



THEORY OF MACHINE LEARNING

LECTURE 16

REGULARIZATION, STABILITY



SUMMARY OF GRADIENT DESCENT

- Convergence with error $O(\frac{1}{\sqrt{T}})$ after T steps for any L -Lipschitz function
- “Noisy gradient oracle” \rightarrow stochastic gradient descent
- Error of $O(1/T)$ for “smooth” convex functions (derivative is M -
Lipschitz), assuming step size $< \frac{1}{2M}$
- If function is also *strongly convex* with parameter μ , convergence bound improves to roughly $\exp(-\frac{\mu}{M}T)$ (extends to Polyak-Lojasiewicz)
- Nesterov’s “acceleration”, preconditioning via the Hessian, or by using first order proxies (AdaGrad), momentum

IMPROVEMENTS, GENERALIZATIONS

$$\exp\left(-\frac{1}{M} \cdot T\right)$$

$$\omega_{t+1} = \left(\omega_t - \eta \nabla f(\omega_t) \right) - \delta (\omega_t - \omega_{t+1})$$

$$\exp\left(-\frac{1}{\sqrt{k}} \cdot T\right)$$

1960s. Polyak's "heavy ball" method (momentum)

- Originally designed for strongly convex functions - achieves \sqrt{k} in exponent

- Second order methods, first order "proxies" (AdaGrad)

- Theme: avoid "slow" convergence - take large steps when possible

- Non-convex functions - "slip out" of local minima



- Perturbed gradient descent -- if you're not moving much via gradient descent, just make a "random jump" to a point in a neighborhood

- Can prove formally that you get out of "bad saddles" [Chi Jin, Rong Ge, M. Jordan '17].



you end up at "approximately" local minima.

$$\nabla^2 f \approx 0$$

$$\epsilon_t - \delta I$$

MANY VARIANTS OF GD



ML Hipster

@ML_Hipster

"Oh sure, going in that direction will totally minimize the objective function" —Sarcastic Gradient Descent.

CHOOSING LOSS FUNCTIONS

earlier: $\frac{\text{binary loss} \rightarrow \text{logistic loss}}{(Hw)}$

- Saw that "smoother" loss functions lead to "faster" optimization

- Utility versus niceness

Fit

dataset: $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$

$$\text{Loss}(w) = \left(\frac{1}{m} \cdot \sum_{i=1}^m \ell(w, x_i, y_i) \right) + \theta \|w\|^2$$

- Today's topic

$$\ell'(w, x_i, y_i) = \ell + \theta \|w\|^2$$

- "Nice" loss functions come with added benefit: "stability" to input changes

- Example of quadratic

- Stability is a form of "simplicity" \Rightarrow generalization

(Stability \Rightarrow concentration.)

Common

Regularizers: ℓ_2 regularizer

- entropy - neg.
- Log-barrier.
- "Self-concordant" fns.

STABILITY OF A LOSS MINIMIZATION ALGORITHM

(deterministic training procedures).

- Given examples $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$, do loss minimization

$$\underline{L(w)} = \frac{1}{m} \cdot \sum_{i=1}^m \ell(w; x_i, y_i) \longrightarrow \operatorname{argmin}_w L(w) = w^*.$$

- Can be viewed as map from examples \rightarrow parameters w

$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\} \longrightarrow w^*(S)$$

- How does changing a single (x_i, y_i) change the w^* ?

$$S' = \{(x_1, y_1), (x_2, y_2), \dots, (x_{m+1}, y_{m+1}), (x'_m, y'_m)\} \rightarrow w^*(S')$$

* Loss function is $\delta^{(m)}$ -stable (for training sets of size m) if $\forall x, y, x'_m, y'_m$ as above,

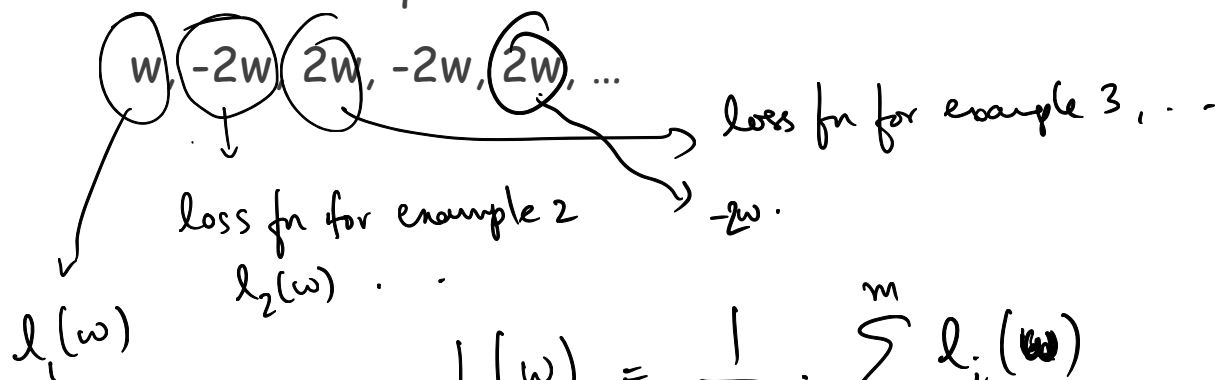
$$\|w^*(S') - w^*(S)\| \leq \delta(m).$$

$$|L(w^*(S')) - L(w^*(S))| \leq \frac{1}{m}.$$

UNDERSTANDING STABILITY -- LINEAR FUNCTIONS

- Suppose we are optimizing over $w \in [-1, 1]$

- Consider sequence of functions:



$$L(w) = \frac{1}{m} \cdot \sum_{i=1}^m l_i(w)$$

$$\arg\min_w L(w) = \begin{cases} -1 & \text{if } m \text{ is even} \\ 1 & \text{if } m \text{ is odd} \end{cases}$$

In general, changing one of the l_i 's can significantly change the w^* .

Theorem.

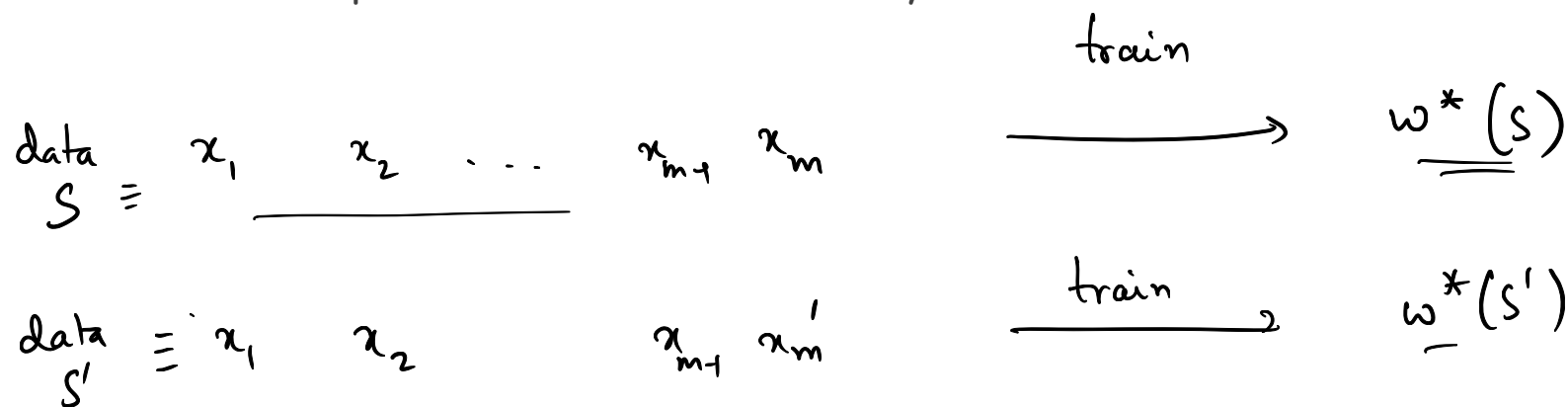
Suppose $L(w)$ is α -strongly convex and we replace $l_i(w)$ with $l'_i(w)$, such that $\|l'_i - l_i\| \leq G$, say. Then $|w^*(S') - w^*(S)| \leq \frac{G}{m\alpha}$

STABILITY IMPLIES GENERALIZATION

- Suppose our loss fn $\ell(w; x)$ is $\delta(m)$ -stable; then generalization error is $\leq \delta(m)$.

■ Recall the notion of "generalization gap"

■ Can we phrase it in terms of stability?



Generalization:

(data i from distr \mathcal{D})

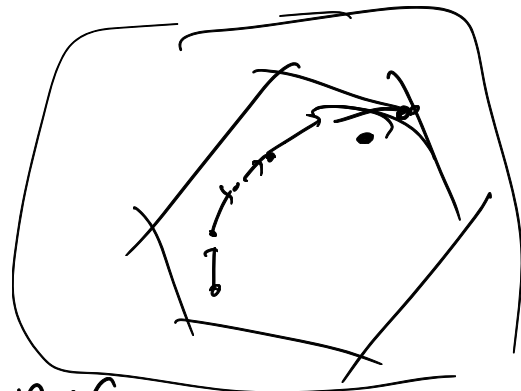
$$\text{loss on test of } w^*(S) := \mathbb{E}_{x'_m \sim \mathcal{D}} \ell(w^*(S), x'_m)$$

$$\left| \underline{\underline{\text{loss on test}}} - \underline{\underline{\text{loss on training}}} \right|$$

$$\mathbb{E}_{i \in [m]} \ell(w^*(S), x_i) \rightsquigarrow \textcircled{1} .$$

$$\textcircled{1} \equiv \underbrace{\mathbb{E}_{i \sim [m]} \ell(w^*(S), x_i) \approx \mathbb{E}_{i' \sim [m]} \mathbb{E} \ell(w^*(S \setminus x_i + x_{i'}), x_i)}_{\delta(m)}$$

Regularization:



- (1) Improves generalization gap: ("fake win" sometimes, because you might have high error before reg.)
- (2) Improves rate of convergence.

CONCENTRATION BOUNDS AND STABILITY

- [Talagrand '80s], [Boucheron, Lugosi, Massart], [Efron-Stein 60s]

Chernoff: $X_1, X_2, \dots, X_n \sim \text{iid from } \mathcal{D} \text{ with support } [a, b] \text{ \& mean } \mu$

$$\Pr \left[\left| \frac{1}{n} (X_1 + \dots + X_n) - \mu \right| > t \right] \leq e^{-\frac{t^2 n}{(a-b)^2}}$$

Talagrand: Let f be any $\frac{n}{\delta}$ -stable, \bar{u} ,
 $|f(x_1, x_2, \dots, x_n) - f(x_1, x_2, \dots, x_i', x_{i+1}, \dots, x_n)| \leq \delta$.

Then if X_1, \dots, X_n are iid r.v.s; then

$$\Pr \left[|f(x_1, \dots, x_n) - \text{median}| > t \right] \leq e^{-\frac{t^2 n}{\delta^2}}$$