

# Lecture 10 Posterior Sampling (Part II) and Objective Perturbation

Yu-Xiang Wang



**COMPUTER SCIENCE**

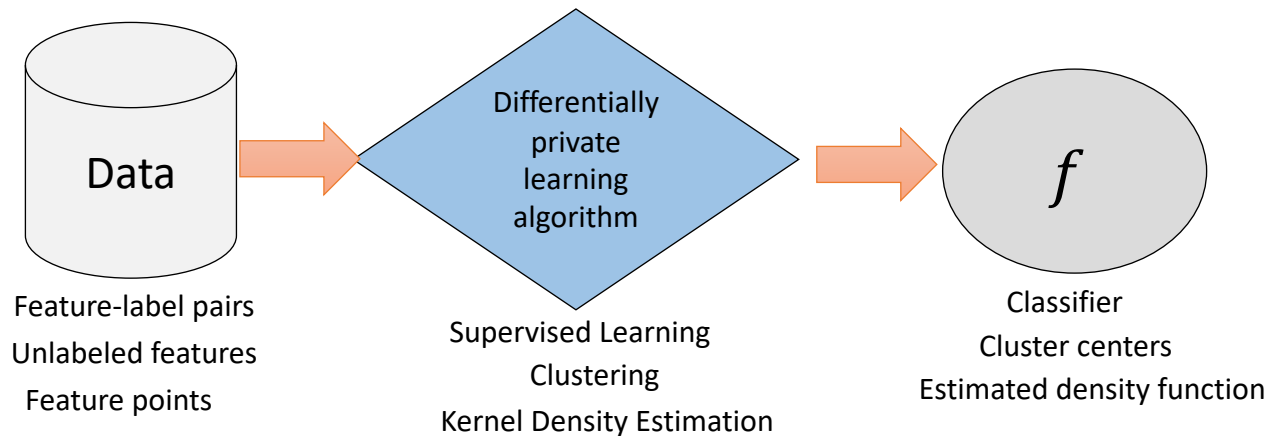
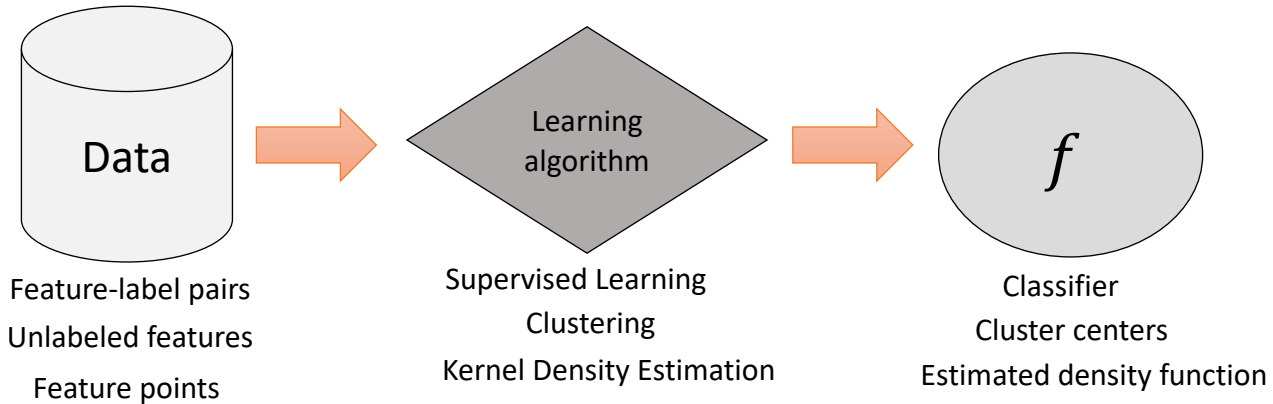
UC SANTA BARBARA

*Computing. ReInvented.*

# Recap: Last lecture

- Introduction to Differentially Private Machine Learning
  - Problem setup and notations
  - Examples
  - A learning theoretic study of the problem
- Posterior sampling mechanism
  - When the loss functions are bounded
  - A new analysis of the when they are not (not covered)

# Recap: Differentially Private Machine Learning



# Recap: Statistical Learning Jargons / Notations

- Data space / data points

$$(x_1, y_1), \dots, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y} = \mathcal{Z}$$

- Hypothesis class  $\mathcal{H}$
- Loss function, risk and empirical risk

$$\hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \ell(h, (x_i, y_i)) \quad R(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell(h, (x, y))]$$

- Goal of learning
- Realizable vs Agnostic

# Recap: Learning with differential privacy in theory

- Finite hypothesis class

	ERM	Exponential mechanism
Excess risk (Agnostic)	$\sqrt{\frac{\log( \mathcal{H} /\beta)}{n}}$	$\sqrt{\frac{\log( \mathcal{H} /\beta)}{n}} + \frac{\log( \mathcal{H} /\beta)}{n\epsilon}$
Excess risk (Realizable)	$\frac{\log( \mathcal{H} /\beta)}{n}$	$\frac{\log( \mathcal{H} /\beta)}{n} + \frac{\log( \mathcal{H} /\beta)}{n\epsilon}$

- Continuous hypothesis class (bounded VC-dim)
  - No DP algorithm gives consistent learning
  - A Packing lower bound
- However, there are weak assumptions we can add
  - Lipschitz loss functions
  - data-distributions with bounded density

# Recap: Packing lower bound in the problem of learning threshold function

- Assume  $K$ -data distributions with disjoint support
- Assume a DP algorithm  $A$  **works** on  $(K-1)$  of them
  - Expected error  $< 0.05 \Rightarrow$  need to stay within support w.p.  $0.9$
- Implies that it **does not work** on the remaining one using the definition of DP
  - Probability of outside support  $> 0.9 \Rightarrow$  Expected error  $> 0.45$

# Recap: Alternative analysis of the exponential mechanism.

- Bayesian learning:
  - Posterior sampling is private for free, if the log-likelihood function is bounded (or clipped).
  - Utility:  $(1 + 4B/\epsilon)$  multiplicative factor in efficiency.
- More advanced analysis by Minami et al:
  - For strongly log-concave prior

# This lecture

- Minami et al's analysis of posterior sampling
- Differentially Private Empirical Risk Minimization
  - Output perturbation
  - Objective perturbation



# Optional readings

- “Differentially Private Empirical Risk Minimization” by Chaudhuri, Monteleone, Sarwate.
  - <https://www.jmlr.org/papers/volume12/chaudhuri11a/chaudhuri11a.pdf>
- “Private Convex Empirical Risk Minimization and High-dimensional regression” by Kifer, Smith and Thakurta
  - <http://proceedings.mlr.press/v23/kifer12/kifer12.pdf>

# Improved analysis of exponential mechanism with strong convexity

- Assume  $\pi(\theta) = e^{-r(\theta)}$  where  $r(\theta)$  is  $\mu$ -strongly convex, i.e., the prior is strongly log-concave.

- Assume that the loss-function is Lipschitz

$$| -\log p(x|\theta) + \log p(x|\theta') | \leq L \|\theta - \theta'\|$$

- Then  $\hat{\theta} \sim P(\theta|x_1, \dots, x_n) \propto e^{(-\tau \sum_{i=1}^n \log p(x_i|\theta) - r(\theta))}$

obeys  $(\epsilon, \delta)$ -DP if  $\tau = O\left(\frac{\epsilon\sqrt{\mu}}{L\sqrt{\log(1/\delta)}}\right)$

# Idea of the proof for the improved analysis of EM

- The privacy loss random variable is
- Apply the tail bound lemma
- The strong log-concavity + Lipschitz assumption ensures that  $\hat{\theta}$  satisfies a “Log-Sobolev Inequality”
  - Which ensures a subgaussian-like tail bound for all Lipschitz functions of  $\hat{\theta}$
  - And a bound on the KL-divergence.

# The problem of convex empirical risk minimization

- Data
- Hypothesis class
- Loss function

# A specific class of problems of interest: generalized Linear models

- Loss function is of a particular form
- Examples:
  - Linear regression
  - Logistic regression
  - Support vector machine

# (Detour) Convex Optimization 101

- Convex functions
- Strongly convex functions
- Optimality condition for differentiable convex functions

# (Detour) Convex Optimization 101

- Lipschitz constant of a function
- Smoothness constant of a function

# Can we apply Laplace / Gaussian mechanism?

- What is the global sensitivity for the minimizer?



# Let's regularize and make it strongly convex

- Regularized objective function
  
- Now what is the global sensitivity again?

# Output perturbation mechanism

1. Solve for the optimal solution
2. Add noise to the optimal solution

# Utility analysis of the output perturbation mechanism

- Let's appeal to the smoothness of the loss function

# Objective perturbation

---

**Algorithm 1** Release  $\hat{\theta}^P$  via Obj-Pert (Kifer et al., 2012)

---

**Input:** Dataset  $D$ , noise parameter  $\sigma$ , regularization parameter  $\lambda$ , loss function  $L(\theta; D) = \sum_i \ell(\theta; z_i)$ , convex and twice-differentiable regularizer  $r(\theta)$ , convex set  $\Theta$ .

**Output:**  $\hat{\theta}^P$ , the minimizer of the perturbed objective.

Draw noise vector  $b \sim \mathcal{N}(0, \sigma^2 I)$ .

Compute  $\hat{\theta}^P$  according to (1).

---

$$\hat{\theta}^P = \operatorname{argmin}_{\theta \in \Theta} L(\theta; D) + r(\theta) + \frac{\lambda}{2} \|\theta\|_2^2 + b^T \theta, \quad (1)$$

- Remark:
  - Equivalent to adding a random synthetic data point with linear loss
  - One could also get pure DP if we choose  $b$  to be:

KKT condition of Objective  
Perturbation implies a relationship  
between  $b$  and the optimal solution.

$$\hat{\theta}^P = \operatorname{argmin}_{\theta \in \Theta} L(\theta; D) + r(\theta) + \frac{\lambda}{2} \|\theta\|_2^2 + b^T \theta,$$

# Change of variable trick: distribution via the Jacobian

**Theorem:** Let  $X, Y$  be random variables in  $\mathbb{R}^d$  and  $Y = g(X)$  satisfying that function  $g$  is bijective (one-to-one map) and differentiable. Then the probability density function of  $Y$  is:

$$f_Y(y) = f_X(g^{-1}(y)) \left| \det \left[ \frac{\partial g^{-1}(y)}{\partial y} \right] \right|$$

- Apply this theorem to our problem:

# The privacy loss random variable of the objective perturbation

$$\begin{aligned}
 \log \frac{\Pr(\mathcal{A}(D) = \hat{\theta}^P)}{\Pr(\mathcal{A}(D_{\pm z}) = \hat{\theta}^P)} &= \log \frac{|\det(\nabla b(\hat{\theta}^P; D))|}{|\det(\nabla b(\hat{\theta}^P; D_{\pm z}))|} \frac{\nu(b(\hat{\theta}^P; D); \sigma)}{\nu(b(\hat{\theta}^P; D_{\pm z}); \sigma)} \\
 &= \log \frac{|\det(\nabla b(\hat{\theta}^P; D))|}{|\det(\nabla b(\hat{\theta}^P; D_{\pm z}))|} + \log \frac{e^{-\frac{1}{2\sigma^2} \|b(\hat{\theta}^P; D)\|_2^2}}{e^{-\frac{1}{2\sigma^2} \|b(\hat{\theta}^P; D_{\pm z})\|_2^2}} \\
 &= \underbrace{\log \frac{|\det(\nabla b(\hat{\theta}^P; D))|}{|\det(\nabla b(\hat{\theta}^P; D_{\pm z}))|}}_{(*)} + \underbrace{\frac{1}{2\sigma^2} \left( \|b(\hat{\theta}^P; D_{\pm z})\|_2^2 - \|b(\hat{\theta}^P; D)\|_2^2 \right)}_{(**)}.
 \end{aligned}$$

# Bounding (\*): the difference of the log determinants

**Lemma 25** (Determinant of Rank-1 perturbation). *For invertible matrix  $A$  and vector  $c, d$  of compatible dimension*

$$\det(A + cd^T) = \det(A)(1 + d^T A^{-1}c).$$

- Iterative application of Lemma 25

$$\begin{aligned} \left| \det(\nabla b(\hat{\theta}^P; D_{\pm z})) \right| &= \left| \det\left(\nabla b(\hat{\theta}^P; D) \mp \nabla^2 \ell(\hat{\theta}^P; z)\right) \right| \\ &= \left| \det\left(\nabla b(\hat{\theta}^P; D) \mp \sum_{k=1}^d \lambda_k u_k u_k^T\right) \right| \end{aligned}$$

Recall:

$$b(\hat{\theta}^P; D) = -\left(\nabla \hat{\mathcal{L}}(\hat{\theta}^P; D) + \nabla r(\hat{\theta}^P) + \lambda \hat{\theta}^P\right).$$

$$\nabla b(\hat{\theta}^P; D) = -\left(\nabla^2 \hat{\mathcal{L}}(\hat{\theta}^P; D) + \nabla^2 r(\hat{\theta}^P) + \lambda I_d\right).$$

- For generalized linear models



Recall:

$$b(\hat{\theta}^P; D) = - \left( \nabla \hat{\mathcal{L}}(\hat{\theta}^P; D) + \nabla r(\hat{\theta}^P) + \lambda \hat{\theta}^P \right).$$

$$\nabla b(\hat{\theta}^P; D) = - \left( \nabla^2 \hat{\mathcal{L}}(\hat{\theta}^P; D) + \nabla^2 r(\hat{\theta}^P) + \lambda I_d \right).$$

# Bounding (\*\*)

$$\begin{aligned} (**) &= \frac{1}{2\sigma^2} \left( \|b(\hat{\theta}^P; D_{\pm z})\|_2^2 - \|b(\hat{\theta}^P; D)\|_2^2 \right) \\ &= \frac{1}{2\sigma^2} \left[ \mp \nabla \ell(\hat{\theta}^P; z) \right] \left[ 2b(\hat{\theta}^P; D) \mp \nabla \ell(\hat{\theta}^P; z) \right] \end{aligned}$$

- For generalized linear models:

# Putting everything together

- For general loss functions

# Next lecture

- Utility analysis of ObjPert
- Noisy Gradient Descent