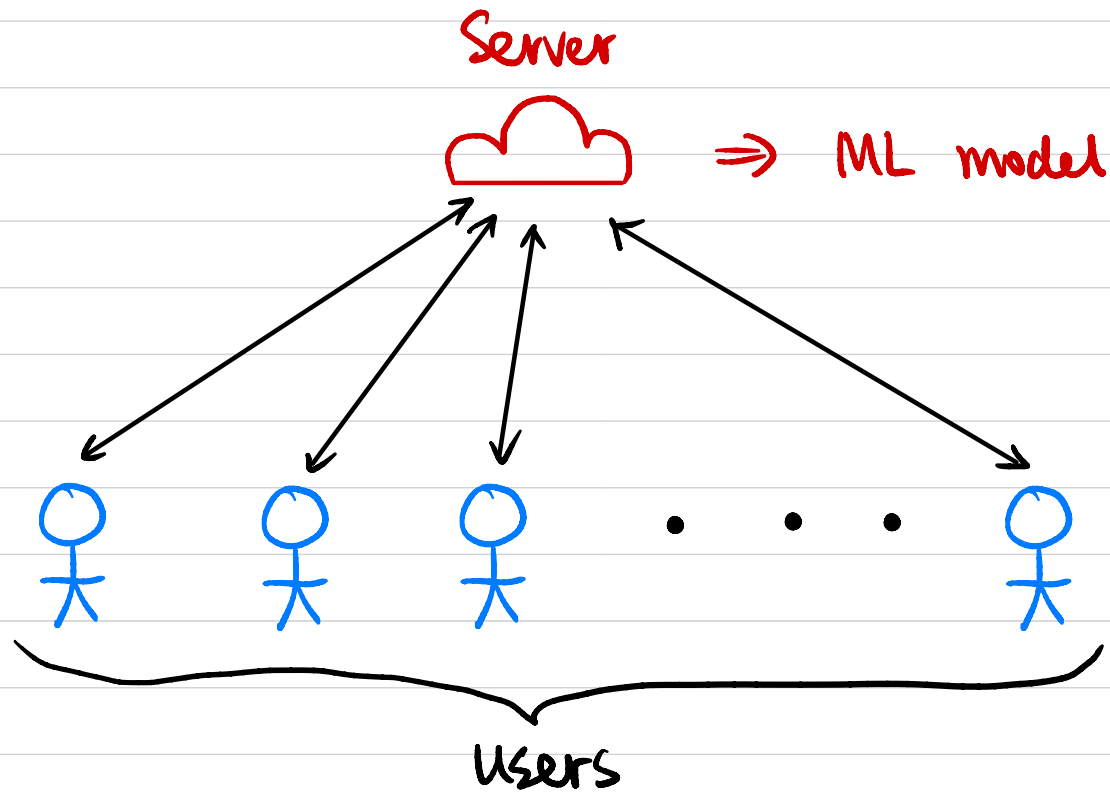


CSCI 1515 Applied Cryptography

This Lecture:

- Federated Learning
- Differential Privacy
- Elliptic Curve Cryptography

Federated Learning (FL)



Application: Google mobile keyboard prediction

Machine Learning Background

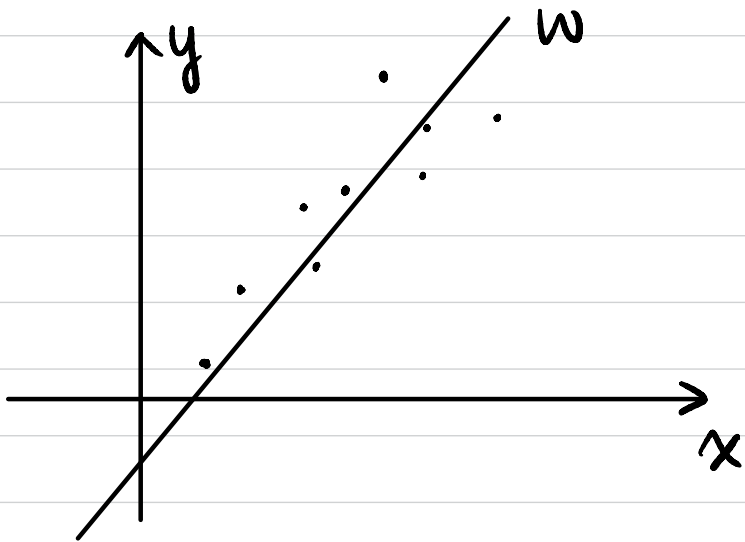
Linear Regression

Data Points (\vec{x}, y)

ML Model: coefficient vector \vec{w}

$$g(\vec{x}) = \langle \vec{x}, \vec{w} \rangle$$

Goal: Find \vec{w} that minimizes $L(\vec{w})$.



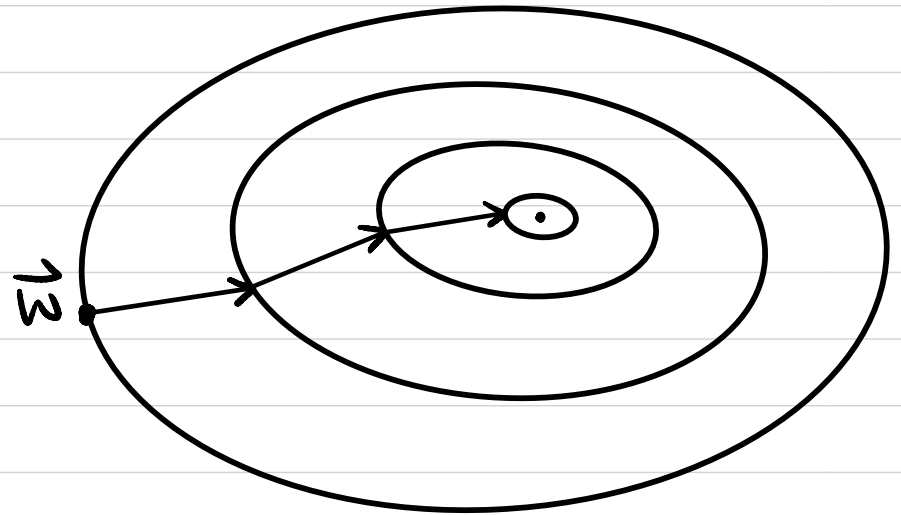
FL for Linear Regression

Stochastic Gradient Descent (SGD)

- \vec{w} initialized with arbitrary value
- Given a data point (\vec{x}_i, y_i) :

$$\vec{w} \leftarrow \vec{w} - \eta \cdot \nabla L_i(\vec{w})$$

$$\vec{w} \leftarrow \vec{w} - \eta \cdot (\langle \vec{x}_i, \vec{w} \rangle - y_i) \cdot \vec{x}_i$$



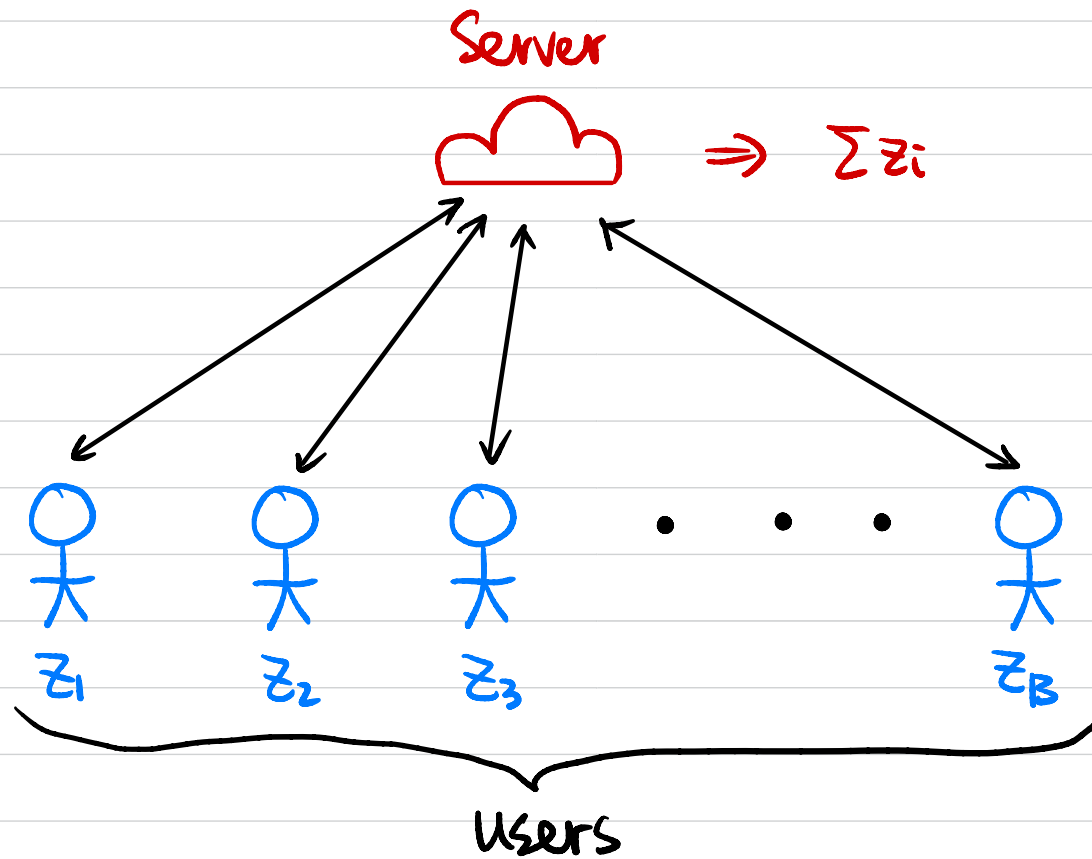
Batch SGD:

$$\vec{w} \leftarrow \vec{w} - \frac{\eta}{B} \cdot \sum_{i \in [B]} \nabla L_i(\vec{w})$$

$$\vec{w} \leftarrow \vec{w} - \frac{\eta}{B} \cdot X_B^T \cdot (X_B \cdot \vec{w} - Y_B)$$

$$\begin{bmatrix} \vec{x}_i^T \\ 1 \end{bmatrix} \cdot \left(\begin{bmatrix} -\vec{x}_i \\ 1 \end{bmatrix} \begin{bmatrix} \vec{w} \\ 1 \end{bmatrix} - \begin{bmatrix} y_i \end{bmatrix} \right)$$

Secure Aggregation



Potential Issues?

FL for Logistic Regression

SGD:

$$\vec{w} \leftarrow \vec{w} - \eta \cdot \nabla L_i(\vec{w})$$

$$\vec{w} \leftarrow \vec{w} - \eta \cdot (f(\langle \vec{x}_i, \vec{w} \rangle) - y_i) \cdot \vec{x}_i$$

Batch SGD:

$$\vec{w} \leftarrow \vec{w} - \frac{\eta}{B} \cdot \sum_{i \in [B]} \nabla L_i(\vec{w})$$

$$\vec{w} \leftarrow \vec{w} - \frac{\eta}{B} \cdot X_B^T \cdot (f(X_B \cdot \vec{w}) - Y_B)$$

$$\begin{bmatrix} 1 \\ \vec{x}_i^T \end{bmatrix} \cdot \left(f \left(\begin{bmatrix} 1 \\ \vec{x}_i \end{bmatrix} \begin{bmatrix} -1 \\ 3 \end{bmatrix} \right) - \begin{bmatrix} y_i \end{bmatrix} \right)$$

Differential Privacy

Name	Age	Gender	Race	Weight	ZIP	Disease
Alice						
Bob						
Charlie						
David						
Emily						
Fiona						

Want to make the (sensitive) data public / available to others
(e.g. for medical study).

Attempt 1: "Anonymize" the data.

Delete personally identifiable information (PII): name, DOB, ...

Attempt 2: Only answer aggregate statistics queries.

Privacy Guarantee?

Access to the output shouldn't enable one to learn anything about an individual compared to one without access.

Is this possible?

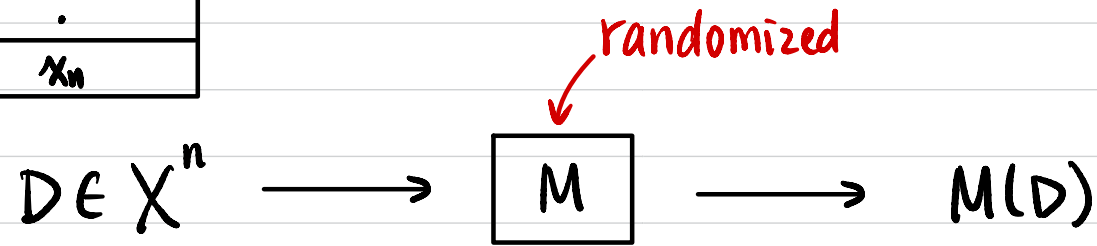
Privacy Guarantee?

Access to the output shouldn't enable one to learn ^{much more} ~~anything~~ about an individual compared to one ~~without~~ access.

with access to the output computed on a database without the individual.

Differential Privacy

x_1
x_2
\vdots
x_n

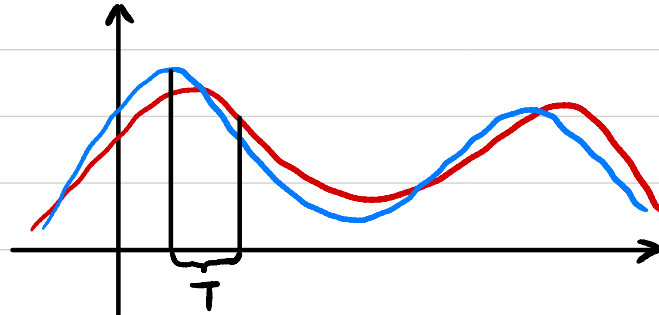


Def ϵ -Differential Privacy for a randomized mechanism:

\forall neighboring datasets D_1 & D_2 (differing in one row),

$\forall T \subseteq \text{range}(M)$,

$$\Pr[M(D_1) \in T] \leq e^\epsilon \cdot \Pr[M(D_2) \in T]$$



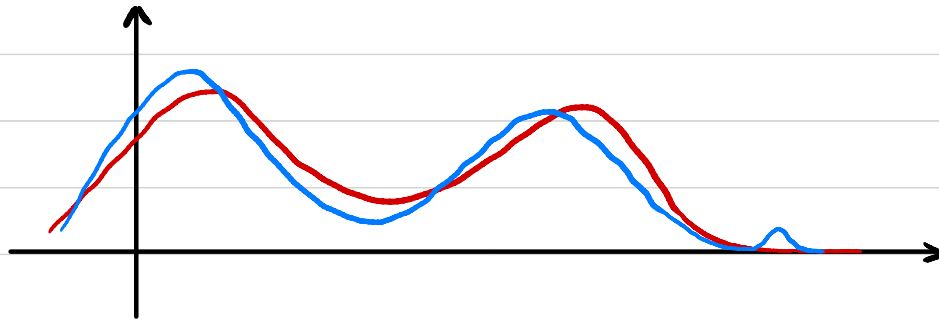
Differential Privacy

Def (ϵ, δ) - Differential Privacy for a randomized mechanism:

\forall neighboring datasets D_1 & D_2 (differing in one row),

$\forall T \subseteq \text{range}(M)$,

$$\Pr[M(D_1) \in T] \leq e^\epsilon \cdot \Pr[M(D_2) \in T] + \delta$$



Is a bigger ϵ better for privacy, or worse?

Is a bigger δ better for privacy, or worse?

Randomized Response

Counting query: What percentage of individuals satisfy predicate P ?

For each row x_i :

① Sample $b \leftarrow \{0, 1\}$

② If $b=0$, then $y_i := P(x_i)$

Otherwise, $y_i \leftarrow \{0, 1\}$

$M(D) := (y_1, y_2, \dots, y_n)$

Thm Randomized Response is $\ln 3$ -DP.

How to estimate the query output?

How to make the mechanism more private?

Laplace Mechanism

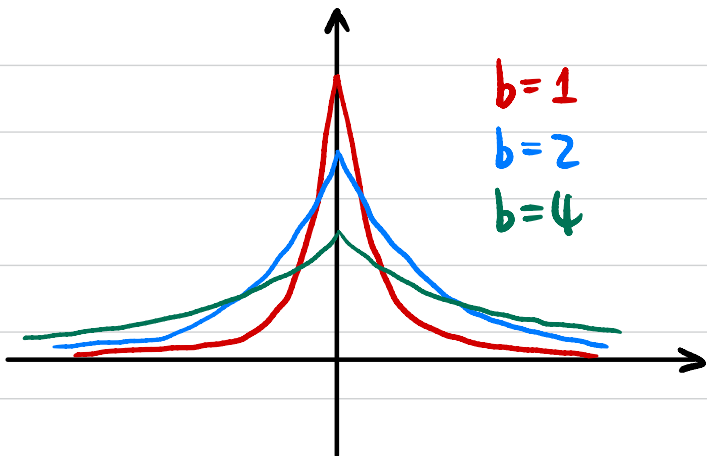
Def Sensitivity of a function $f: X^n \rightarrow \mathbb{R}$

$$\Delta f := \max_{D_1 \sim D_2} |f(D_1) - f(D_2)|$$

Laplace Mechanism: $M(D) = f(D) + \text{Lap}(\Delta f / \epsilon)$

Thm The Laplace Mechanism is ϵ -DP.

Laplace distribution:



probability distribution function

$$\text{PDF}(x) = \frac{1}{2b} \cdot \exp\left(-\frac{|x|}{b}\right)$$

For $X \sim \text{Lap}(b)$, $\Pr[|X| \geq bt] \leq \exp(-t)$

Is a bigger b better for privacy, or worse?

Composition Theorems

Thm (post-processing) If $M: X^n \rightarrow Y$ is (ϵ, δ) -DP,

$f: Y \rightarrow Z$ is an arbitrary randomized function,

then $f \circ M: X^n \rightarrow Z$ is also (ϵ, δ) -DP.

Thm (group privacy) If $M: X^n \rightarrow Y$ is $(\epsilon, 0)$ -DP,

then M is $(k \cdot \epsilon, 0)$ -DP for groups of size k .

Thm (composition) If $M_i: X^n \rightarrow Y$ is (ϵ_i, δ_i) -DP $\forall i \in [k]$,

then $M(D) := (M_1(D), \dots, M_k(D))$ is $(\sum_{i \in [k]} \epsilon_i, \sum_{i \in [k]} \delta_i)$ -DP.

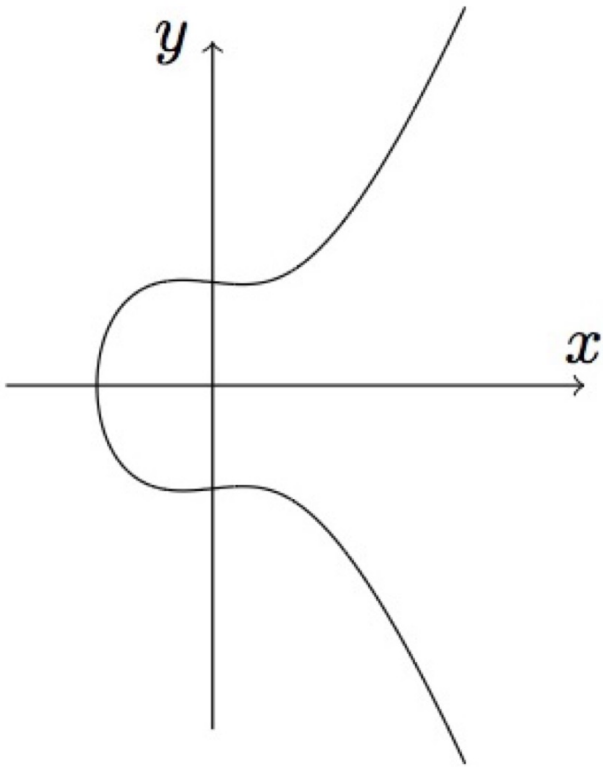
Elliptic Curve Cryptography

Cyclic group G of order q with generator g where DLOG/CDH/DDH holds.

↑
How large is q ? (128-bit security)

- Integer groups: $q \sim 2048$ bits
- Elliptic Curve groups: $q \sim 256$ bits
 - ↳ Additional structure: bilinear pairings

Elliptic Curves



$$y^2 = x^3 + ax + b$$

$$(4a^3 + 27b^2 \neq 0)$$

Example: $y^2 = x^3 - x + 9$

points: $(0, \pm 3)$

$$(1, \pm 3)$$

$$(-1, \pm 3)$$

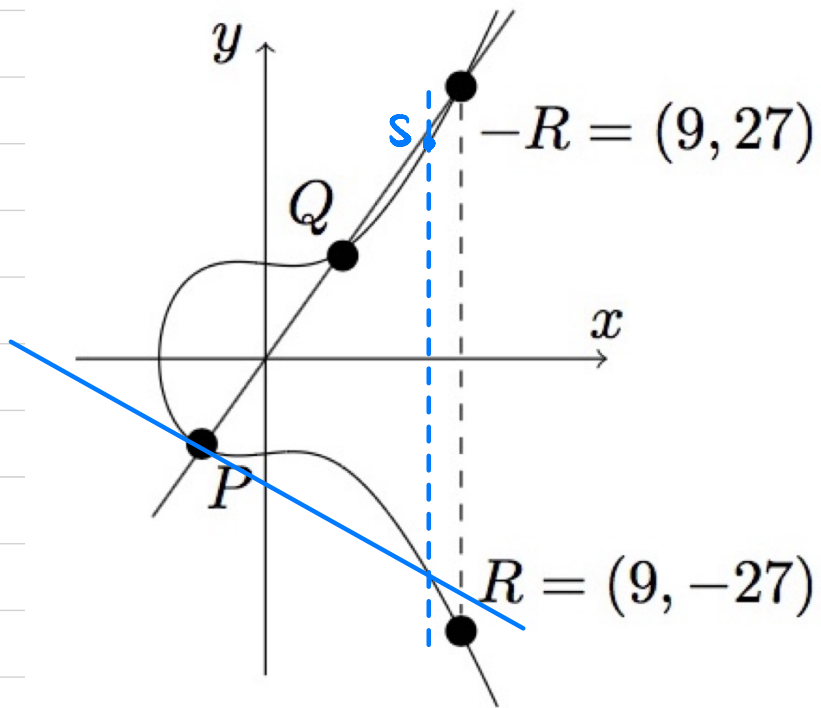
How to find rational points $(x, y) \in \mathbb{Q}^2$ on the curve?

$$x = \frac{s}{t}, \quad y = \frac{u}{v}$$

$$s, t, u, v \in \mathbb{Z}$$

Elliptic Curves

How to find rational points $(x, y) \in \mathbb{Q}^2$ on the curve?



Example: $y^2 = x^3 - x + 9$

① Chord method

$$R := P \oplus Q$$

$$\begin{aligned} P &= (-1, -3) \\ Q &= (1, 3) \end{aligned} \Rightarrow y = 3x$$

\Downarrow

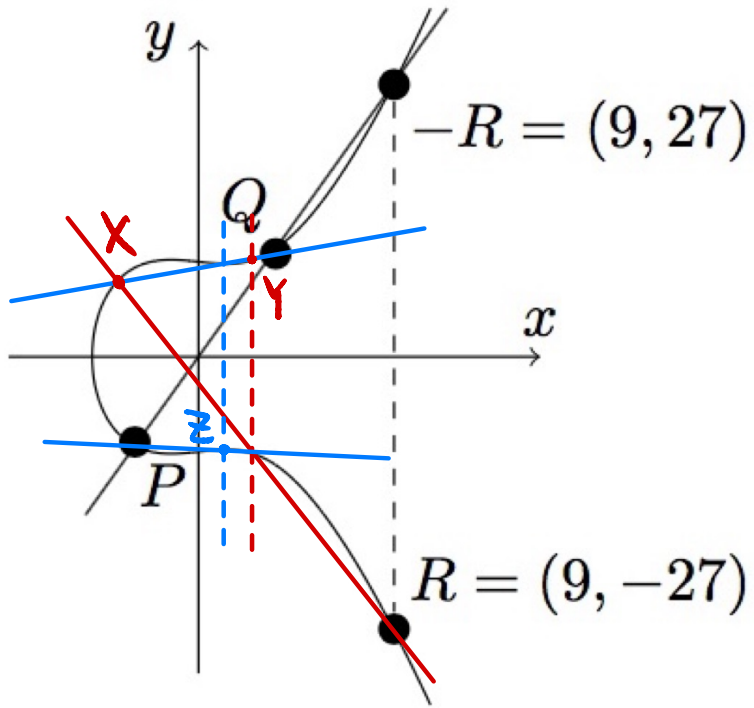
$$\begin{aligned} (3x)^2 &= x^3 - x + 9 \\ x^3 - 9x^2 - x + 9 &= 0 \end{aligned}$$

Why is the third root rational?

② tangent method

$$S := P \oplus P$$

Elliptic Curves



$$R := P \oplus Q$$

$$(P \oplus Q) \oplus X = P \oplus (Q \oplus X)$$

$$R = P \oplus Q$$

$$Z = Q \oplus X$$

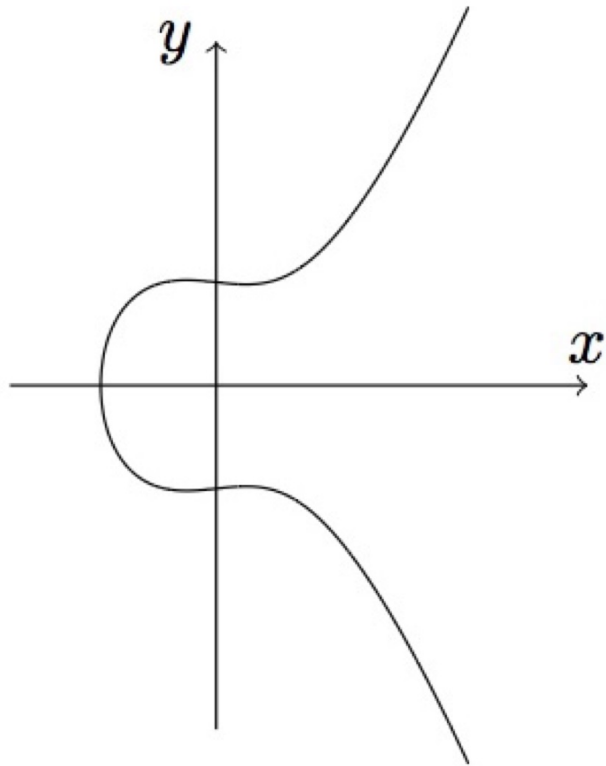
$$Y = R \oplus X$$

$$Y = P \oplus Z$$

$$P \oplus Q = Q \oplus P$$

Example: $y^2 = x^3 - x + 9$

Elliptic Curves over Finite Fields



$$y^2 = x^3 + ax + b$$

$$(4a^3 + 27b^2 \neq 0)$$

Finite field \mathbb{F}_p , $p > 3$ prime
 $\{0, 1, \dots, p-1\}$, $+$, \cdot , inverse

Elliptic curve E defined over \mathbb{F}_p : E/\mathbb{F}_p .

$$a, b \in \mathbb{F}_p$$

(x, y) is a point on the curve if

$$x, y \in \mathbb{F}_p$$

$$y^2 = x^3 + ax + b \text{ over } \mathbb{F}_p$$

Point at infinity: O

Example: $y^2 = x^3 + 1$ over \mathbb{F}_{11} .

$$E/\mathbb{F}_{11} = \{O, (-1, 0), (0, \pm 1), (2, \pm 3), (5, \pm 4), (7, \pm 5), (9, \pm 2)\}$$

Elliptic Curves over Finite Fields

Group properties:

① Closure: $\forall g, h \in G, g \circ h \in G$

② Existence of an identity:

$$\exists e \in G \text{ st. } \forall g \in G, e \circ g = g \circ e = g.$$

③ Existence of inverse:

$$\forall g \in G, \exists h \in G \text{ st. } g \circ h = h \circ g = e$$

④ Associativity:

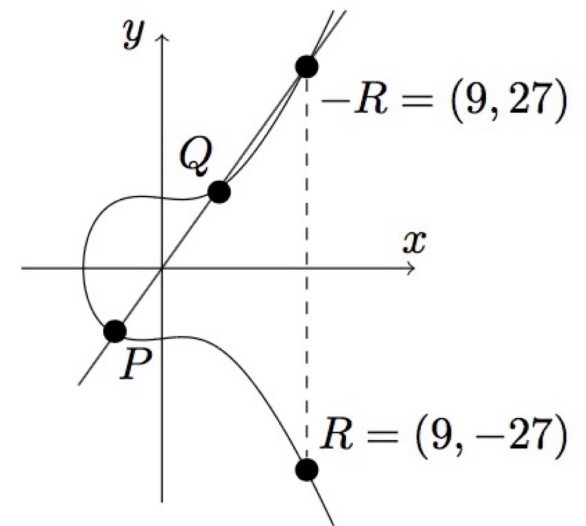
$$\forall g_1, g_2, g_3 \in G, (g_1 \circ g_2) \circ g_3 = g_1 \circ (g_2 \circ g_3)$$

⑤ Commutativity (abelian):

$$\forall g, h \in G, g \circ h = h \circ g$$

SEA algorithm: count number of points on E/\mathbb{F}_p in time $\text{polylog}(p)$.

How to compute g^a for $a \in \mathbb{Z}_q$?



Elliptic Curve Cryptography

- Curve secp256r1 (P256)
 - prime $p = 2^{256} - 2^{224} + 2^{192} + 2^{96} - 1$
 - $y^2 = x^3 - 3x + b$ b : 255-bit
 - Number of points on the curve is prime (close to p)
 - Generator point G
- Curve secp256k1
- Curve 25519