P[\exists (M(D_2), T)]$, D_1 and D_2 differ in k rows.
 - 4.5 Mechanisms for DP:** 1. Laplace mechanism: $P[\exists (M(D), T)] = e^{-\epsilon} P[\exists (M(D'), T)]$, D and D' differ in 1 row. Sensitivity $S(q)$ is the smallest number so that for $V(D), D'$: $q(D) - q(D') \leq S(q)$. Def. ϵ -sensitivity: $\Delta = \max_{D, D'} |q(D) - q(D')| \leq S(q) \cdot \epsilon \Rightarrow \epsilon$ -diff. privacy. Counting query: $S(q) = \sum_{x \in D} q(x)$. Average unbounded mechanism: Clipping 2. Exponential mechanism: D : R : set of outputs with respect to best element \mathbf{d} . Scoring function: $u: R \rightarrow \mathbb{R}$ with sensitivity Δ_u : $\Delta_u = \max_{D, D'} |u(D, R)|$. Exp. mechanism: given input D , set of output R , scoring function u samples output R' with $P[R' = r] = \exp(\frac{u(r)}{\Delta_u})$.
 - 4.6 Unit of Privacy:** One person. Person - data. Some examples: 1. U.S. broadband coverage dataset (unit: one device, $\epsilon = 0.01$) 2. Post-secondary employment outcomes (unit: person, $\epsilon = 3$) 3. Approximate differential privacy: $VTC(Y) P[\exists (M(D), T)] \leq e^\epsilon P[\exists (M(D'), T)]$. Typical choice: $\epsilon = \frac{n}{10} = 1$.
 - 4.7 Properties:** same as DP (seq. composition, post-processing, parallel composition). Mechanism 1. Gaussian mechanism - for a query function $f(x)$: $G(x) = f(x) + \mathcal{N}(0, \sigma^2)$. Gaussian M satisfies (ϵ, δ) -diff. privacy: $\epsilon = \frac{2\sigma}{\Delta} \sqrt{2 \ln(1 + \frac{1}{\delta})}$, where Δ is the sensitivity of $f(x) = \mathcal{N}(0, \sigma^2)$. Gaussian distribution: $\mathcal{N}(0, \sigma^2)$. Vector-valued functions $\mathbf{f}: D \rightarrow \mathbb{R}^k$: examples: ML, histograms. Sensitivity Δ_f : use maximum magnitude $\Delta_f = \max_{D, D'} \|\mathbf{f}(D) - \mathbf{f}(D')\|_2$. Laplace mechanism: $\|\mathbf{f}(D) - \mathbf{f}(D')\|_2 \leq \Delta_f \cdot \epsilon$. Δ_f sensitivity generally smaller \rightarrow more privacy. Examples: 1. LinkedIn audience engagement api (user-month analysis) 2. Urban mobility data (unit: whether an individual made a trip from A to B during week w). Advanced composition theorem: if each m ϵ -DP then ℓ -fold adaptive composition satisfies (ϵ', δ') -DP. For $\ell \geq 2, \epsilon' = 2e^{\ell} \sqrt{2 \ln(1 + \frac{1}{\delta'})}$. If $m \geq \ell$, then $\ell \cdot \epsilon' - \Delta_m \leq \epsilon' \leq \ell \cdot \Delta_m$. Advanced better than Sequential with bigger ℓ . Local sensitivity: $LS(f, x) = \max_{x' \in \mathbb{R}^d} |f(x) - f(x')|$. Local sensitivity on an actual dataset \rightarrow cannot be used to achieve DP. Propose test - release analysis: proposed on upper bound on local sensitivity + framework test whether the dataset is sufficiently far from a dataset with higher local sensitivity: $AL(f, \epsilon, \delta) = \max_{x \in \mathbb{R}^d} |f(x, b) - f(x, b')|$. Test: $\epsilon \cdot AL(f, \epsilon, \delta) \leq \frac{1}{2} \log(\frac{1}{\delta})$. PTR satisfies (ϵ, δ) -DP.
 - 4.8 Smooth Sensitivity:** use an approximation of local sensitivity to calibrate noise. Considers local sensitivity of all entries in the neighborhood of x : $\mathbb{E}[f(x) - f(x')] = \mathbb{E}[f(x) - \mathcal{N}(f(x), \sigma^2)]$. Larger than local \rightarrow Real differential privacy (RDP): $D_\epsilon = \frac{1}{\epsilon}$. In $E[\ln(\frac{1}{\epsilon})]$, $\Delta \cdot \ln(\frac{1}{\epsilon})$. Randomized mechanism M satisfies (ϵ, δ) -RDP if for $x, x' \in D$: $|M(x) - M(x')| \leq \Delta$. Any M which is (ϵ, δ) -RDP also is (ϵ, δ) -DP for any $\epsilon, \delta = \epsilon - \delta$. To achieve this we can use Gaussian mechanism: $f(x) = f(x) + \mathcal{N}(0, \sigma^2) \Rightarrow \Delta^2 = \frac{\Delta^2 \cdot \epsilon^2}{2\epsilon^2 - \Delta^2} \cdot \Delta^2 \cdot \epsilon^2 \cdot \delta^2 \leq \Delta^2 \cdot \epsilon^2 \cdot \delta^2$. Stronger than RDP, but all \times . Conversion to (ϵ, δ) -DP: $\epsilon' = \epsilon \cdot \ln(\frac{1}{\epsilon})$, $\delta' = \delta \cdot \ln(\frac{1}{\epsilon})$. Stronger than RDP, but all \times . Conversion to (ϵ, δ) -DP: $\epsilon' = \epsilon \cdot \ln(\frac{1}{\epsilon})$, $\delta' = \delta \cdot \ln(\frac{1}{\epsilon})$.
 - 4.9 Settings:** Players P_1, P_2, \dots, P_n want to compute $f(x_1, x_2) = (y_1, y_2, \dots, y_n)$. They can communicate. Assumption: ideal world: T - completely trusted, incorruptible party. Each P_i sends x_i to T - T computes $f(x_1, x_2)$ and return result. In real world: there is no T , but there is a protocol π , which is secure if it emulates the ideal execution. So we simulate adversary \mathcal{A} : $\mathcal{A}(M, E) \rightarrow \mathcal{I} \stackrel{\text{def}}{=} \mathcal{P}(L, E) \rightarrow \mathcal{I} \subseteq \mathcal{I}$ is ideal world. Types of adversaries:
 - 1. Semi-honest (follows π , but nothing)
 - 2. Malicious (deviate from π)**4.10 Attacks on ML:** 1. Membership inference 2. Shadow models (feature or stat. properties); member or not. 3. Nodal inversion (feature or stat. properties; feature) 4. Property inference (reconstruct properties from training data; SVD; gradient from L1)
 - 4.11 Training data extraction & Model extraction:**
 - 4.12 DP-SGD:** Add DP noise to computation of gradient.
 - 4.13 Comp. avg. gradient:** $g_t = \frac{1}{|B|} \sum_{i \in B} \nabla_{w_i} L_i$, $\Delta = \frac{2\sigma^2 \ln(1 + \frac{1}{\delta})}{\epsilon^2}$.
 - 4.14 Clipping:** scaling vector-element wise division by $\|z\|_2$ norm results in $\|z\|_2$ norm (sensitivity) of 1.
 - 4.15 Clipping threshold:** C : median of the norm of the unscaled gradient.
 - 4.16 Privacy accounting:** strong composition: $\Delta^* = \frac{2\sigma^2}{\epsilon^2}$, each step is (η, η) -private, $q = \frac{\Delta}{\eta}$. For T iterations, the training is $(\eta + T \log(\frac{1}{\delta}), \eta)$ -private.
 - 4.17 2 moments accountant:** better method to analyze privacy cost, tighter bound $(\eta + T \sqrt{\eta}, \delta)$ -private. Same privacy can train longer.
 - 4.18 Algorithm:** Sample mini-batch B : For each $i \in B$: $g_i = \nabla_{w_i} L_i$. Clip gradient: $\tilde{g}_i = \max(\frac{1}{C}, \frac{|g_i|}{C}) \cdot g_i$. Add noise, avg. gradients: $\tilde{g}_t = \frac{1}{|B|} \sum_i \tilde{g}_i + \mathcal{N}(0, \sigma^2 \cdot C^2)$. Update weights: $w_{t+1} = w_t - \eta \tilde{g}_t$. Output w_{T+1} , compute ϵ, δ .
 - 4.19 Issues:** fixed choice of optimiser, params to tune, batch size and selection influence privacy guarantee.
 - 4.20 PATE (Private Aggregation of Teacher Ensembles):** Independent of ML algo and of optimiser, loss, batch size.
 - 4.21 Main phases:** 1. teacher training (data in disjoint set, each used for train, $f_i(x)$ without privacy) 2. aggregation (label count n_j for class j is # of teacher models that predict class, label with highest count is agg. output; add Lap noise: $f(x) = \text{argmax}_j f_j(x) + \mathcal{N}(0, \frac{1}{n})$) 3. knowledge transfer (to get around limit number of user queries within; learns how to mimic agg. output of teacher ensemble; control privacy by limiting how many data records labeled).
 - 4.22 Sparse vector technique:** idea: perform a stream of queries with sensitivity Δ ; release only the identity of the first query which passes a test. Above threshold algorithm: inputs q_1, \dots, q_n , Δ , ϵ , $\hat{f} = T + \mathcal{N}(0, \frac{1}{\Delta})$. For i in q : $\hat{f}(q_i) \geq \Delta$ \rightarrow q_i is a hit. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a miss. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ \rightarrow q_i is a true negative. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a false negative. $\hat{f}(q_i) < \Delta$ \rightarrow q_i is a false positive. $\hat{f}(q_i) = \Delta$ \rightarrow q_i is a true positive. $\hat{f}(q_i) > \Delta$ $\rightarrow</$

① Introduction

1. Privacy types:

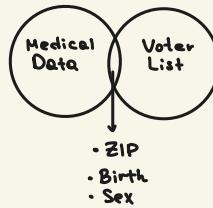
- 1. Location
- 2. State of Body and Mind
- 3. Behavior + Action
- 4. Social life
- 5. Media

2. Privacy properties

Properties		Threats
a) HARD:	{ 1. Unlinkability 2. Anonymity and Pseudonymity 3. Plausible deniability 4. Undetectability 5. Confidentiality	Linkability Identifiability Non-repudiation
b) Soft	{ 1. Content awareness 2. Policy and consent compliance	Detectability Disclosure of information content Unawareness policy and consent Noncompliance

3. Example: Medical Records (Group Insurance Commission)

- Voter list was bought in addition (← Adversary has side info)



4. Example: NYC TAXI Data

- fare data and trip data

5. Example: Netflix (ranking competition)

6. Example: Apple Quicktype

7. Example: Genomics (SNPs)

8. Example: Sugarbeet auction

② k-anonymity

① Aim: publish a lot of data without harming the privacy

② Approach: 1. Naive Anonymization

- remove all identifiers

- linked attacks (example: medical records)

③ k-anonymity:

- 3 types of attributes:

- Identifiers (PII)

- Quasi-identifiers (can be identifying in combination)

def Let $T(A_1, \dots, A_n)$ be a table. A quasi-id. of T is a set of attributes $\{A_i, \dots, A_j\} \subseteq \{A_1, \dots, A_n\}$ whose release must be controlled.

- Sensitive attributes

- Idea: create anonymity sets, so that each individual cannot be distinguished from at least $k-1$ others. (its about quasi-identifiers)

- def** k-anonymity: Let $RT(A_1, \dots, A_n)$ be a released table and

Q_{RT} be the quasi-id. RT is said to satisfy k-anonymity if and only if for each $Q_1 \in Q_{RT}$ each sequence of values in $RT[Q_{1RT}]$ appears with at least k occurrences in $RT[Q_{1RT}]$

- Approaches to make data k-anonymous:

1. Generalization

2. Suppression

- Find optimal k-anonymized table with generalization and suppression is

NP-Hard.

④ Attacks against k-anonymity

1. Homogeneity attack

2. Background knowledge attack

⑤ l-diversity

- A Q_1 -block is l-diverse if it contains at least l well-represented

values for the sensitive attribute S. A table is l-diverse if every Q_1 -block is l-diverse.

- Types:

1. Probabilistic

- Freq. of the most frequent value in an equivalence class is bounded by $\frac{1}{e}$

2. Entropy

• entropy of the distribution of sensitive values in each equivalence class is at least $\log_2(e)$

• $H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$, X - freq. of sensitive values in each eq. class level of uncertainty

3. Recursive (c, ϵ) - diversity

• $r_i < c (r_1 + r_2 + \dots + r_m)$, r_i - freq. of the i -th most freq. value

• idea: the most freq. value r_i does not appear too frequently.

- Checking for 1-diversity straightforward, but creating NP-complete.
- Attacks:

1. Skewness

• Overall distribution of sensitive values

2. Semantic Similarity

3. Probabilistic Inference

⑥ m-invariance (for multi releases)

1. every eq. class E must have at least m rows.

2. the values for sensitive attributes s must all be different

3. set of distinct sensitive values in each eq. class must be the same.

- methods:

- generalize quasi-identifiers

- counterfeit tuples

• attacks the same

⑦ t-closeness

• idea: distribution of sensitive attribute values within each eq. class should be $\boxed{\text{"close"}}$ to the distribution of senset. attrib. in the entire table.

- methods:

1. for numeric attributes: normalize the distance by range
$$\frac{|i-j|}{\text{range}}$$

2. for categorical: define distance based on a hierarchy

• attacks: - background knowledge

③ Differential Privacy

① Setting:

- database D
- query S with the response R_S
- mechanism M: S → R_S
 - ↳ 1. give statistically useful responses (utility)
 - 2. without leaking sensitive info (privacy)

② Naive mechanism

- a) always report correct and try to limit query types

↳ difference attack (aggregating over large groups)

- b) add noise

↳ independent noise → attack with repeated queries (median filter)

2. too little noise → Dinur-Nissim reconstruction attack

- Such mechanisms are called blatantly non-private,
- Adding O(\sqrt{n}) noise is blatantly non-private.

③ Differential privacy

- D, A: D → T

- Adversary gets access to T and details of A

- adversary should learn nothing new about an individual after seeing A and T, regardless any other background knowledge

• def A mechanism M: X → Y is ϵ -diff. private (ϵ -DP) if for any two neighboring D₁ and D₂ (the same except one row):

$$\forall T \subseteq Y: \Pr[M(D_1) \in T] \leq e^\epsilon \Pr[M(D_2) \in T]$$

- Smaller ϵ ⇒ more privacy; for small ϵ : $e^\epsilon \approx 1 + \epsilon$

Properties:

1. Safety against post-processing
2. Sequential composition ($\epsilon_1 + \epsilon_2 + \epsilon_3 = k \cdot \epsilon$)
3. Parallel composition ($X = x_1, x_2, \dots; M(x_1), M(x_2) = \epsilon$ -diff.)
4. Group privacy: $\Pr[M(D) \in T] \leq e^{k \cdot \epsilon} \Pr[M(D_2) \in T]$, D₁ and D₂ differ in k rows

④ Mechanisms for DP:

1. Laplace mechanism:

$A(S) + \eta$, η from $\text{Lap}(\mu, b)$, $f(x) = \frac{1}{2b} \exp\left(\frac{-|x-\mu|}{b}\right)$, $\sigma^2 = 2b^2$

[def] sensitivity $S(q)$ is the smallest number so that for D_1, D_2 :

$$|q(D_1) - q(D_2)| \leq S(q)$$

[def] ℓ_1 -sensitivity: $\Delta_1 = \max_{D_1, D_2 : d(D_1, D_2) \leq 1} |q(D_1) - q(D_2)|$

- so $b = \frac{\epsilon}{2} \Rightarrow \epsilon$ -diff. privacy

- Counting query: $S(q) = 1$
- Sum query and Average unbounded

- Mechanism: Clipping

2. Exponential mechanism:

- D, R - set of outputs from which we want to select best element

- [def] Scoring function: $u: D \times R \rightarrow \mathbb{R}$ with sensitivity Δu

$$\Delta u = \max_{R, D, D'} |u(D, R) - u(D', R')|$$

- Exp. mechanism:

- given input d , set of output R , scoring function u
samples output $r \in R$ with $P_r \sim \exp\left(\frac{E \cdot u(d, r)}{2\Delta u}\right)$

⑤ Unit" of Privacy:

- One person
- Person-day

⑥ Some examples:

1. U.S. Broadband coverage dataset (unit - one device, $\epsilon = 0.2$)

2. Post-secondary employment outcomes (unit - person, $\epsilon = 3$)

⑦ Approximate differential privacy

$$\forall T \subseteq Y: \Pr(M(D) \in T) \leq e^\epsilon \Pr(M(D') \in T) + \delta \quad (\frac{1-\epsilon}{\delta} - \text{DP})$$

- Typical choice: $\delta < \frac{1}{n^2}$, $|D| = n$

- Properties: same as DP (Seq. composition, Post-processing, parallel composition)

- Mechanism

1. Gaussian mechanism

- For a query function $f(x)$: $G(x) = f(x) + \mathcal{N}(\sigma^2)$

$$\frac{2\sigma^2(\ln \frac{1.25}{\delta})}{\epsilon^2}$$

- Gaussian M satisfies (ϵ, δ) -diff. privacy if $\delta = \frac{2\sigma^2(\ln \frac{1.25}{\delta})}{\epsilon^2}$, where
 - σ is the sensitivity of $f(x)$
 - $N(\delta^2)$ - gaussian distribution

⑧ Vector-valued functions

$$f: D \rightarrow \mathbb{R}^k$$

- examples: ML, histograms

- sensitivity Δ_1^q : use maximum magnitude

$$\Delta_1^q = \max_{D_i, D_j: d(D_i, D_j) \leq 1} |q(D_i) - q(D_j)|, \quad || \cdot || - \text{chosen norm}$$

- Magnitude of a vector:

$$1. L_1 \text{ norm: } ||V||_1 = \sum |V_i| \quad \leftarrow \text{Laplace mechanism}$$

$$2. L_2 \text{ norm: } ||V||_2 = \sqrt{\sum V_i^2} \quad \leftarrow \text{Gaussian}$$

- L_2 sensitivity generally smaller \Rightarrow more privacy

⑨ Examples:

1 LinkedIn's audience engagements api (user-month-analyst)

2. Urban mobility data (unit: whether an individual made a trip from A to B during week w)

⑩ Advanced composition:

1. k-fold adaptive composition

- seq. of mechanism

- each m_i is chosen adaptively, based on the outputs of m_{i-1}

- examples: loops, recursive function

• Advanced composition theorem:

- if each m_i ϵ -DP then k-fold adaptive composition

satisfies (ϵ', δ') -DP for $\delta' \geq 0$, $\epsilon' = 2\epsilon \sqrt{2k \cdot \log(\frac{1}{\delta'})}$

• if m_i (ϵ, δ) -DP, then $(\epsilon', k\epsilon + \delta')$ -DP, $\delta' \geq 0$ and $\epsilon' = 2\epsilon \sqrt{2k \cdot \log(\frac{1}{\delta'})}$

• Advanced better than Sequential with bigger k

⑪ Local sensitivity

$LS(f, x) = \max_{x' \sim d(x, x)} |f(x) - f(x')|$ - max. difference on an actual dataset

• cannot be used to achieve DP

⑫ Propose - test - release

- analyst proposes an upper bound on local sensitivity

- framework test whether

$$1. A(f, x, k) = \max_{y: d(x, y) \leq k} LS(f, y)$$

$$2. D(f, x, b) = \underset{k}{\operatorname{argmin}} A(f, x, k) > b$$

$$3. \text{Test: } D(f, x, b) + \text{Lap}\left(\frac{1}{\epsilon}\right) < \frac{\log\left(\frac{2}{\delta}\right)}{2\epsilon}$$

$$4. f(x) + \text{Lap}\left(\frac{1}{\epsilon}\right)$$

- PTR satisfies (ϵ, δ) -DP

(13) Smooth Sensitivity

- use an approximation of local sensitivity to calibrate noise

- considers local sensitivity of all datasets in the neighborhood of the x

$$\bullet \beta = \frac{\epsilon}{2\log(2/\delta)}, S = \max e^{-\beta k} A(f, x, k), \text{ release } f(x) + \text{Lap}\left(\frac{2S}{\epsilon}\right)$$

- larger than local

(14) Renyi differential privacy (RDP)

$$\bullet D_\alpha(P||Q) = \frac{1}{\alpha-1} \ln E_{x \sim Q} \left(\frac{P(x)}{Q(x)} \right)^\alpha, \alpha=0 \Rightarrow \text{DP}$$

- Randomized mechanism M satisfies $(\alpha, \bar{\epsilon})$ -RDP if for $\forall x, x'$

$$D_\alpha(M(x)||M(x')) \leq \bar{\epsilon}$$

- Any M which is $(\alpha, \bar{\epsilon})$ -RDP also is (ϵ, δ) -DP for any δ , $\epsilon = \bar{\epsilon} + \frac{C_1(\bar{\epsilon})}{\alpha-1}$

- To achieve this we can use Gaussian mechanism:

$$F(x) = f(x) + \mathcal{N}(0, \sigma^2), \sigma^2 = \frac{\Delta f^2 \alpha}{2\bar{\epsilon}}, \Delta f - L_2 \text{ sensitivity}$$

- $G(x) = (f_1(x), f_2(x))$ satisfies $(\alpha, \bar{\epsilon}_1 + \bar{\epsilon}_2)$ -RDP, so

lower privacy cost with RDP Compose

(15) Zero-concentrated differential privacy

$$\bullet D_\alpha(M(x)||M(x')) \leq \rho_d, \text{ all } \alpha \in (1, \infty)$$

- stronger than RDP, cut off α

- Conversion to (ϵ, δ) -DP: $\delta > 0$ and $\epsilon = \rho + 2\sqrt{\rho \log\left(\frac{1}{\delta}\right)}$

$$\bullet F(x) = f(x) + \mathcal{N}(0, \sigma^2), \sigma^2 = \frac{\Delta f^2}{2\rho};$$

$$\bullet G(x) = (f_1(x), f_2(x)) : (\rho_1 + \rho_2) - 2 \text{CDP}$$

(16) Sparse vector technique

- Idea: perform a stream of queries with sensitivity 1; release only the identity of the first query which passes a test

- Above threshold algorithm

Inputs: $q_1, \dots, q_n, D, T, \epsilon$

- $\hat{T} = T + \text{Lap}(\frac{2}{\epsilon})$
- for i, q in queries:
- $v_i = \text{Lap}(\frac{u}{\epsilon})$
- if $q(0) + v_i \geq \hat{T}$: return i
- return None

Application: auto-Clip (helps to find right bound)

It is possible to return multiple values

(17) Local DP

idea: each part adds noise

$$x_i \stackrel{\text{before}}{\geq} \boxed{\sum} \rightarrow + \quad x_i + y_i \stackrel{\text{after}}{\rightarrow} \boxed{\sum} \rightarrow$$

Randomized response

- With $p = \frac{1}{2+\epsilon}$ report true, with $\frac{1}{2-\epsilon}$ report $y_i = 1-x_i$ (flipped value)
- $x_i = \frac{y_i - 1/\epsilon + 1}{2\epsilon}$, so amount $M(x) = \sum x_i$ - unbiased estimator
- Randomized response satisfies ϵ -DP if: $\epsilon = \frac{e^{\epsilon}-1}{e^{\epsilon}+1}$
- accuracy lower than for central DP

Unary encoding

- use one-hot encoding
- $\Pr[B'_i[i] = 1] = \begin{cases} p & \text{if } B[i] = 1 \\ q & \text{if } B[i] = 0 \end{cases}$
- this M satisfies ϵ -DP, $\epsilon = \ln(\frac{p(1-q)}{(1-p)q})$
- then aggregator for each position:
 $A[i] = \frac{\sum B'_i[i] - nq}{p-q}$

Example of this M : Google RAPPOR

- hash the client value v onto the Bloom filter B
- $B'_i = \begin{cases} 1 & \text{with } p = 1/\epsilon \\ 0 & \text{with } p = 1/\epsilon \\ B_i & \text{with } 1-\epsilon \end{cases}$
- $P(S_i = 1) = \begin{cases} q & B_i = 1 \\ p & B_i = 0 \end{cases}$
- Send S to server

Example: Chrome browser homepages

(18) Local DP at Apple: Count mean sketch

Idea: encode the data element (instead of one-hot encoding)

Collection data over time

- difference to RAPPOR: much larger ϵ , no consideration of longitudinal privacy.
- Attack against: Pool inference attack

④ Secure multi-party computation

① Settings

- Players: $P = \{P_1, P_2, \dots, P_n\}$ wants to compute $f(x_1, \dots, x_n) = (y_1, y_2, \dots, y_n)$
They can communicate
- Assumption: ideal world:
 - T - completely trusted, incorruptible party
 - each P_i sends x_i to T
 - T computes $f(x_1, \dots, x_n)$ and return result
- [in real world there is no T , but there is a protocol π , which is secure if emulates the ideal execution] so we simulate adversary
- $\Pr[A(\pi, E) = x] \cong \Pr[S(f, E) = x] \leftarrow$ in ideal world
- Types of adversaries:
 1. Semi-honest (follows π , but sniffing)
 2. Malicious (deviate from π)

② Yao's Garbled circuits (GC) protocol

1. Idea: lookup table

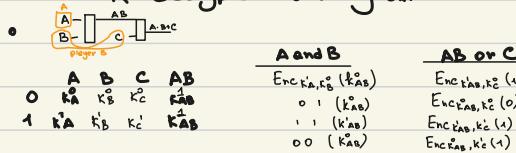
2. Oblivious transfer

• R generates S_k, P_k and take random key P_k'

• R sends $\begin{cases} (P_k, P_k') & \text{for } \delta=0 \\ (P_k', P_k) & \text{for } \delta=1 \end{cases}$

• S returns $(\text{Enc}_{P_k}(m_0), \text{Enc}_{P_k}(m_1))$.

• R decrypts $e\delta$ using S_k .



• A sends to B: 1. circuit and its key

• B via OT ask keys for its input

• B decrypt and sends output to A

• Security: Semi-honest model (OT secure in semi-honest)

• We lose security if we evaluate a circuit more than once with diff. outputs

③ BGW

- Idea: arithmetic functions with more than two inputs \Rightarrow arithmetic circuits
- we can represent any polynomial with a combination of addition and multiplication gates
- we operate on int.

- idea: - secret S , that we split into n shares

- only $t+1$ ^{called threshold} is enough to reconstruct s

• Shamir's secret sharing

I. - n players, threshold $t \in [1, n]$

- select finite field $\mathbb{Z}_q, q > n$

- secret is $s \in \mathbb{Z}_q$

- select n distinct elements $P_1, \dots, P_n \in \mathbb{Z}_q \setminus \{0\}$

- sample random polynomial p with degree t and $p(0) = s$

- compute the share of player $i: s_i = p(P_i)$

- reconstruct: Lagrange Polynomial Interpolation with any $t+1$, and evaluate $p(0) = s$

II. • with multiplication degree of polynomial $p \leq 2t$

↳ degree reduction step

$$p_c(x) = \sum_{i=1}^{2t+1} a_i c_i$$

• each player holds one of c_i

↳ • each player generates and distributes shares c_{ij} of their c_i

$$\bullet \text{each computes } C_j = \sum_{i=1}^{2t+1} a_i c_{ij}$$

• C_j is a share of c in a secret sharing with t

Example

• 5 and 2, \mathbb{Z}_{11}

• 3 parties

$\circ S_1(x) = 5 + 4x$

$\circ S_2(x) = 2 + 9x$

$\circ S_3(x) = 1 + 0 \cdot 9 \cdot x$

• new shares $0, 1, 2$

$\circ S_1(x) = 0 + 3x$ for $x \in \{0, 1, 2\}$

$\circ S_2(x) = 5 + 2x$ for $x \in \{0, 1, 2\}$

$\circ S_3(x) = 1 + 0 \cdot 9 \cdot x$ for $x \in \{0, 1, 2\}$

• BGW: for every wire in the arithmetic circuit the players holds shares $[Vw]$ of the value on that wire; semi-honest model; $2t+1 \leq n$; adversary can corrupt up to t parties.

④ Zero-knowledge proofs

- GMW compiler: run π and prove in zero-knowledge that every msg is the result of running π honestly
- can transform any SMC π from semi-honest to malicious model.

5 Privacy Metrics

① Privacy Metrics

1. Entropy (high \rightarrow more uncertainty)

$$\text{priv} = H(x) = - \sum_{\text{guess}} p(x) \log_2 p(x)$$

adversary estimated
ex. users

2. Expected estimation error

$$\text{PrvAE} = \sum_{\text{observations}} p(x|y) d(x, x^*)$$

3. Amount of Leaked information

$$\text{PrvAI} = I(S)$$

4. Relative entropy (number of bits of prob. info revealed to adversary)

$$\text{PrvRE} = D_{KL}(X^* || X)$$

5. Adv. success rate

② Utility metrics (how useful private is vs non-private)

1. General (Statistics, distance, confidence intervals,
(Wasserstein, Kullback, total-variation))

2. Task-specific (error, classification (Acc, Precision, Recall, F1, Auc))

③ Performance (Many steps)

- Time
- Communication
- Storage

⑥ Privacy for ML

① Background

- Logistic class more informative, err measures how far the prediction is from true label
- SGD: sample batch in each round to estimate the gradient
 - Sample mini-batch B from training ex.
 - For each $i \in B$ compute loss l_i
 - Compute average gradient $\bar{g}_t = \frac{1}{|B|} \sum_{i \in B} \nabla_{\theta_t} l_i$
 - Update weights $w_{t+1} = w - \eta g_t$

② Attacks on ML:

1. Membership inference

2. Shadow models (\rightarrow features or stat. properties; member or not)

3. Model Inversion (\rightarrow features or stat. properties; face ex.)

4. Property Inference (reconstruct properties from training data; SNP+gen.info/text from LLM)

5. Training data extraction

6. Model extraction

③ DP-SGD

- Add DP noise to computation of gradient

$$\text{Comp. avg. gradient } \bar{g}_t = \frac{1}{|B|} \sum_{i \in B} \nabla_{\theta_t} l_i + \mathcal{N}(0, \delta^2), \delta^2 = \frac{2S^2 \ln(\frac{1.25}{\epsilon})}{\epsilon^2}$$

- Clipping: scaling vector: element wise division by L₂ norm results in L₂ norm (sensitivity) of 1

- Clipping threshold C : median of the norms of the unclipped gradient

$$\text{• Privacy accounting: strong composition: } \delta^2 = \frac{2 \ln(\frac{1.25}{\epsilon})}{\epsilon^2}, \text{ each step is } (\eta q, \eta p)\text{-private, } q = \frac{|B|}{N},$$

for T iterations, the training is $(q\sqrt{T}\log(\frac{1}{\epsilon}), T\eta\delta)$ -private

2. moments accountant: better method to analyze privacy cost, tighter bound

$(q\sqrt{T}, \delta)$ -private; (same privacy, can train longer)

• Algorithm:

- Sample mini-batch B

- For each $i \in B$: l_i :

- For each $i \in B$: $g_t(x_i) = \nabla_{\theta_t} l_i$

$$- \text{Clip gradient: } \bar{g}_t(x_i) = \max\left(0, \frac{\min(g_t(x_i), C)}{C}\right)$$

$$- \text{Add noise, avg. gradients: } \tilde{g}_t = \frac{1}{|B|} \sum_{i \in B} (\bar{g}_t(x_i) + \mathcal{N}(0, \delta^2 C))$$

$$- \text{Update weights: } w_{t+1} = w - \eta \tilde{g}_t$$

- Output w_{t+1} , compute (ϵ, δ)

- Issues: fixed choice of optimizer, params to tune, batch size and selection influence the privacy guarantee.

(4) PATE (Private Aggreg. of teacher ensembles)

- Independent of ML algo and of optimizers, loss f, batch
- Main phases:
 1. teacher training (data in n disjoint sets, each used for train. $f_i(x)$ without privacy)
 2. aggregation (Label count $n_j(x)$ for class j is # of teacher models that predict class, label with highest count is agg. output; add Lap. noise: $f(x) = \text{argmax} \{f_{n_j}(x) + \text{Lap}(\frac{1}{f})\}$)
a single agg. is $(x, 0)$ -DP
 3. knowledge transfer (to get around limit number of user queries within, learns how to mimic agg. output of teacher ensemble; control privacy by limiting how many data records labeled.)

7 Federated learning

① Federated SGD:

- Clients $\{C_1, \dots, C_n\}$, input $\{(x_i^c, y_i^c)\dots\}$
- Server chooses loss L , initializes weights w , chooses η
- For iteration $t \in [0, \dots, n]$:
 - Server samples a random fraction of clients, C_r , sends them weights w
 - Each client $c \in C_r$: compute loss for all samples $\hat{L}_t^c = L(h_t; x_i^c, y_i^c)$
 - Compute gradient vector $\hat{g}_t^c = \nabla h_t \hat{L}_t^c$
 - Send \hat{g}_t^c to server
 - Server compute avg gradient $g_t = \frac{1}{|C_r|} \sum_c g_t^c$
 - Server update $w_{t+1} = w - \eta g_t$

② Attacks

- Adversary models: external, server, clients.

1. Honest-but-curious server: target membership inference (server records all model updates (gradients) from all C at all epochs, after training membership for each C)
2. Honest-but-curious server: gradient inversion (attacker steals training data from C ; updates dummy inputs and labels to min distance between gradients or add more params so client grad. memorizes training)
3. Malicious user or server: gradient ascent (is single datapoint part of dataset, report inverse gradient, so next iter. someone's SGD will abruptly reduce grad. of loss)

③ Privacy FL

1. FL with central DP (ex. Gboard)

- protect against honest-but-curious server

- algorithm (combination of DPSGD and federated SGD, users clip grad and server avg. grads and applies Gaussian noise)

2. FL with local DP

- each user adds DP to their grad.

- threat model: malicious server, users

- drawbacks: need very large # users and amount of noise

3. FL with Secure Aggregation

- malicious server, so use secure multi-party computation to agg. grad.

- drawbacks: no DP (membership attacks)

4. FL with distributed DP and secure aggregation

1. Hybrid privacy-preserving FL

- each iteration: one epoch of DPSGD (central DP, not local) on each client
- gradients shared with threshold homomorphic encryption
- assume t honest clients
- accuracy is high, but secure agg. allows to use less noise

2. Distributed discrete gaussian for FL with secure aggregation

- Secure agg. is a black box → can plug in any SMC scheme, including ones for malicious Server
- Client discretizes data to reduce communication overhead (fewer bits in discretization; reduced accuracy, reduced communication; Client apply discrete Gaussian noise, achieves $\frac{1}{2}\epsilon^2$ -concentrated DP)

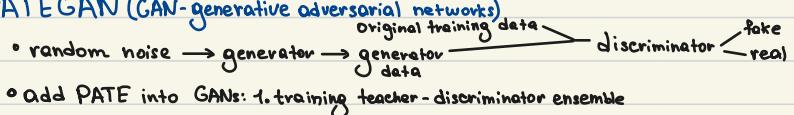
④ Synthetic data generation

- one column: distribution of data, histograms, #. of bins = # of unique values; add Lap. noise, parallel composition; sample points
- multi-column: n-way marginal distribution; ok for small data; for big: low accuracy (noise overwhelms the signal)

1. PrivBayes (mitigate low cell counts)

- represent correlations in Bayesian Network or Markov random fields (better acc)
- (derive low-dim tables; generate DP version of it, sample)
- DP Bayesian Network: input d attributes $A_1 \dots A_d, k > 0$
 - initialize empty Bayesian network N
 - consider all possible $(k+1)$ attribute combination $A_{i_1} \dots A_{i_k}, A_j$ and evaluate mutual info $I(A_{i_1} \dots A_{i_k}; A_j)$
 - Add DP noise to mutual info values and select combination with largest noisy mutual info.
 - Add selected $A_{i_1} \dots A_{i_k} \rightarrow A_j$ into N
 - Repeat with $A_j \& N$

2. PATEGAN (GAN-generative adversarial networks)



- random noise → generator → generator → discriminator → fake
original training data → teacher → student → real
- Add PATE into GANs: 1. training teacher-discriminator ensemble

2. training generator and student-discriminator (add noise on aggregation)

⑤ Implementation challenges

- both implementation determine value ranges for columns from raw data \rightarrow break DP
- Win approach of NISI [MS]:
 1. measure 1, 2, 3-way marginals using gaussian mechanism
 2. find synthetic dataset that has those marginals, using graphical model
 3. generate data samples and measure performance (lower score - better)
- future ?: other data types, image\time data ; with rare events.

8 Censorship

① Censorship phases:

1. fingerprinting (destinations [IP, ports, domains], content, flow properties [packet lengths, time], protocol semantics).
2. direct censorship (user-side [preinstalled software], publisher-side censorship, degrade performance, block destinations, corrupt routing info, corrupt flow content [modify html], corrupt protocol semantics)

② Censorship resistance systems (CRS)

- phases:
 1. **Communication establishment** (obtain credentials, CRS-specific info)
 - mechanisms - high churn access (freq. credentials change) - rate-limited access (proof of work)
 - active probing resistance (obfuscate aliveness or service)
 - true-based access (social graph)
 2. **Conversation**
 - Client connects to CRS (censor should not be able to block, fingerprint link)
 - CRS server connects to publisher
 - mechanism:
 1. **Focus on link user-CRS**:
 - mimicry (transform traffic so „looks“ as whitelisted, content mimicry → high traffic volumes for censor; imitate protocol, introduce randomness → censors typically permit unknown traffic; flow mimicry)
 - tunnelling - Covert channel (hide content), traffic manipulation
 - destination obfuscation
 2. **publisher-CRS link**:
 - content redundancy (store content on multiple servers)
 - distributed Content storage (chunks)

④ Scramblesuit/DBFS4*

Tor	VPN
	Socks
#	
TCP	
IP	

- pseudo-random payload: traffic indistinguishable from randomness → no useful fingerprints for DP)
- polymorphic: flow characteristic change
- Shared Secret: server only speaks protocol if client proves knowledge of secret

⑤ Domain fronting/MEEK

- domain fronting: using different domain names for diff. protocols in the stack
- not free • Google, Azure capable of this

⑥ Snowflake:

- idea: proxy traffic through tmp proxies over WebRTC
- volunteers run snowflake proxies (e.g. Tor bridge), broker manages and gives proxy addresses to clients
 - use domain fronting → high cost of blocking
 - Only communication with broker is domain-fronted → lower cost compared to meek

⑦ Conjure

- idea: instead of hiding Tor traffic, integrate in into network core
- ISP deploys Conjure station: Client connects to phantom proxy: unused address in ISP address range
- station injects packets as if there was true proxy available
- hard to block: censor need to block all websites served by this ISP
- PF_Ring: interface that allows Conjure station to tap ISP traffic