

Lecture 10 Posterior Sampling (Part II) and Objective Perturbation

Yu-Xiang Wang



COMPUTER SCIENCE

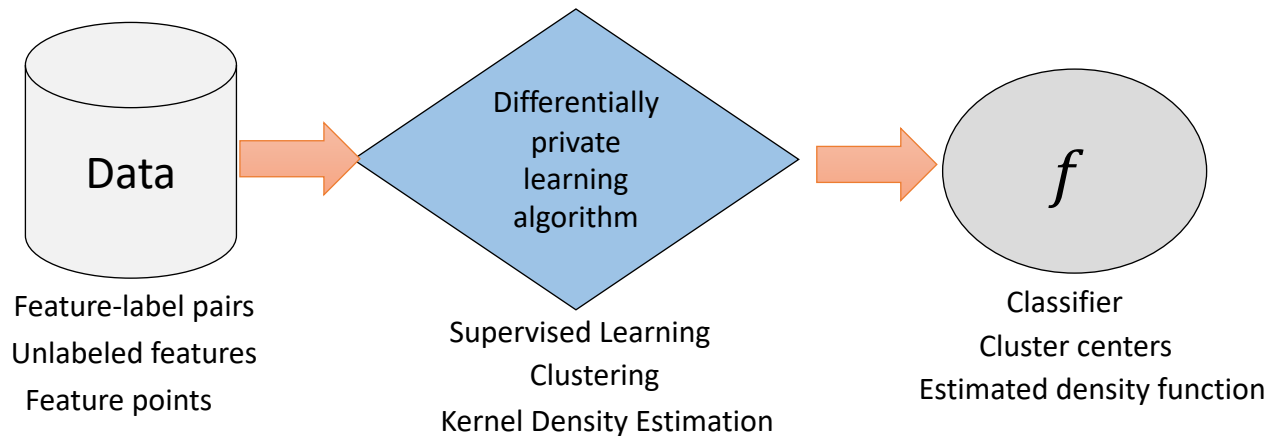
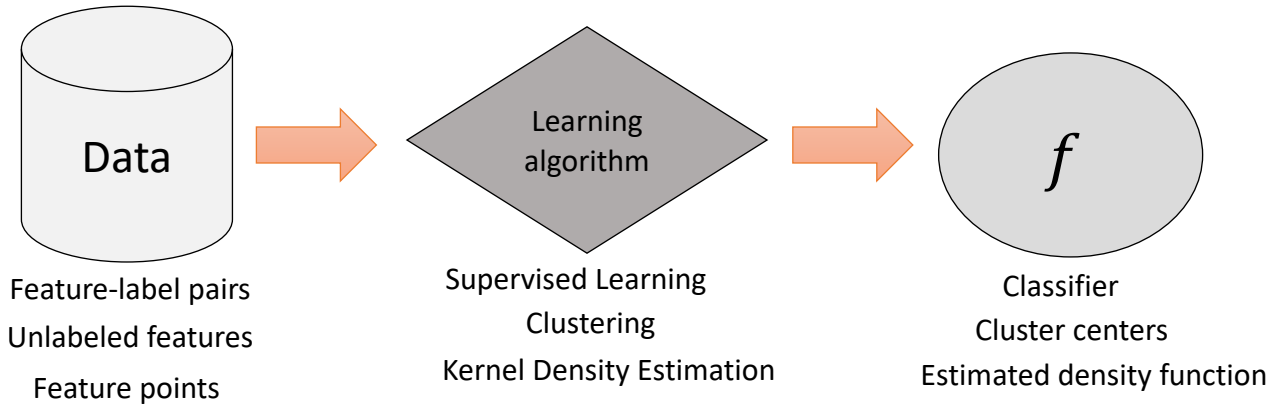
UC SANTA BARBARA

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

$$\hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \ell(h, (x_i, y_i))$$

$$R(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell(h, (x, y))]$$

- Goal of learning

$$h^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} R(h)$$

Excess Risk

$$R(\hat{h}) - R(h^*) \leq \epsilon$$

- Realizable vs Agnostic

$$R(h^*) = 0$$

$$R(h^*) > 0$$

Recap: Learning with differential privacy in theory

$|\mathcal{H}| < +\infty$

- Finite hypothesis class

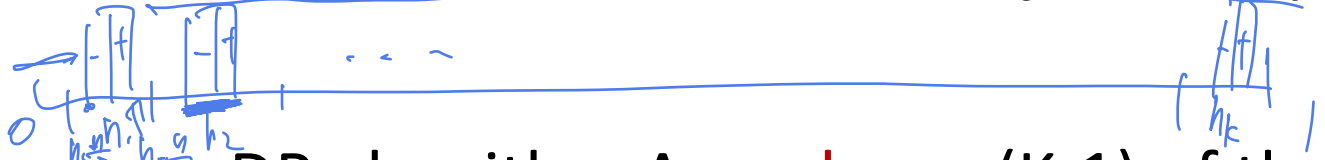
	ERM	Exponential mechanism
Excess risk (Agnostic)	$\sqrt{\frac{\log(\mathcal{H} /\beta)}{n}}$	$\sqrt{\frac{\log(\mathcal{H} /\beta)}{n}} + \frac{\log(\mathcal{H} /\beta)}{n\epsilon}$
Excess risk (Realizable)	$\frac{\log(\mathcal{H} /\beta)}{n}$	$\frac{\log(\mathcal{H} /\beta)}{n} + \frac{\log(\mathcal{H}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

- 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

$$\chi = [0, 1]$$

- Assume K-data distributions with disjoint support



- Assume a DP algorithm A **works** on (K-1) of them
 - Expected error < 0.05 => 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 => Expected error > 0.45

$$\begin{aligned}
 P(A(D_i) \notin [h_i - \frac{\eta}{3}, h_i + \frac{\eta}{3}]) &\geq \sum_{i=2}^K P(A(D_i) \in [h_i - \frac{\eta}{3}, h_i + \frac{\eta}{3}]) \\
 &\geq \sum_{i=2}^K e^{-n\epsilon} P(A(D_i) \in [h_i - \frac{\eta}{3}, h_i + \frac{\eta}{3}]) \\
 &\geq \sum_{i=2}^K e^{-n\epsilon} \cdot 0.9 = K e^{-n\epsilon} \cdot 0.9
 \end{aligned}$$

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 μ -strongly convex, i.e., the prior is strongly log-concave.

- Assume that the loss-function is Lipschitz

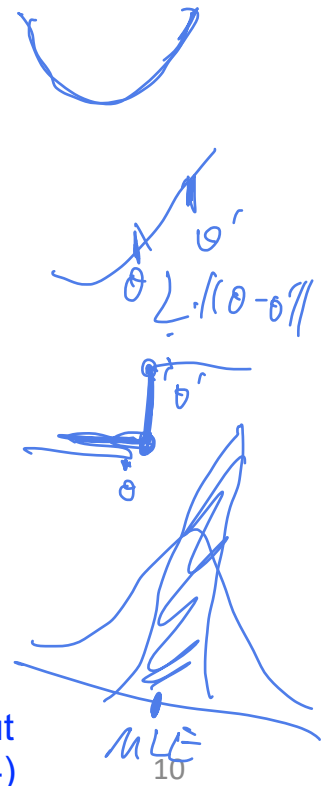
$$| -\log p(x|\theta) + \log p(x|\theta') | \leq L \|\theta - \theta'\|$$

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

$\sum_{i=1}^n \log p(x_i|\theta) \leq O(n)$

$\mu = \sqrt{n}$



Idea of the proof for the improved analysis of EM

- The privacy loss random variable is

$$P(\theta|D) \approx \frac{e^{-\tau \sum \log p(x_i|\theta) - \tau(\theta)}}{\sum_{\theta'} e^{-\tau \sum \log p(x_i|\theta') - \tau(\theta')}}$$

$$\log \frac{P(\theta|D)}{P(\theta|D')} = -\tau \sum_i (\log p(x_i|\theta) - \log p(x_i|\theta')) + \tau (\theta - \theta')$$

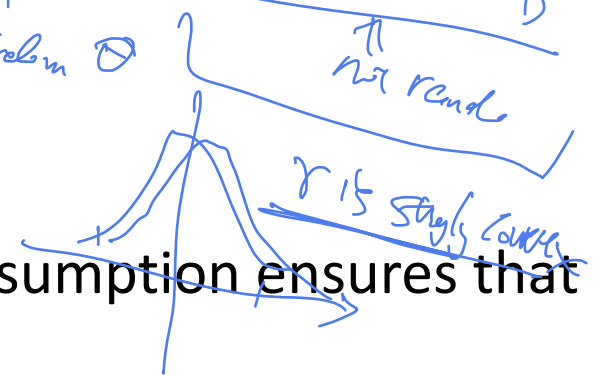
$$= \underbrace{\log p(x_i|\theta) - \log p(x_i|\theta')}_{\text{Random } \theta} - \log Z_D + \log Z_{D'}$$

$$\int_0^1 e^{-\tau \sum (\log p(x_i|\theta) - \log p(x_i|\theta'))} d\theta$$

- Apply the tail bound lemma

$$E \left(\log \frac{P(\theta|D)}{P(\theta|D')} \right) = D_{KL}(P(\theta|D) || P(\theta|D'))$$

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

- Hypothesis class

$$\hat{y} = \theta(x, \theta) = x^T \theta$$

- Loss function

$$l(\theta, (x, y)) \quad l \text{ needs to be convex in } \theta$$

$$\underbrace{f(\theta(x^T \theta, y))}_{\text{GLM}} \quad f \text{ is convex} \Rightarrow l \text{ is convex in } \theta$$

GLM

A specific class of problems of interest: generalized Linear models

- Loss function is of a particular form

$$l(\theta, (x, y)) = f(\theta^T x, y)$$

- Examples:

- Linear regression

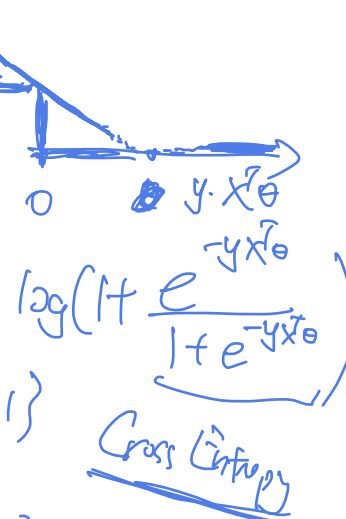
$$(y - x^T \theta)^2$$

- Logistic regression

$$\gamma = \{-1, 1\}$$

- Support vector machine

$$\gamma = \{0, 1\}$$



$$l(\theta, (x, y)) = \max \left\{ 0, 1 - \theta^T x \cdot y \right\}$$

(-1) large loss



(Detour) Convex Optimization 101

- Convex functions

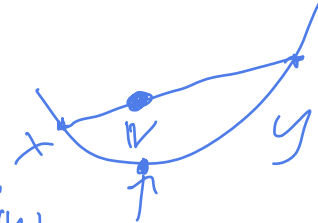
f is convex iff $\forall x, y, 0 \leq \alpha \leq 1$

$$f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$$

For f exists first order conditions

f is convex iff, $\forall x, y$

$$f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle$$



- Strongly convex functions

$$f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle + \frac{1}{2} \|y-x\|^2$$

f is μ -strongly convex

- Optimality condition for differentiable convex functions

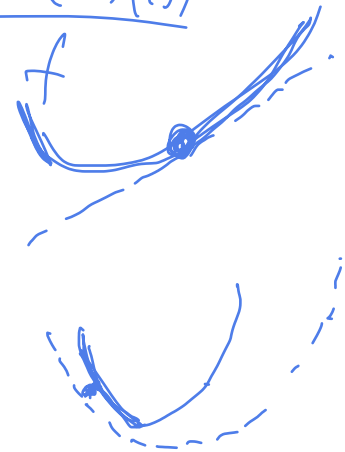
$\min_{x \in \mathcal{C}} f(x)$

$x^* \in \text{argmin } f(x)$

global

iff $\nabla f(x^*) = 0$

local information



(Detour) Convex Optimization 101

- Lipschitz constant of a function

$$f \text{ is } L\text{-Lipschitz if } \forall x, y \quad f(x) - f(y) \leq L \|x - y\|_2$$

- Smoothness constant of a function

f is β -Smooth
if $\forall x, y,$

$$f(y) \leq f(x) + \langle y - x, \nabla f(x) \rangle + \frac{\beta}{2} \|y - x\|^2$$

λ -Strong convex

$$f(y) \geq f(x) + \langle y - x, \nabla f(x) \rangle + \frac{\lambda}{2} \|y - x\|^2$$

Let's regularize and make it strongly convex

- Regularized objective function

$$\theta_{\lambda}^* = \underset{\theta}{\operatorname{argmin}} \underbrace{\sum_{i=1}^n \ell_i(\theta)}_{L(\theta)} + \frac{\lambda}{2} \|\theta\|^2$$

$$\underline{\theta}_{\lambda}^* = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^n \ell_i(\theta) + \ell(\theta) + \frac{\lambda}{2} \|\theta\|^2$$

- Now what is the global sensitivity again?

By Strong Convexity

$$\cancel{L(\theta_{\lambda}^*)} + \frac{\lambda}{2} \|\theta_{\lambda}^*\|^2 \geq \cancel{L(\theta_x^*)} + \frac{\lambda}{2} \|\theta_x^*\|^2 + \frac{\lambda}{2} \|\theta_x^* - \theta_{\lambda}^*\|^2$$

$$\cancel{L(\theta_{\lambda}^*)} + \frac{\lambda}{2} \|\theta_{\lambda}^*\|^2 \geq \cancel{L(\theta_{\lambda}^*)} + \frac{\lambda}{2} \|\theta_{\lambda}^*\|^2 + \frac{\lambda}{2} \|\theta_x^* - \theta_{\lambda}^*\|^2$$

$$+ \ell(\theta_x^*) \quad \quad \quad \underline{+ \ell(\theta_{\lambda}^*)}$$

add up

$$\lambda \|\theta_x^* - \theta_{\lambda}^*\|^2 \leq \ell(\theta_x^*) - \ell(\theta_{\lambda}^*) \leq L \|\theta_x^* - \theta_{\lambda}^*\|$$

$$\|\theta_x^* - \theta_{\lambda}^*\| \leq \frac{L}{\lambda}$$

Output perturbation mechanism

1. Solve for the optimal solution

$$\theta_{\lambda}^* = \underset{\theta}{\operatorname{argmin}} L(\theta) + \frac{\lambda}{2} \|\theta\|^2$$

2. Add noise to the optimal solution

Release $\theta^{\text{P}} = \theta_{\lambda}^* + \mathcal{N}(0, G^2)$

$$G = \frac{L}{\lambda \epsilon} \sqrt{2 \log \frac{2.5}{\delta}}$$

(ES)-DP

Utility analysis of the output perturbation mechanism

$\sum_{i=1}^n \ell_i$ is $\beta \Sigma_{i=1}^n \ell_i$
 $\sum_{i=1}^n \ell_i$ is $n\beta$ -smooth

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

$$(*) \quad \underline{L(\hat{\theta}^P) - L(\theta^*)} = \underline{L(\hat{\theta}^P) - L(\theta_\lambda^*)} + \underline{L(\theta_\lambda^*) - L(\theta^*)}$$

$$\hat{\theta}^P = \theta_\lambda^* + N(0, \sigma^2 I_d)$$

$$L(\hat{\theta}^P) + \lambda \|\hat{\theta}^P\|^2 \leq L(\theta_\lambda^*) + \frac{\lambda}{2} \|\theta_\lambda^*\|^2 + \underbrace{\frac{n\beta}{2} \|\hat{\theta}^P - \theta_\lambda^*\|^2}_{\text{noise}}$$

$$L(\hat{\theta}^P) - L(\theta_\lambda^*) \leq \underbrace{\frac{\lambda}{2} \|\theta_\lambda^*\|^2 - \frac{\lambda}{2} \|\hat{\theta}^P\|^2}_{\leq 0} + \underbrace{\frac{n\beta \sigma \sqrt{d}}{2} \sqrt{2 \log \frac{1}{\delta}}}_{\geq}$$

$$L(\theta_\lambda^*) \leq L(\theta^*)$$

$$+ \underbrace{\frac{1}{2} \|\theta_\lambda^*\|^2}_{\geq}$$

$$+ \frac{n\beta}{2} \|\theta_\lambda^*\|^2$$

$$L(\theta_\lambda^*) + \frac{1}{2} \|\theta_\lambda^*\|^2$$

$$\leq L(\theta^*) + \frac{1}{2} \|\theta_\lambda^*\|^2$$

$$L(\hat{\theta}^P) - L(\theta^*) \leq \frac{1}{2} \|\theta_\lambda^*\|^2 + \frac{n\beta \sigma \sqrt{d}}{2} \sqrt{2 \log \frac{1}{\delta}}$$

(*) \leq

$$+ \frac{n\beta}{2} \|\theta_\lambda^*\|^2 \geq \frac{n\beta}{2} \frac{L(\theta_\lambda^*)}{\lambda} \geq \frac{n\beta}{2} \frac{L(\theta^*)}{\lambda}$$

Choose λ optimally

$$\leq O\left(\frac{n\beta L(\theta^*) \sqrt{d}}{\epsilon}\right)$$

$$O\left(\frac{1}{\epsilon}\right)$$

Objective perturbation

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

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

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).

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

- Remark:

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

$$P(b) \propto e^{-\frac{\|b\|_2^2}{\sigma^2}}$$

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

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

$$\nabla L(\theta; D) + \nabla r(\theta) + \lambda \theta + b = 0$$

$$\underline{b}(\theta; D) = -\nabla L(\theta; D) - \nabla r(\theta) - \lambda \theta$$

$$\underline{b}(\theta; D) = -\nabla L(\theta; D) - \underline{\ell}(\theta) - \nabla r(\theta) - \lambda \theta$$

$$\underline{b} \sim \mathcal{N}(0, \sigma^2 I)$$

$$\log \frac{P(\theta|D)}{P(\theta|D')}$$

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:

$$b = \theta^{-1} \begin{pmatrix} 0 \\ 0 \end{pmatrix} = b(\theta; D)$$

$$P(A(D) = \hat{\theta}^0) = \frac{1}{(2\pi\sigma)^d} e^{-\frac{\|b(\hat{\theta}^0; D)\|^2}{2\sigma^2}} \left| \det(-\nabla^2 \ell(\hat{\theta}^0; D) - \nabla^2 r - \lambda I_d) \right|$$

$b \sim N(b, \sigma^2 I_d)$

The privacy loss random variable of the objective perturbation

\mathcal{V} :

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

$$\begin{aligned} &\log\left(1 + \frac{\beta}{\lambda}\right)^d \\ &= d \log\left(1 + \frac{\beta}{\lambda}\right) \\ &\approx \frac{d\beta}{\lambda} \end{aligned}$$

$$\frac{L^2}{2\sigma^2} + O\left(\frac{L^2 \beta}{\sigma^2}\right)$$

(ϵ, δ) -DP

Bounding (*): the difference of the log determinants

$\lambda_i \leq \beta$ ℓ is β -smooth

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

Recall:

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

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

$$\left| \det(\nabla b(\hat{\theta}^P; D_{\pm z})) \right| = \left| \det(\nabla b(\hat{\theta}^P; D) \mp \nabla^2 \ell(\hat{\theta}^P; z)) \right|$$

$$\leftarrow \left(\frac{d^T \nabla b(\hat{\theta}^P; D)}{\lambda} \right) \cdot \left(1 + \frac{\beta}{\lambda} \right)^d = \left| \det(\nabla b(\hat{\theta}^P; D) \mp \sum_{k=1}^d \lambda_k u_k u_k^T) \right|$$

$$= \left| \det(\nabla b(\hat{\theta}^P; D) \mp \sum_{k=1}^d \lambda_k u_k u_k^T) \right|$$

$$= \left| \det(\nabla b(\hat{\theta}^P; D) \mp \sum_{k=1}^d \lambda_k u_k u_k^T) \right| \cdot \left(1 + \frac{u_k^T A^{-1} u_k}{\lambda} \right) \leq \left| \det(\nabla b(\hat{\theta}^P; D)) \right| \cdot \left(1 + \frac{\beta}{\lambda} \right)^d$$

- For generalized linear models

GLM $\nabla_{\theta} \ell(\theta^T x, y) = \underbrace{f'(\theta^T x, y)}_{\in \mathbb{R}} \cdot \underbrace{x}_{\in \mathbb{R}^d}$

$\leq \log\left(1 + \frac{\beta}{\lambda}\right)$

$$\nabla^2 \ell(\theta; z) = U \Lambda U^T$$

$\left| \det(\nabla b(\hat{\theta}^P; D)) \right| \cdot \left(1 + \frac{\beta}{\lambda} \right)^d$

$x x^T$ -scale

Recall:

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

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

Bounding (**)

$$x^2 - y^2 = (x+y)(x-y)$$

$$(**) = \frac{1}{2\sigma^2} (\|b(\hat{\theta}^P; D_{\pm z})\|_2^2 - \|b(\hat{\theta}^P; D)\|_2^2)$$

$$= \frac{1}{2\sigma^2} [\mp \nabla \ell(\hat{\theta}^P; z)] [2b(\hat{\theta}^P; D) \mp \nabla \ell(\hat{\theta}^P; z)]$$

$$b \sim \mathcal{N}(0, \sigma^2 L_d)$$

$$= \frac{1}{2\sigma^2} \|\nabla \ell(\hat{\theta}^P; z)\|_2^2 + \frac{1}{2\sigma^2} \nabla \ell(\hat{\theta}^P; z)^T 2b(\hat{\theta}^P; D)$$

$$\leq L^2$$

$$\leq 2 \|\nabla \ell(\hat{\theta}^P; z)\|_2 \|b\|_2$$

$$\leq 2L \sigma \sqrt{2 \log \frac{1}{\delta}}$$

$$\|b\|_2^2 \sim \chi^2(d)$$

- For generalized linear models:

$$b^T x \cdot \text{scalar}$$

Second term

$$\leq \sigma \cdot 2 \sqrt{2 \log \frac{1}{\delta}}$$

$$\sim \mathcal{N}(0, \sigma^2 \|x\|^2)$$

Putting everything together

- For general loss functions

Next lecture

- Utility analysis of ObjPert
- Noisy Gradient Descent