

Regularization and Stability

- ERM algorithm

$$w = \operatorname{argmin}_{w \in \mathbb{R}^d} L_S(w)$$

- Regularized Loss Minimization (RLM)

$$w = \operatorname{argmin}_{w \in \mathbb{R}^d} L_S(w) + R(w)$$

where $R: \mathbb{R}^d \rightarrow \mathbb{R}$.

- The function R is called a regularizer.

• Two interpretations :

1) Similar to SRM or MDL

2) Stability

• Regarding 1), assume R is non-negative. We can consider

$$H_n = \{ w \in \mathbb{R}^d : R(w) \leq n \}$$

Then

$$H_1 \subseteq H_2 \subseteq H_3 \subseteq \dots$$

$$H = \bigcup_{n=1}^{\infty} H_n$$

- By far the most common regularizer is

$$R(w) = \lambda \|w\|_2^2$$

where $\lambda > 0$ and

$$\|w\|_2^2 = \sum_{i=1}^d w_i^2 \quad \left(\begin{array}{l} \text{Euclidean} \\ \text{norm} \\ \text{squared} \end{array} \right)$$

It is called L2-regularizer
or Tikhonov regularizer.

Ridge regression

- Linear least squares with L2-regularizer is called

Ridge regression:

$$\hat{w} = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \quad \lambda \|w\|_2^2 + \frac{1}{m} \sum_{i=1}^m (w^T x_i - y_i)^2$$

where $x_1, x_2, \dots, x_m \in \mathbb{R}^d$
and $y_1, y_2, \dots, y_m \in \mathbb{R}$.

- Ridge regression has a solution

$$\hat{w} = (\lambda m I + A)^{-1} b$$

where $A = \sum_{i=1}^m x_i x_i^T$

$$b = \sum_{i=1}^m y_i x_i$$

$$i=1$$

- To derive the solution, let

$$f(w) = \lambda \|w\|_2^2 + \sum_{i=1}^m (w^T x_i - y_i)^2.$$

- The solution satisfies

$$\nabla f(w) = 0$$



$$2\lambda w + \frac{2}{m} \sum_{i=1}^m (w^T x_i - y_i) x_i = 0$$



$$m\lambda w + \left(\sum_{i=1}^m x_i x_i^T \right) w = \sum_{i=1}^m y_i x_i$$

$$\begin{array}{c} \Updownarrow \\ (\mu\lambda I + A) w = b \end{array}$$

- A is symmetric positive semi-definite matrix.

Its eigen values are non-negative:

$$A = V^T D V$$

$V^T = V^{-1}$, D is non-negative diagonal matrix

$$D = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_d \end{pmatrix}$$

$$\mu\lambda I + A = \mu\lambda V^T V + V^T D V$$

$$= V^T (m\lambda I + D) V$$

$$= V^T D' V$$

$$D' = \begin{pmatrix} \lambda_1 + \lambda_m & & 0 \\ & \ddots & \\ 0 & & \lambda_n + \lambda_m \end{pmatrix} \leftarrow \begin{array}{l} \text{strictly} \\ \text{positive} \\ \text{diagonal} \\ \text{entries} \end{array}$$

So eigenvalues of $(m\lambda I + A)$ are positive and thus


$m\lambda I + A$ is invertible:

$$(m\lambda I + A)^{-1} = (V^T D' V)^{-1}$$

$$= V^{-1} (D')^{-1} (V^T)^{-1}$$

$$\dots^{-1}$$

$$= V' (D') V$$


This exists.

- L2 regularization is used also for logistic regression and also in neural networks.
-

Stability

- Let $Z = X \times Y$ be the space of labeled samples.
(This simplifies the notation)
- Let \mathbb{D} be a distribution over Z .

- Let $S = (z_1, z_2, \dots, z_m)$ be an i.i.d. sample from D .
- Let H be a space of predictors.
- Let $l: H \times Z \rightarrow \mathbb{R}$ be a loss function
- Let $A: Z^* \rightarrow H$ be a learning algorithm.
- Overfitting happens when

$$L_D(A(S)) - L_S(A(S))$$

is large positive number.

- More specifically we look at

$$\mathbb{E} \left[L_D(A(S)) - L_S(A(S)) \right]$$

(The expectation is w.r.t. random choice of $S \sim D^m$)

- This quantity can be expressed in alternative way, using stability.
- Suppose $z' \sim D$ is a single example independent of S

- For any $i = 1, 2, \dots, m$ we form "alternative" sample

$$S_i = (z_1, \dots, z_{i-1}, z', z_{i+1}, \dots, z_m)$$

z_i is replaced by z'

- Note: S_i is an i.i.d. sample from \mathcal{D}

Definition (An-Average-Replace-One-Stability)

Let $\varepsilon: \mathbb{N} \rightarrow \mathbb{R}$ be a non-increasing function.

Let $A: \mathcal{Z}^* \rightarrow \mathcal{H}$ be an algorithm.

Algorithm A is called on-average-
replace-one-stable if for any
distribution \mathbb{D} and any $m \geq 1$,

$$\mathbb{E} \left[\ell(A(S_I), z_I) - \ell(A(S), z_I) \right] \leq \epsilon(m)$$

where $S \sim \mathbb{D}^m$

and $I \sim \text{Uniform}(\{1, 2, \dots, m\})$.

The function $\epsilon(m)$ is
called a rate.

Theorem:

(X1.1) imitation from above definition

with ...

$$\mathbb{E}[\ell(A(S_I), z_I) - \ell(A(S), z_I)] \\ = \mathbb{E}[L_D(A(S)) - L_S(A(S))]$$

Proof:

$$\begin{aligned} \bullet \mathbb{E}[L_D(A(S))] &= \mathbb{E}[L_D(A(S_I))] \\ &= \mathbb{E}[\ell(A(S_I), z_I)] \end{aligned}$$

↑ ↑
These are independent

$$\bullet \mathbb{E}[\ell(A(S))] = \mathbb{E}[\ell(A(S_I))] + \mathbb{E}[\ell(A(S), z_I)]$$

$$+ [L_S(\cdot, \cdot)] - \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \ell(A(s), z_i) \right]$$

$$= \mathbb{E} \left[\ell(A(s), z_I) \right]$$

$I \sim \text{Uniform}(\{1, 2, \dots, m\})$

Strongly convex functions

Definition:

Let $C \subseteq \mathbb{R}^d$ be a convex set.

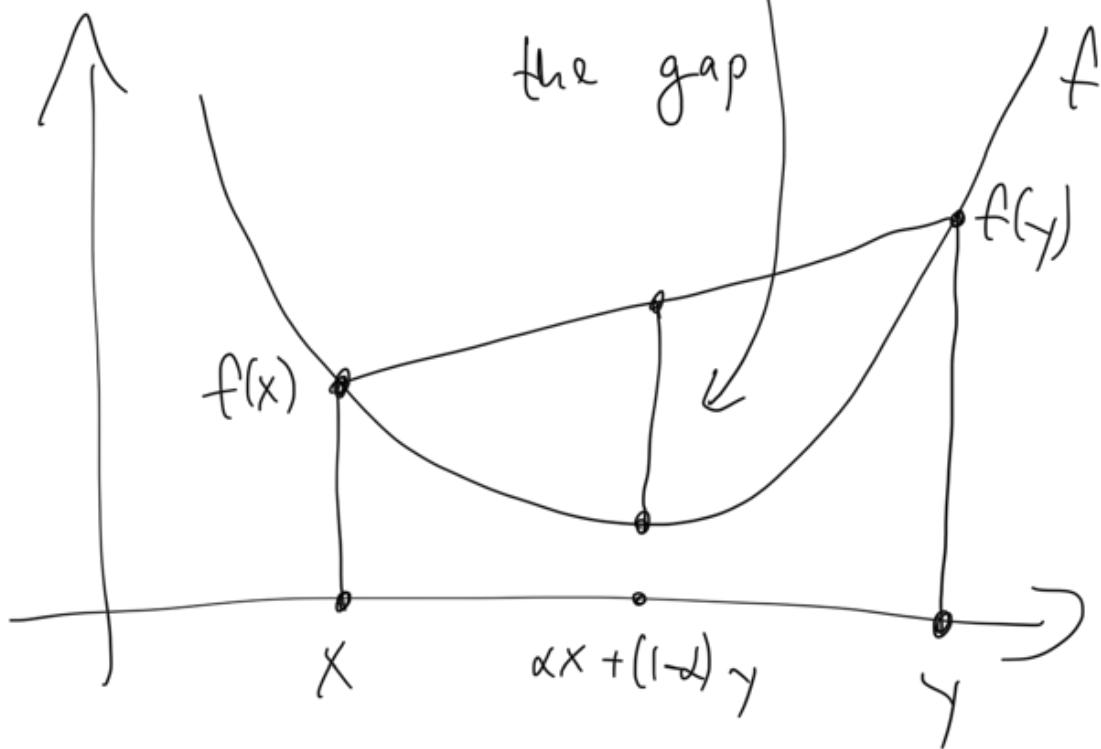
Let $\lambda \geq 0$ be a real number.

A function $f: C \rightarrow \mathbb{R}$ is called

λ -strongly convex if

for every $\alpha \in [0, 1]$ and
every $x, y \in \mathcal{C}$

$$f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y) - \frac{\lambda}{2} \alpha(1-\alpha) \|x - y\|_2^2$$



Note: Any convex function is

0-strongly convex

Lemma:

1) The function $f(x) = \lambda \|x\|_2^2$ is (2λ) -strongly convex ($\lambda \geq 0$).

2) If f is λ -strongly convex and g is convex, then $f+g$ is λ -strongly convex

3) If f is λ -strongly convex and x^* is a minimizer of f then for any x

$$f(x) - f(x^*) \geq \frac{\lambda}{2} \|x - x^*\|_2^2$$

root:

1) $f(x) = \|x\|_2^2$ satisfies

(2λ) -strong convexity condition
with equality:

$$\lambda \|\alpha x + (1-\alpha)y\|_2^2 = \alpha \lambda \|x\|_2^2 + (1-\alpha) \lambda \|y\|_2^2 - \lambda \alpha (1-\alpha) \|x - y\|_2^2$$

2) A more general fact is true: If f is λ_1 -strongly convex and g is λ_2 -strongly convex, then $f+g$ is $(\lambda_1 + \lambda_2)$ -strongly convex.

3) By definition of λ -strong convexity. For $\alpha \in (0,1)$

$$\frac{f(x^* - \alpha(x-x^*)) - f(x^*)}{\alpha} \leq f(x) - f(x^*) - \frac{\lambda}{2}(1-\alpha)\|x-x^*\|_2^2$$

$\underbrace{\hspace{10em}}_{\forall \alpha > 0}$

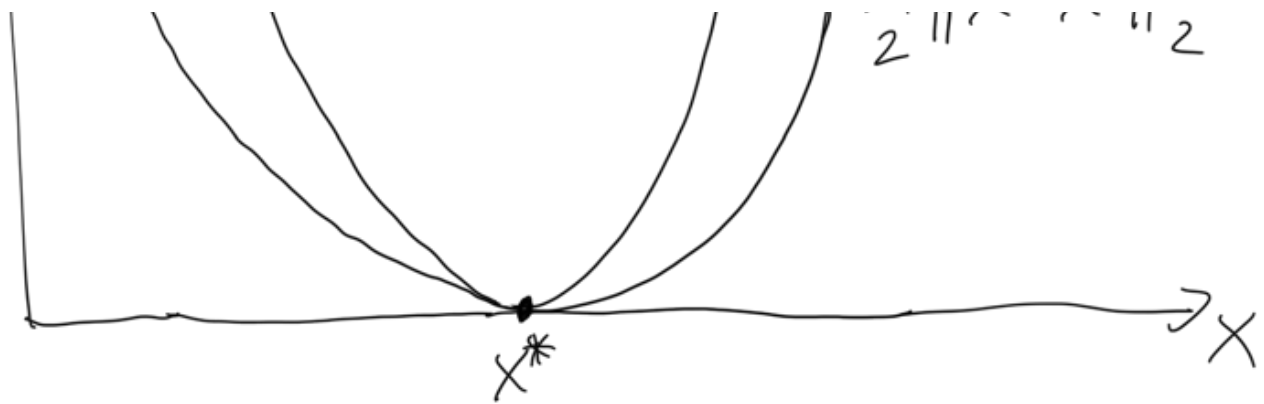
• Thus, for any $\alpha \in (0, 1)$

$$0 \leq f(x) - f(x^*) - \frac{\lambda}{2}(1-\alpha)\|x-x^*\|_2^2$$

• Thus

$$0 \leq f(x) - f(x^*) - \frac{\lambda}{2}\|x-x^*\|_2^2$$





Stability of RLM

- RLM with L_2 -regularizer

$$A(s) = \underset{w \in C}{\operatorname{argmin}} L_S(w) + \lambda \|w\|_2^2$$

- $C \subseteq \mathbb{R}^d$ is a convex set
(Subset of domain of L_S)

Lemma: Let A be RLM,

$$a(\lambda) = \underset{w \in C}{\operatorname{argmin}} L_S(w) + \lambda \|w\|_2^2.$$

$A(s)$ - argument $\rightarrow s$
 $w \in C$

If A is on-average-replace-one-stable
with rate $\epsilon(m)$ then
for all $w \in C$,

$$\mathbb{E}[L_D(A(s))] \leq L_D(w) + \lambda \|w\|_2^2 + \epsilon(m)$$

where $S \sim \mathcal{D}^m$.

Proof:

$$\bullet \mathbb{E}[L_D(A(s))] \leq \mathbb{E}[L_S(A(s))] + \epsilon(m)$$

$$\bullet L_S(A(s)) \leq L_S(A(s)) + \lambda \|A(s)\|_2^2$$

$$\nearrow \leq L_S(w) + \lambda \|w\|_2^2$$

since $A(S)$ minimizes
 $L_S(w) + \lambda \|w\|_2^2$



-
- Note: The result is similar to non-uniform PAC learning: Instead of a "high-probability result", lemma proves an upper bound "in-expectation".
 - If C is bounded, then $\lambda \|w\|^2$ is bounded and we get a "uniform" PAC-like result
(This is NOT uniform convergence, though.)
 - It remains to show that RLM is on-average-replace-one-stable

We will need convex analysis for that.

• Let $f_S: \mathcal{C} \rightarrow \mathbb{R}$,

$$f_S(w) = L_S(w) + \lambda \|w\|_2^2.$$

• Note: If $L_S: \mathcal{C} \rightarrow \mathbb{R}$ is convex then f_S is (2λ) -strongly convex.

Lemma:

Assuming $L_S: \mathcal{C} \rightarrow \mathbb{R}$ is convex then

$$\lambda \|A(s_i) - A(s)\|_2^2 \leq \frac{\ell(A(s_i), z_i) - \ell(A(s), z_i)}{m} + \frac{\ell(A(s), z') - \ell(A(s_i), z')}{m}$$

Proof:

• Let $v = A(s_i)$ and $u = A(s)$

• We lower and upper bound

$$f_S(v) - f_S(u).$$

• Since f_S is (2λ) -strongly convex and u is the minimizer of f_S , for any $w \in C$:

$$f_S(v) - f_S(u) \geq \lambda \|u - v\|^2$$

• $f_S(v) - f_S(u) =$

$$= L_S(v) + \lambda \|v\|_2^2 - (L_S(u) + \lambda \|u\|^2)$$

$$= L_{S_i}(v) + \lambda \|v\|_2^2 - (L_{S_i}(u) + \lambda \|u\|^2)$$

$$+ \frac{l(v, z_i) - l(u, z_i)}{m} + \frac{l(u, z') - l(v, z')}{m}$$

$$\leq \frac{l(v, z_i) - l(u, z_i)}{m} + \frac{l(u, z') - l(v, z')}{m}$$

Since $v = \underset{w \in C}{\operatorname{argmin}} L_{S_i}(w) + \lambda \|w\|_2^2$

• Combine upper and lower bounds



Lipschitz losses

KUW

Theorem: ↙ First coordinate

Suppose $l: H \times Z \rightarrow \mathbb{R}$ is

ρ -Lipschitz in the first coordinate.

Let A be RLM algorithm

$$A(s) = \operatorname{argmin}_{w \in C} L_s(w) + \lambda \|w\|^2.$$

Then A is on-average-replace-one-stable

with rate $\epsilon(m) = \frac{2\rho^2}{\lambda m}$

Proof:

• Let $v = A(s_i)$ and $u = A(s)$

• Since l is ρ -Lipschitz,

$$l(v, z_i) - l(u, z_i) \leq \rho \|u - v\|$$

and

$$l(u, z'_i) - l(v, z'_i) \leq \rho \|u - v\|$$

• By previous lemma,

$$\lambda \|u - v\|_2^2 \leq \frac{2\rho \|u - v\|_2}{m}$$

• Therefore

$$\|u - v\|_2 \leq \frac{2\rho}{\lambda m}$$

• Thus

$$\ell(A(s_i), z_i) - \ell(A(s), z_i)$$

$$\leq \ell(v, z_i) - \ell(u, z_i)$$

$$\leq \rho \|u - v\|_2$$

$$\leq \frac{2\rho^2}{\lambda m} \quad \circ$$

• Take expectation over S and I :

$$\mathbb{E} \left[\ell(A(s_I), z_I) - \ell(A(s), z_I) \right] \leq \frac{2\rho^2}{\lambda m} \circ$$



Corollary: ↙ First coordinate

Let $l: C \times Z \rightarrow \mathbb{R}$ be convex and β -Lipchitz in the first coordinate.

Let $A: Z^* \rightarrow \mathbb{R}$ be the RLM algorithm,

$$A(s) = \underset{w \in C}{\operatorname{argmin}} L_s(w) + \lambda \|w\|_2^2.$$

Then, for any $w \in C$

$$E[L_D(A(s))] \leq L_D(w) + \lambda \|w\|_2^2 + \frac{2\beta^2}{\lambda m}$$

Corollary:

Let $C \subseteq \mathbb{R}^d$ be a convex set such that for every $w \in C$,

$$\|w\|_2 \leq B.$$

Let $l: C \times Z \rightarrow \mathbb{R}$ be convex and ρ -Lipschitz in the first coordinate.

Let $A: Z^* \rightarrow \mathbb{R}$ be the RLM algorithm

$$A(s) = \operatorname{argmin}_{w \in C} L_S(w) + \lambda \|w\|_2^2$$

where $\lambda = \sqrt{\frac{2\rho^2}{B^2 m}}$ and $m = |S|$.

Then,

$$\mathbb{E}[L_D(A(s))] \leq \inf_{w \in C} L_D(w) + \rho B \sqrt{\frac{\delta}{m}}.$$

As δ decreases, $\sqrt{\frac{\delta}{m}}$ decreases

NOTE: The term $\frac{1}{m} \sum_{i=1}^m \ell(w; x_i, y_i)$ decreases with sample size.

So if $m \geq \frac{8\beta^2 B^2}{\epsilon^2}$ then

$$\mathbb{E}[L_D(A(s))] \leq \inf_{w \in C} L_D(w) + \epsilon.$$

Smooth Non-negative Loss

Lemma:

Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be β -smooth.

Suppose f is also non-negative.

Then, for all $x \in \mathbb{R}^d$

$$\| \nabla f(x) \|^2 \leq 2\beta f(x)$$

$$\| \nabla f(x) \| \leq L \| x - y \|$$

Proof:

- Since f is β -smooth,
for any $x, y \in \mathbb{R}^d$,

$$f(y) \leq f(x) + \nabla f(x)^\top (y-x) + \frac{\beta}{2} \|x-y\|^2$$

- Substitute

$$y = x - \frac{1}{\beta} \nabla f(x)$$

and use

$$f(y) \geq 0$$

- We get

$$0 \leq f(x) + \nabla f(x)^T \left(-\frac{1}{\beta} \nabla f(x) \right) + \frac{\beta}{2} \left\| \frac{1}{\beta} \nabla f(x) \right\|^2$$



$$0 \leq f(x) - \frac{1}{\beta} \|\nabla f(x)\|^2 + \frac{1}{2\beta} \|\nabla f(x)\|^2$$



$$\|\nabla f(x)\|^2 \leq f(x)$$

etc