

Lecture 15.

- Basic Machine Learning
 - Optimization
 - Convexity
 - Gradient Descent
- Private Machine Learning
 - DP Gradient Descent

Announcement: Recitation in person this
Fri!

Review basic ML stuff.

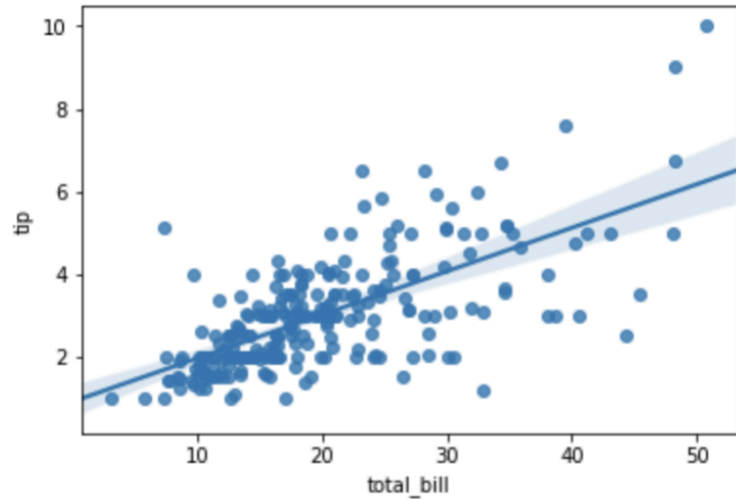
HW2 (?)

TCS 434 7pm Friday

Linear Regression

feature $\rightarrow x$ $y \leftarrow$ label

	total_bill	tip
0	16.99	1.01
1	10.34	1.66
2	21.01	3.50
3	23.68	3.31
4	24.59	3.61
5	25.29	4.71
6	8.77	2.00
7	26.88	3.12
8	15.04	1.96
9	14.78	3.23
10	10.27	1.71
11	35.26	5.00
12	15.42	1.57
13	18.43	3.00
14	14.83	3.02
15	21.58	3.92



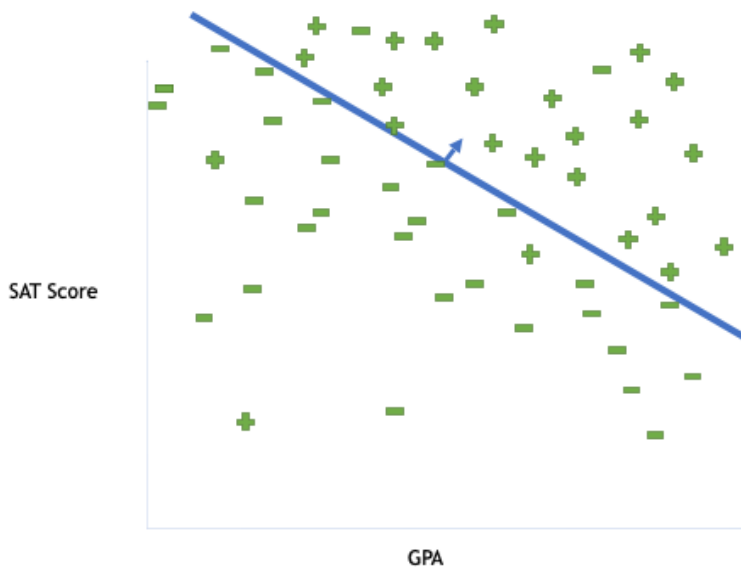
labeled example (x_i, y_i)

$x_i =$ total bill

$y_i =$ tip.

$$\min_{w \in \mathbb{R}} \sum_{i=1}^n \underbrace{(w x_i - y_i)^2}_{\text{Squared loss.}}$$

Linear Classification



$$x_i = (\text{SAT score}_i, \text{GPA}_i)$$

$$y_i = (\text{"college success"}) \\ \text{in } \{\pm 1\}$$

$$z_i = y_i \langle w, x_i \rangle$$

\uparrow \uparrow
 ± 1 want to match sign of y_i

$$\text{Logistic Loss} = \sum_{i=1}^n \ln(1 + \exp(-z_i))$$

$$\text{Hinge Loss } L(w; x) = \frac{1}{n} \sum_{i=1}^n (1 - z_i)_+$$

$$(a)_+ = \begin{cases} a & \text{if } a > 0 \\ 0 & \text{o/w.} \end{cases}$$

(Private) Optimization.

Given a data set $\mathcal{X} = (x_1, \dots, x_n)$

loss function: l

feasible set of parameters: $C \subseteq \mathbb{R}^d$
(weights)

Empirical Risk Minimization (ERM):

$$\min_{w \in C} \underbrace{L(w; \mathcal{X})}_{\text{Empirical Risk}} = \frac{1}{n} \sum_{i=1}^n l(w; x_i)$$

+ Optional
Regularization

e.g., $\lambda \|w\|_1$

Often: $C = \mathbb{R}^d$

Empirical Risk: $L(w; x) = \frac{1}{n} \sum_{i=1}^n l(w; x_i)$

Population Risk: $L(w; P) = \mathbb{E}_{x' \sim P} [l(w; x')]$

Usually, ER is a good proxy for PR.

Differential Privacy \Rightarrow Preventing overfitting.

Examples of loss functions:

$$x = ((x_1, y_1), \dots, (x_n, y_n))$$

Squared loss $L(w; x) = \frac{1}{n} \sum_{i=1}^n (\langle w, x_i \rangle - y_i)^2$

Hinge Loss $L(w; x) = \frac{1}{n} \sum_{i=1}^n (1 - y_i \langle w, x_i \rangle)_+$

(used in support vector machine)

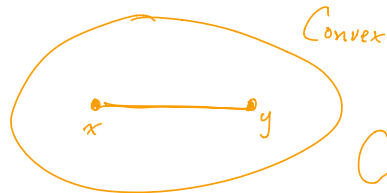
$$(a)_+ = \begin{cases} a & \text{if } a > 0 \\ 0 & \text{o/w.} \end{cases}$$

- 2 Criteria: ① Captures predictive accuracy
② Easy to optimize.

Convexity. (Sets and functions)

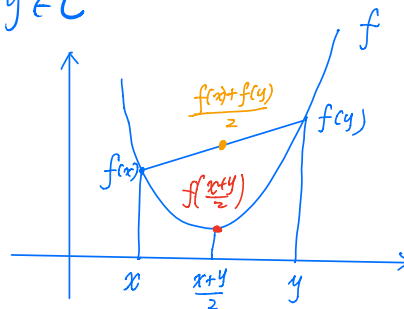
- $C \subseteq \mathbb{R}^d$ is convex if $\forall x, y \in C, t \in [0, 1]$
 $t \cdot x + (1-t) \cdot y \in C$

Non-Convex

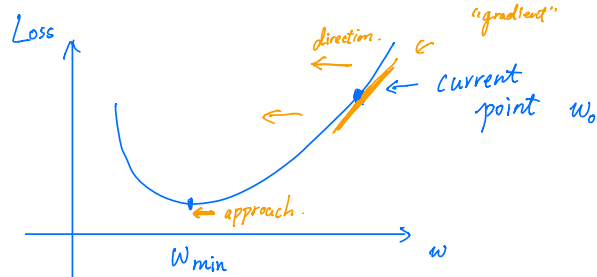


- $f: C \rightarrow \mathbb{R}$ is convex if $\forall x, y \in C$

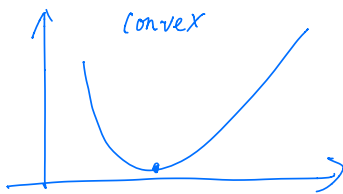
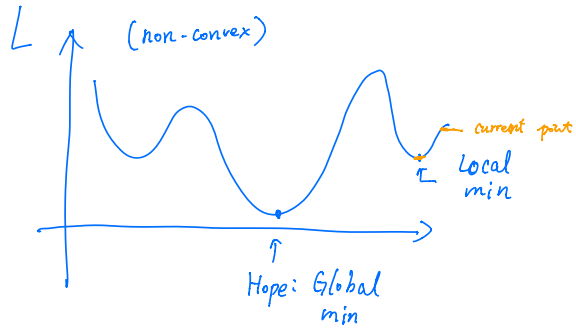
$$f\left(\frac{x+y}{2}\right) \leq \frac{f(x) + f(y)}{2}$$



"Local Search"



Loss from Neural
↓
Networks.



$$l: C \times X \mapsto \mathbb{R}$$

$l(w, x)$ measures "loss"

$$L: C \mapsto \mathbb{R}$$
$$L(w) = \frac{1}{n} \sum_{i=1}^n l(w, x_i)$$

Gradients.

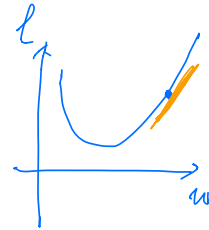
1-dim example ($w \in \mathbb{R}$)

$$l(w; (x_i, y_i)) = (wx_i - y_i)^2$$

In 1-d, gradient = derivative

$$\nabla_w l(w; (x_i, y_i)) = 2(wx_i - y_i) x_i$$

gradient
w.r.t.
parameter w



← "Chain Rule"

Multi-dim example ($w, x_i \in \mathbb{R}^d$)

$$l(w; (x_i, y_i)) = (\langle w, x_i \rangle - y_i)^2$$

$$\nabla_w l(w; (x_i, y_i)) = 2(\langle w, x_i \rangle - y_i) x_i$$

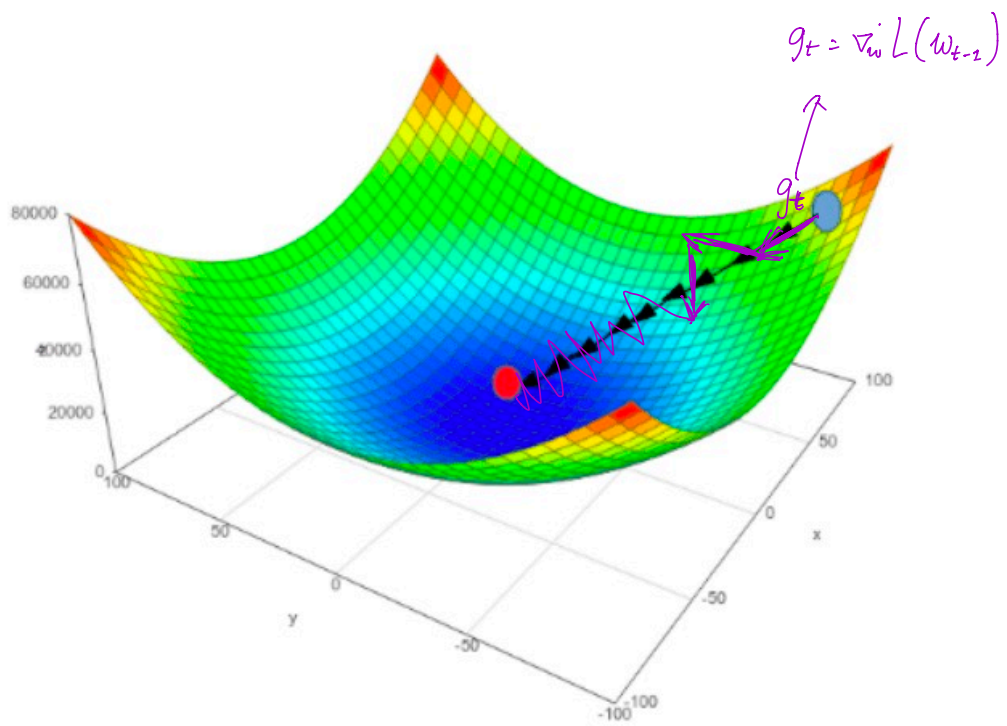
ERM minimizes

$$L(w; \mathcal{X}) = \frac{1}{n} \sum_{i=1}^n l(w; (x_i, y_i))$$

$$\nabla_w L(w; \mathcal{X}) = \frac{1}{n} \sum_{i=1}^n \nabla_w l(w; (x_i, y_i))$$

arg over all gradients

in \mathbb{R}^d



Projected Gradient Descent (PGD)

$$\text{PGD} \left(\underset{\substack{\text{Loss} \\ \downarrow}}{L}, \underset{\substack{\text{Feasible} \\ \text{Set} \\ \downarrow}}{C}, \underset{\substack{\text{Learning} \\ \text{Rate} \\ \downarrow}}{\eta}, \underset{\substack{\text{\# rounds} \\ \downarrow}}{T} \right):$$

$$\text{Init: } w_0 \in C \quad (\text{any point})$$

$$\text{For } t = 1, \dots, T:$$

gradient: $g_t = \nabla L(w_{t-1})$
(Backpropagation)

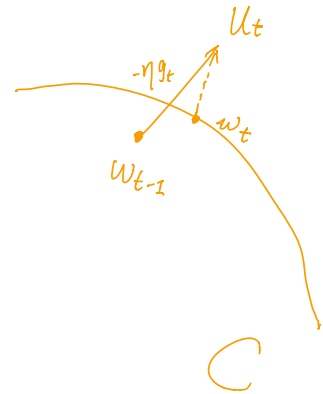
$$u_t \leftarrow w_{t-1} - \eta g_t$$

projection: $w_t \leftarrow \underset{w \in C}{\text{argmin}} \|w - u_t\|_2$

$$\text{Output} = w_T \quad \leftarrow \text{Last iterate}$$

or

$$\frac{1}{T} \sum_{t=1}^T w_t \quad \leftarrow \text{Average iterate.}$$



Robustness to noise in gradient estimation. (g_t)

Two sources of noise:

→ For efficiency:

Sample a minibatch $B \subseteq \{1, 2, \dots, n\}$

gradient estimate $\tilde{g}_t = \frac{1}{|B|} \sum_{i \in B} \nabla_w L(w_{t-1}, x_i)$

50 ↓ *50 million* ↓

→ For privacy: Add Gaussian Noise

$$\tilde{g}_t = g_t + N(0, \beta^2 I_d)$$

↑ from Gaussian mech.

In both cases,

\tilde{g}_t is an unbiased estimate of g_t

$$\mathbb{E}[\tilde{g}_t] = g_t$$