

Lecture 11 Noisy Gradient Descent

Yu-Xiang Wang



COMPUTER SCIENCE

UC SANTA BARBARA

Computing. ReInvented.

Recap: Last lecture

- Convex empirical risk minimization
- Output perturbation
- Objective perturbation

Recap: Convex ERM and optimality conditions

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

- Convex ERM:

$$\min_{\theta \in \Theta \subset \mathbb{R}^d} \sum_{i=1}^n \ell(\theta, (x_i, y_i)) \quad \text{=: } \mathcal{L}(\theta)$$

- Optimality condition: gradient = 0

$$\nabla \mathcal{L}(\theta^*) = 0$$

- Assumptions: Lipschitzness, Smoothness

ℓ is L -Lipschitz, β -smooth

Logistic Regress 1-Lip $\frac{1}{4}$ -Smooth if $\|x\|_2 \leq 1$

Recap: Output perturbation

$$\theta^* = \underset{\theta}{\operatorname{argmin}} L(\theta) \quad \sum$$

- Stability of the output via regularization

$$\|\theta_{\lambda}^* - \theta_a^*\|_2 \leq \frac{\|\nabla L(\theta^*)\|}{\lambda} \leq \frac{L}{\lambda}$$

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

- Privacy: from Gaussian mechanism

$$G = \frac{L}{\lambda \epsilon} \sqrt{2 \ln \frac{1.45}{\epsilon}} \quad (\epsilon < 1)$$

- Utility:

- Last time: under smoothness (has a small error ☹)
- Let's do it again.

$$L = \sum \ell_i \quad (n\beta + 1)\text{-Smooth}$$

Recap: Utility of Output perturbation

β

• Smooth losses

$$\begin{aligned}
 L(\hat{\theta}^P) + \frac{\lambda}{2} \|\hat{\theta}^P\|^2 &\leq L(\theta_\lambda^*) + \frac{\lambda}{2} \|\theta_\lambda^*\|^2 + \frac{n\beta + 1}{2} \|\theta_\lambda^* - \hat{\theta}^P\|^2 \stackrel{\leq \frac{\beta}{2} \sigma^2 \cdot d}{\leq} \\
 &\stackrel{\geq 0}{\geq} \leq L(\theta^*) + \frac{\lambda}{2} \|\theta^*\|^2 + \frac{(n\beta + 1)}{2} \frac{\beta \sigma^2 d}{\lambda} \log(\cdot) \quad \uparrow \text{noise} \\
 &= L(\theta^*) + \frac{\lambda}{2} \|\theta^*\|^2 + \frac{(n\beta + 1)}{2} \frac{\beta \sigma^2 d}{\lambda} \log(\cdot) \\
 &\stackrel{\text{choose } \lambda \text{ optimally}}{=} L(\theta^*) + \frac{(n\beta + 1)}{2} \frac{\beta \sigma^2 d}{\lambda} \log(\cdot) \geq \frac{\lambda^2 \epsilon^2}{8} \log(\cdot) \\
 &= L(\theta^*) + \frac{n^{\frac{1}{3}} d^{\frac{1}{3}} \beta^{\frac{1}{3}} L^{\frac{2}{3}}}{\epsilon^{\frac{2}{3}}} \|\theta^*\|^{\frac{2}{3}} (\log(\cdot))^{\frac{1}{3}} \quad (\log(\cdot))
 \end{aligned}$$

• Lipschitz losses

$$\begin{aligned}
 L(\hat{\theta}^P) &\leq L(\theta_\lambda^*) + nL \|\theta_\lambda^* - \hat{\theta}^P\|_2 \quad \uparrow \text{noise} \\
 &\leq L(\theta^*) + \frac{\lambda}{2} \|\theta^*\|^2 + \frac{nL \sqrt{d} L \sqrt{\frac{1}{8} \frac{\beta \sigma^2 d}{\lambda}} \sqrt{\log \frac{d}{\beta}}}{\epsilon} \\
 &= L(\theta^*) + o\left(\frac{n^{\frac{1}{3}} d^{\frac{1}{3}} L (\log(\cdot))^{\frac{1}{4}}}{\epsilon^{\frac{1}{4}}}\right)
 \end{aligned}$$

if $n\beta < 1$ $\Rightarrow \frac{1}{2} \|\theta^*\|^2 + \frac{\beta \sigma^2 d}{\lambda} \log(\cdot)$

if $n\beta \gg 1$ $\Rightarrow \frac{\beta \sigma^2 d}{\lambda} \log(\cdot)$

$\frac{1}{\sqrt{n}} \sqrt{d} L \log(\cdot)$

Recap: Objective perturbation

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

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

$\nabla L(\hat{\theta}^P) + \lambda \hat{\theta}^P + b = 0$

- Privacy analysis $\frac{\log(1+\beta)}{\lambda}$

- For GLM

$P(\text{PLRV} \geq \frac{2\beta}{\lambda} + \frac{L^2}{2G^2} + \frac{L \sqrt{2 \log \frac{2d}{\delta}}}{2G}) \leq \delta$

(Handwritten notes: $\lambda \geq \beta$, $\frac{\epsilon}{2}$, $\frac{\epsilon}{2}$)

if $\lambda \geq \frac{\epsilon}{4\beta}$

We can choose $G = \frac{L}{\frac{\epsilon}{2}} \sqrt{2 \log \frac{2d}{\delta}}$

- For General smooth learning problems

$P(\text{PLRV} \geq \frac{2d\beta}{\lambda} + \frac{L^2}{2G^2} + \frac{L \sqrt{2 \log \frac{2d}{\delta}}}{2G}) \leq \delta$

$\lambda = \frac{d\epsilon}{4\beta}$ $G = \frac{dL}{\frac{\epsilon}{2}} \sqrt{2 \log \frac{2d}{\delta}}$

This lecture

- Utility analysis of objective perturbation
- Noisy Gradient Descent
- Privacy amplification by sampling and NoisySGD

Readings

- Chaudhuri et al. / Kifer et al. (continuing)
- Bassily et al. (2014) Private empirical risk minimization: Efficient algorithms and tight error bounds. In *FOCS*. <https://arxiv.org/abs/1405.7085>
 - For the NoisySGD algorithm
 - For NoisyGD just refer to this lecture note.

Utility analysis of objective perturbation

$$J(\theta) := L(\theta) + \frac{\lambda}{2} \|\theta\|^2$$

$$\hat{\theta}^P = \underset{\theta}{\operatorname{argmin}} J(\theta) + \underline{b^T \theta}$$

$$J(\hat{\theta}^P) + \underline{b^T \hat{\theta}^P} \leq J(\theta_\lambda^*) + \underline{b^T \theta_\lambda^*} + b^T (\theta_\lambda^* - \hat{\theta}^P)$$

$$\leq J(\theta_\lambda^*) + b^T (\theta_\lambda^* - \hat{\theta}^P) \leftarrow$$

$$= L(\theta_\lambda^*) + \frac{\lambda}{2} \|\theta_\lambda^*\|^2 + \|b\|_2 \|\theta_\lambda^* - \hat{\theta}^P\|_2 \leq L(\theta_\lambda^*) + \frac{\lambda}{2} \|\theta_\lambda^*\|^2 + \frac{\|b\|_2^2}{\lambda}$$

$$\theta_\lambda^* = \underset{\theta}{\operatorname{argmin}} J(\theta)$$

$$\theta^* = \underset{\theta}{\operatorname{argmin}} L(\theta)$$

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

Recall J is λ -strongly convex

$$\left(\|\theta_\lambda^* - \hat{\theta}^P\|_2 \leq \frac{\operatorname{Lip}(b^T \theta)}{\lambda} = \frac{\|b\|_2}{\lambda} \right)$$

For GLM: $\mathcal{O}\left(\frac{\sqrt{d}}{\varepsilon} \sqrt{\log \frac{1}{\varepsilon}} \|\theta^*\| \sqrt{\log \frac{d}{P}}\right)$

for general smth $\mathcal{O}\left(\frac{dL}{\varepsilon} \sqrt{\log \frac{1}{\varepsilon}} \|\theta^*\|\right)$

$$\leq L(\theta^*) + \frac{\lambda}{2} \|\theta^*\|^2 + \frac{\sigma^2}{\lambda} \sqrt{\log \frac{d}{P}} \stackrel{\text{chosen } \lambda}{=} L(\theta^*) + \sigma \sqrt{\log \frac{d}{P}} \frac{1}{\varepsilon} \|\theta^*\|$$

Checkpoint: Compare the **excess empirical risk** of Output/Objective Perturbation

	Lipschitz losses	Smooth losses	Smooth / Lipschitz GLM
Output Pert	$\frac{d^{1/4} L \ \theta^*\ \log(\frac{1}{\delta})^{1/4}}{n^{1/2} \epsilon^{1/2}}$	$\frac{d^{1/3} \beta^{1/3} L^{2/3} \ \theta^*\ ^{4/3} \log(\frac{1}{\delta})^{1/3}}{n^{2/3} \epsilon^{2/3}}$	Same as left
ObjPert	Not applicable	$\frac{d L \ \theta^*\ \sqrt{\log(\frac{1}{\delta})}}{n \epsilon}$ Lower order terms and dependence on β hidden.	$\frac{\sqrt{d} L \ \theta^*\ \sqrt{\log(\frac{1}{\delta})}}{n \epsilon}$

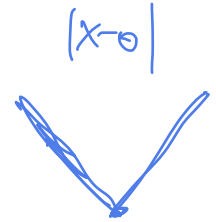
$$\left(\frac{d^{1/2}}{n \epsilon} \right)^{2/3}$$

$$\left(\frac{d^{1/2}}{n \epsilon} \right)$$

- Normalized by $1/n$ to be consistent with prior tables.
- Non-private excess risk is on the order of $\sqrt{d/n}$
- Could be $O(d/n)$

What are not quite satisfactory?

- Require the loss to be twice **differentiable**
 - Convex losses need not be even differentiable
- We did not handle the **constrained** convex ERM
$$\min_{\theta} \sum_i \ell_i(\theta)$$
$$\text{s.t. } \|\theta\|_2 \in \mathcal{B}$$
- They do not handle **non-convex** ERM problems, e.g., those that arise when optimizing deep neural networks



Gradient Descent

- Unconstrained, differentiable optimization problem

$$\min_x f(x)$$

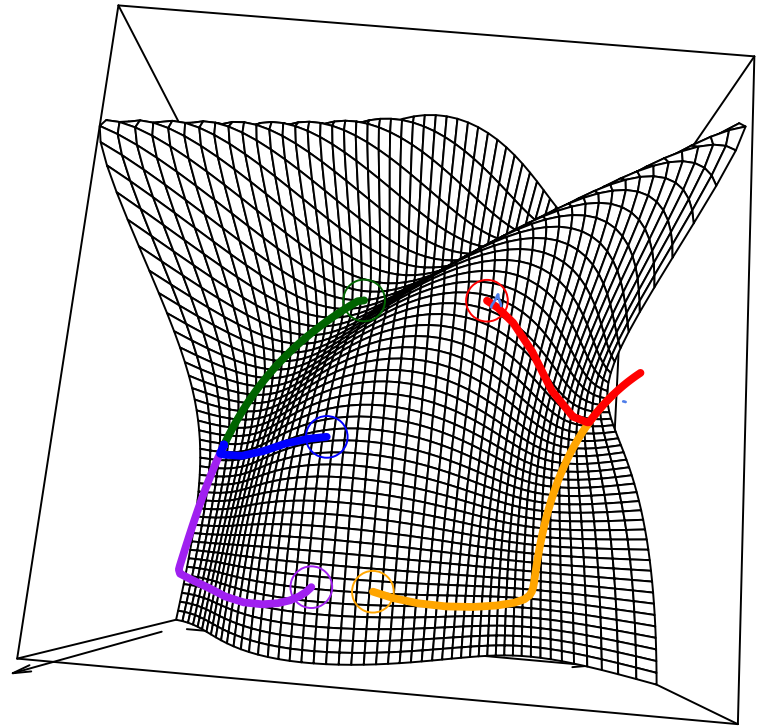
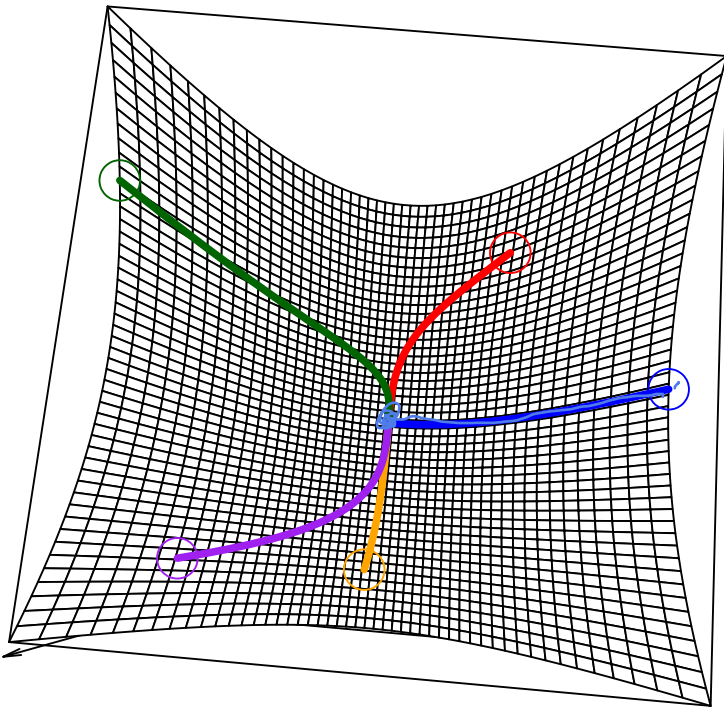
- The algorithm:

Gradient descent: choose initial point $x^{(0)} \in \mathbb{R}^n$, repeat:

$$x^{(k)} = x^{(k-1)} - \eta \cdot \nabla f(x^{(k-1)}), \quad k = 1, 2, 3, \dots$$

Stop at some point $x_t = x_t - \frac{\eta}{\eta} \cdot \nabla f(x_t)$

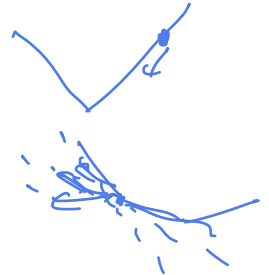
Gradient descent in convex problems vs nonconvex problems



Extensions of Gradient Descent

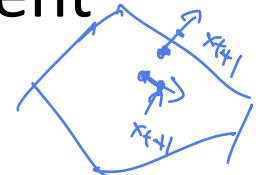
- Non-differentiable case: Subgradient descent

$$\theta_{t+1} = \theta_t - \eta_t \underset{\substack{\uparrow \\ \text{subgradient}}}{g_t}$$



- Constrained case: Projected gradient descent

$$\theta_{t+1} = \underset{\mathcal{C}}{\text{P}}(\theta_t - \eta_t g_t)$$



- Non-smooth penalty function: Proximal gradient descent

$$\theta_{t+1} = \underset{\mathcal{C}}{\text{P}}(\theta_t - \eta_t g_t)$$

$\min \sum_{i=1}^n \ell_i(\theta) + r(\theta) + \sum \lambda \|\theta_i\|?$

- Nonconvex cases: We give up theoretical guarantees but in practice it works (remarkably well)

$$f(\theta) = \sum_{i=1}^n \ell_i(\theta)$$

$$g_t = n \cdot \ell_I(\theta_t)$$

Stochastic gradient descent

$$I_t \sim \text{Unif}(\{1, 2, \dots, n\})$$

$$E[g_t] = \frac{1}{n} \sum_{i=1}^n \nabla \ell_i(\theta_t)$$

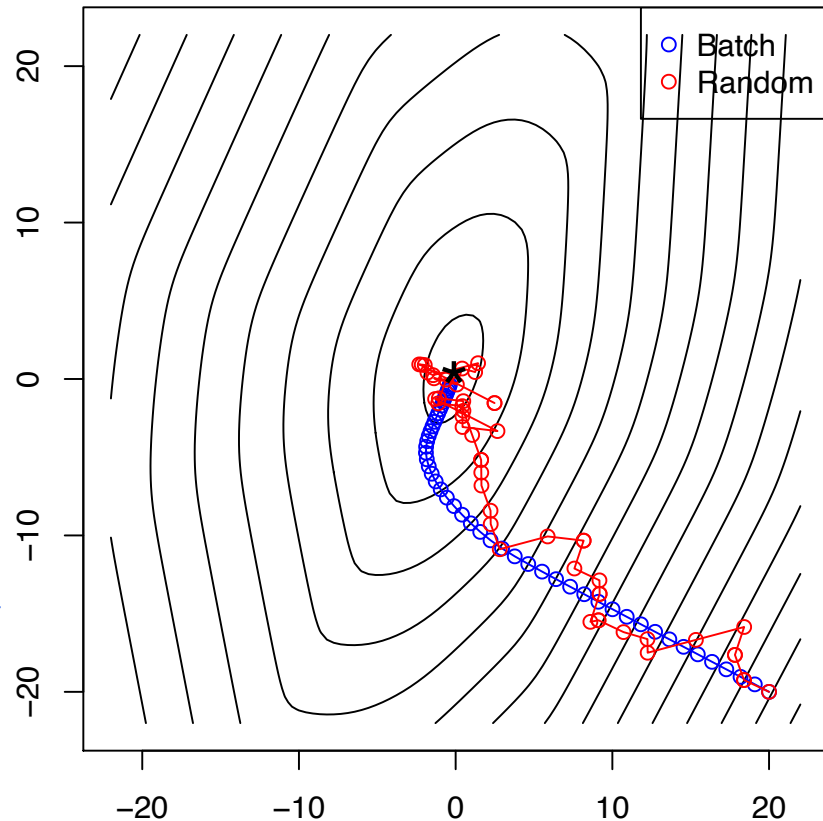
- Update rule:

$$\theta_{t+1} = \theta_t - \eta_t \hat{g}_t$$

- Assumptions:

$$E[g_t | \theta_t] = \nabla f(\theta_t)$$

$$E[\|g_t - \nabla f(\theta_t)\|_2^2 | \theta_t] \leq dG^2$$



The convergence of GD and SGD

$$\theta^* = \underset{\theta}{\operatorname{argmin}} f(\theta)$$

f is β smooth

- GD in Smooth / convex problems

$$f(\theta_T) - f(\theta^*) \leq \frac{\beta \|\theta_1 - \theta^*\|^2}{T}, \quad \eta = \frac{1}{\beta}$$

- GD in general convex problems

$$\min_{t \in \{1, \dots, T\}} [f(\theta_t) - f(\theta^*)] \leq \frac{\|\theta_1 - \theta^*\|_2 \cdot L}{\sqrt{T}}$$

f is L -lipschitz

- SGD in general convex problems

$$\min_{t \in \{1, \dots, T\}} [f(\theta_t) - f(\theta^*)] \leq \frac{\|\theta_1 - \theta^*\|_2 \sqrt{L^2 + d\sigma^2}}{\sqrt{T}}$$

$$E[\|g_t - \nabla f(\theta_t)\|^2] \leq d\sigma^2$$

- SGD in strongly convex problems

- Projected version

$$\min_{t \in \{1, \dots, T\}} [f(\theta_t) - f(\theta^*)] \leq \frac{L^2 + d\sigma^2}{\lambda T}$$

f is λ -Strongly Convex

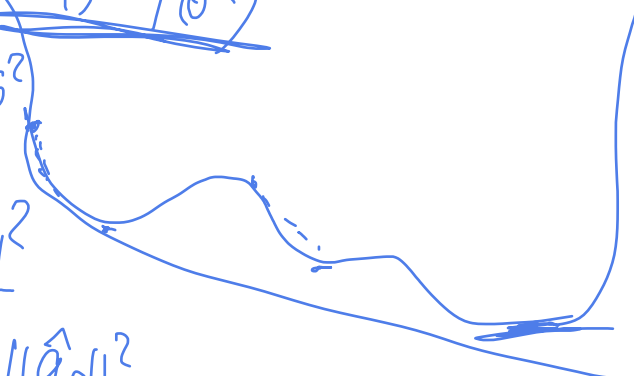
Convergence of stochastic gradient descent (in the smooth / nonconvex case)

Descent Lemma

$$\rightarrow x_{t+1} = x_t - \eta_t \hat{g}_t$$

① $E[\hat{g}_t | x_t] = \nabla f(x_t)$ $f(x_t) - f(x^*)$

② $E[\|\hat{g}_t - \nabla f(x_t)\|^2 | x_t] \leq \sigma^2$



$$f(x_{t+1}) \leq f(x_t) + \langle x_{t+1} - x_t, \nabla f(x_t) \rangle + \frac{\beta \|x_{t+1} - x_t\|^2}{2}$$

$$= f(x_t) - \eta_t \langle \hat{g}_t, \nabla f(x_t) \rangle + \frac{\beta \eta_t^2 \|\hat{g}_t\|^2}{2}$$

$$E[f(x_{t+1}) | x_t] \leq f(x_t) - \eta_t \|\nabla f(x_t)\|^2 + \frac{\eta_t^2 \beta}{2} (\|\nabla f(x_t)\|^2 + d \sigma^2)$$

$\frac{\text{Var}(x)}{2} = E[x^2] - (E[x])^2$

$$\eta_t = \eta \leq \frac{1}{\beta}$$

$$= f(x_t) - \eta \|\nabla f(x_t)\|^2 + \frac{\eta}{2} \beta \|\nabla f(x_t)\|^2 + \frac{\eta^2 \beta \sigma^2 d}{2}$$

$$= f(x_t) - \frac{\eta}{2} \|\nabla f(x_t)\|^2 + \frac{\eta^2 \beta \sigma^2 d}{2}$$

$$E[f(x_{t+1})] \leq E[f(x_t)] - \frac{\eta}{2} E \|\nabla f(x_t)\|^2 + \frac{\eta^2 \beta \sigma^2 d}{2}$$

add up $t=1, \dots, T$

Convergence of stochastic gradient descent (in the smooth / nonconvex case)

- Descent Lemma

add-up

$$\cancel{\mathbb{E} f(x_{t+1})} \leq \mathbb{E} f(x_t) - \frac{\eta}{2} \mathbb{E} [\|\nabla f(x_t)\|^2] + \frac{\beta \eta^2}{2} d G^2$$

$$\cancel{\mathbb{E} f(x_{t+2})} \leq \cancel{\mathbb{E} f(x_{t+1})} - \frac{\eta}{2} \mathbb{E} [\|\nabla f(x_{t+1})\|^2] + \frac{\beta \eta^2}{2} d G^2$$

$$\mathbb{E} f(x_T) \leq \mathbb{E} f(x_0) - \frac{\eta}{2} \mathbb{E} \left[\sum_{t=0}^{T-1} \|\nabla f(x_t)\|^2 \right] + \frac{T \beta \eta^2}{2} d G^2$$

$$\Rightarrow \frac{\eta}{2} \mathbb{E} \left[\sum_{t=0}^{T-1} \|\nabla f(x_t)\|^2 \right] \leq \mathbb{E} f(x_0) - \mathbb{E} f(x_T) + \frac{T \beta \eta^2}{2} d G^2$$

$$\Rightarrow \frac{\mathbb{E} \left[\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x_t)\|^2 \right]}{\min_{x \in \mathcal{K}} \|\nabla f(x)\|^2} \leq \frac{2(\mathbb{E} f(x_0) - \mathbb{E} f(x_T))}{T \eta} + \frac{\beta \eta d G^2}{\frac{\mathbb{E} \left[\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x_t)\|^2 \right]}{\min_{x \in \mathcal{K}} \|\nabla f(x)\|^2}} = \frac{2(\mathbb{E} f(x_0) - \mathbb{E} f(x_T))}{T \eta} + \frac{\beta \eta d G^2}{\frac{\mathbb{E} \left[\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x_t)\|^2 \right]}{\min_{x \in \mathcal{K}} \|\nabla f(x)\|^2}}$$

Noisy Gradient Descent Mechanism for Convex ERM

- The algorithm:

$$\theta_{t+1} = \theta_t + \eta_t \left(\underbrace{\sum_{i=1}^n \nabla l_i(\theta_t)}_{\nabla f(\theta_t)} + \mathcal{N}(0, G^2 I_d) \right)$$

- Privacy analysis:

- A composition of T Gaussian mechanisms

CS: $\nabla f(\theta)$ is L lipschitz

P-CDP: $P = \frac{L^2}{2G^2}$

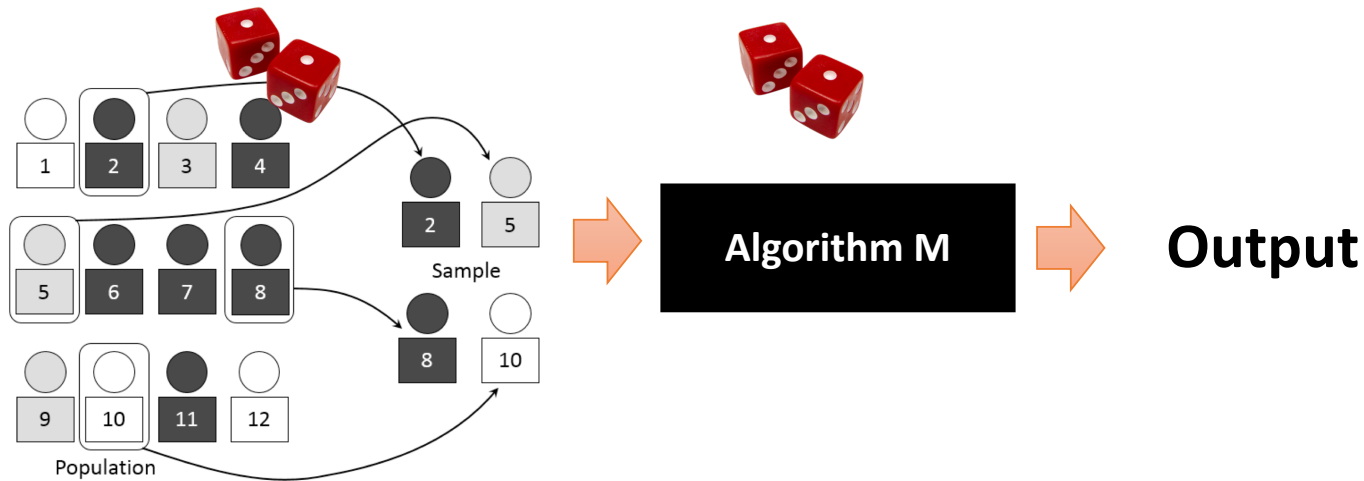
T-gaussian mech

TP-CDP: $P = \frac{TL^2}{2G^2}$

$$\frac{P}{T} = \frac{L^2}{2G^2}$$

$E[g_t | \theta_t] = \nabla f(\theta_t)$
 $E[g_t - \nabla f(\theta_t) | \theta_t] = 0$
 $E[\|g_t - \nabla f(\theta_t)\|^2 | \theta_t] = dG^2$

Privacy Amplification by Sampling



$$\mathcal{M} \circ \text{Sample} : \text{Data} \rightarrow \text{Output}$$

Subsampling Lemma: If \mathcal{M} obeys (ϵ, δ) -DP, then $\mathcal{M} \circ \text{Subsample}$ obeys that (ϵ', δ') -DP with $\delta' = \gamma \delta$

$$\epsilon' = \log(1 + \gamma(e^\epsilon - 1)) = O(\gamma \epsilon)$$

Random subset sampling vs Poisson sampling

The Noisy **Stochastic** Gradient Descent Mechanism (NoisySGD)

- Privacy analysis:
 - A composition of T subsampled gaussian mechanism.

The Noisy **Stochastic** Gradient Descent Mechanism (NoisySGD)

- Utility analysis:
 - A composition of T subsampled gaussian mechanism.

Next lecture

- Differentially private deep learning
- Knowledge transfer model of private learning