

HW 2 $\|W_i^j\|_1 = B_i$ (Problem 4)

F₁, F₂ ...

$$\text{Rad}(F_2) \leq \text{Rad}(F_{L-1}) \cdot \underline{C}$$

Online (Convex) optimization;

- decision set Ω (action space) often conv

For $t = 1, 2 \dots T$.

- learner picks $w_t \in \Omega$ ($w_1 f_1, w_2 f_2 \dots w_{k-1} f_{k-1}$)
- environment picks loss fn f_t : $\Omega \rightarrow \mathbb{R}$.
- learner suffers loss $f_t(w_t)$, $w_t \leftarrow \text{update}(w_{t-1}, f_t)$

$w_t = \text{argmin}_{w \in \Omega} \sum_{s=1}^t f_s(w)$
Follow-the-leader

Key performance measure:

$$\text{regret } R_T(u) = \sum_{t=1}^T f_t(w_t) - \sum_{t=1}^T f_t(u)$$

$$R_T(\Omega) = \max_{u \in \Omega} R_T(u)$$

Goal: want $R_T(u)$ to be $o(T)$

Connection to statistical learning:

$$S = (z_1 \dots z_T) \stackrel{\text{iid}}{\sim} D$$

$$\hat{w} = A(S) \quad \Omega : \text{hypothesis class}$$

$$l(w, z)$$

goal: generate \hat{w} s.t. $L_D(\hat{w}) = \mathbb{E}_{z \sim D} l(\hat{w}, z)$ small.

$$\text{excess loss } L_D(\hat{w}) - L_D(w^*) \quad (= o(\cdot))$$

$$w^* = \underset{w \in \Omega}{\text{argmin}} L_D(w)$$

Differences:

- online learning does not necessarily assume iid assumptions
- sometimes online learning algs are more computationally efficient (b/c it uses fast update rule)

Connection: (online to batch conversion)

given online learning alg w/ good regret guarantees, use it to construct a stat. learning alg w/ good excess loss guarantees.

Alg: input $(z_1 \dots z_T) \stackrel{\text{iid}}{\sim} D$. online learning alg A w/ action space Ω .

for $t = 1, 2, \dots, T$

A outputs w_t

$$f_t(w) = l(w, z_t)$$

$$w_1 \dots w_T$$

$$\hat{w} = \frac{1}{T} \sum_{t=1}^T w_t = \bar{w} \quad L_D(\hat{w})$$

Regret guarantees of A on u:

$$\frac{R_T(u)}{T} = \frac{1}{T} \sum_{t=1}^T \underbrace{\ell(w_t, z_t)}_{\substack{\downarrow \\ \frac{1}{T} \sum_{t=1}^T L_D(w_t)}}} - \frac{1}{T} \sum_{t=1}^T \underbrace{\ell(u, z_t)}_{\substack{\downarrow \text{Hoeffding.} \\ L_D(u)}}} = o(\mathbb{I})$$

Assuming $L_D(w)$ is convex ($\ell(w, z)$ is convex in z) $\Rightarrow o(\mathbb{I})$

$$L_D(\hat{w}) \leq \frac{1}{T} \sum_{t=1}^T L_D(w_t) \leq L_D(u) + \underbrace{\frac{R_T(u)}{T}}_{\text{concentration factors}} + o(\mathbb{I})$$

Thm: Assume $\ell(w, z) \in [0, B]$, is convex in w ,

then w.p. $1 - \delta$,

$$L_D(\hat{w}) \leq L_D(w^*) + \frac{R_T(w^*)}{T} + 2B \cdot \sqrt{\frac{2 \ln(4/\delta)}{T}} \quad \forall w^* \in \mathcal{Z}$$

Pf: Hoeffding \Rightarrow w.p. $1 - \delta/2$:

$$\left| \frac{1}{T} \sum_{t=1}^T \ell(w^*, z_t) - L_D(w^*) \right| \leq B \cdot \sqrt{\frac{2 \ln(4/\delta)}{T}}$$

Azuma \Rightarrow w.p. $1 - \delta/2$: $\ell(w_t, z_t)$ may not be independent from $\ell(w_{t+1}, z_{t+1})$

$$\left| \frac{1}{T} \sum_{t=1}^T \ell(w_t, z_t) - \frac{1}{T} \sum_{t=1}^T L_D(w_t) \right| \leq B \cdot \sqrt{\frac{2 \ln(4/\delta)}{T}}$$

Azuma's inequality: $X_1 \dots X_T \in [-B, B]$. ($B > 0$)

$\forall t, \mathbb{E}[X_t | X_1 \dots X_{t-1}] = 0$. ($X_1 \dots X_T$ is a martingale difference

sequence), w.p. $1 - \delta$:

$$\left| \sum_{t=1}^T X_t \right| \leq B \cdot \sqrt{2T \ln \frac{2}{\delta}}$$

Azuma with $X_t = L(w_t, z_t) - L_D(w_t) \in [-B, B]$

$$\mathbb{E}[X_t \mid \text{all observations up to } t-1 \text{ and } w_t] = 0.$$

$$\Rightarrow \mathbb{E}[X_t \mid X_1 \dots X_{t-1}] = 0.$$

using union bound & algebra

~~X~~

Ex: $C_1 \dots C_T \stackrel{iid}{\sim} U(\pm 1)$

X_t depend on $C_1 \dots C_{t-1}$. $\frac{\text{sign}(X_t)}{|X_t|} : X_t \in [-B, B]$

profit at t : $C_t \cdot X_t = Z_t$

$$\sum_{t=1}^T Z_t \in [\pm B \cdot \sqrt{T}]$$

pf of Azuma's inequality:

$$\mathbb{E}\left[\sum_{t=1}^T X_t\right] = 0 \quad \text{why?}$$

Mtg diff. seq. defn.

$$\mathbb{E}\left[\sum_{t=1}^T X_t\right] = \mathbb{E}_{X_1, \dots, X_T} \mathbb{E}_{X_T} \left[\sum_{t=1}^T X_t\right] = \mathbb{E}\left[\sum_{t=1}^T X_t\right] = 0.$$

$\mathbb{E}[\mid X_1 \dots X_{T-1}]$. $X_T \mid X_1 \dots X_{T-1}$ supported on $[-B, B]$.

verify $\sum_{t=1}^T X_t$ is $T \cdot B^2 - 54$

$$\forall \lambda, \mathbb{E}\left[e^{\lambda \sum_{t=1}^T X_t}\right] = \mathbb{E}_{X_1, \dots, X_T} \mathbb{E}_{X_T} \left[e^{\lambda \sum_{t=1}^{T-1} X_t} \cdot e^{\lambda X_T} \right]$$

\mathbb{Z} supported on $[a, b]$. $\mathbb{E} e^{\lambda Z} \leq e^{\frac{\lambda^2}{2} \cdot \frac{(b-a)^2}{4}}$

$$\leq \mathbb{E} \left[e^{\lambda \sum_{t=1}^T x_t} \cdot e^{\frac{\lambda^2}{2} B^2} \right]$$

$$\leq e^{\frac{\lambda^2}{2} \cdot T \cdot B^2}$$

SG \Rightarrow exponential tail property of SG r.v.'s. \star

Algorithms for online learning / online optimization.

First algorithm: online gradient descent.

Motivation: $l(w, (x, y)) = \ln(1 + e^{-y \langle w, x \rangle})$.

given $(x_i, y_i) \stackrel{iid}{\sim} D$.

can we use online learning + online to batch conversion to develop algs that can output \hat{w} with $L_D(\hat{w})$ small?

- gradient descent.

- convex fn.

Recall: fn f is conv in domain Ω (convex set), if $\forall x, y \in \Omega$,

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



$$\forall x \in \Omega, f(x) \in \mathbb{R}.$$

Ex: 1. $f(w) = \langle a, w \rangle + b$ conv

2. $f(w) = \|w\|_2$.

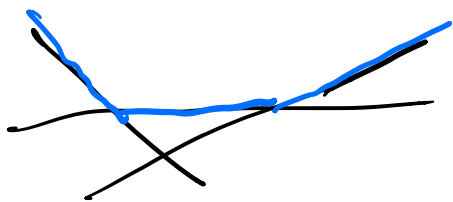
$$3. f(w) = \frac{1}{2} \|w\|_2^2 \quad (1-s.c \Rightarrow \text{cvx})$$

Basic properties of cvx fns:

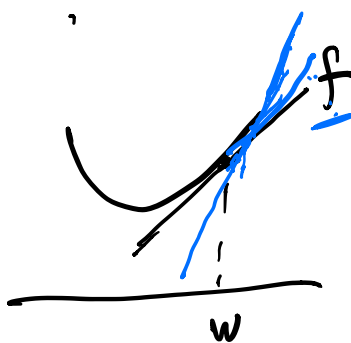
$$1. f, g \text{ cvx} \quad d. \beta \geq 0 \Rightarrow \alpha f + \beta g \text{ cvx.}$$

$$2. f \text{ cvx} \Rightarrow g(x) = f(Ax + b) \text{ cvx.}$$

$$3. \underbrace{f_1 \dots f_n}_{\text{cvx}} \Rightarrow g(x) = \max_{i=1}^n f_i(x) \text{ cvx.}$$



subgradients:



$$f(u) \geq f(w) + \langle g, u-w \rangle$$

for all x in the interior of Ω

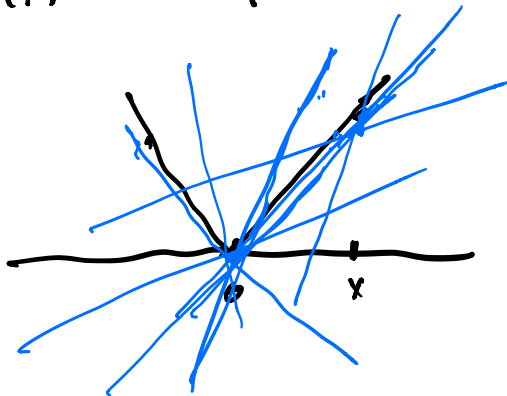
f is cvx on Ω , then for "almost all" pts x in Ω ,
subgradient of f on x

$$\partial f(x) = \{ \underline{g} : \forall y \in \Omega, f(y) \geq f(x) + \langle g, y-x \rangle \} \neq \emptyset$$

when f is differentiable at x . $\partial f(x) = \{ \nabla f(x) \}$

Ex: $f(x) = |x|$

$$\partial f(x) = \begin{cases} \{1\} & , x > 0 \\ \{-1\} & , x < 0 \\ [-1, 1] & , x = 0 \end{cases}$$



Fact: $x^* = \operatorname{argmin}_{x \in \Omega} f(x) \Leftrightarrow 0 \in \partial f(x)$

Online ^(sub) gradient descent algorithm:

Initializer $w_1 \in \Omega$.

For $t=1, 2, \dots, T$:

- choose w_t .
- receive f_t . suffer loss $f_t(w_t)$
- set $g_t \in \partial f_t(w_t)$
- update: $w'_{t+1} \leftarrow w_t - \eta g_t$.

Ω



$$w_{t+1} \leftarrow \Pi_{\Omega}(w'_{t+1})$$

$$= \arg \min_{w \in \Omega} \|w - w'_{t+1}\|_2.$$

