

Handout 3: Learnability and stability in learning problems

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim. To introduce concepts PAC, fitting vs stability trade off, stability, and their implementation in regularization problems and convex problems.

Reading list & references:

- Shalev-Shwartz, S., & Ben-David, S. (2014). Understanding machine learning: From theory to algorithms. Cambridge university press.
– Ch. 2, 3, 13
- Bousquet, O., Boucheron, S., & Lugosi, G. (2003). Introduction to statistical learning theory. In Summer school on machine learning (pp. 169-207). Berlin, Heidelberg: Springer Berlin Heidelberg. (Suitable for PG students)

1. LEARNABLE PROBABLY APPROXIMATELY CORRECT (PAC) LEARNING

Note 1. We formally define the broad learning problem we will work on.

Note 2. Learning algorithms \mathcal{A} use training data sets \mathcal{S} which may miss characteristics of the unknown data-generating process g . Essentially, approximations (aka errors) are inevitable; some characteristics of data generating process g will be missed even if we use a very representative training set \mathcal{S} . We cannot hope that the learning algorithm will find a hypothesis whose error is smaller than the minimal possible error.

Note 3. The PAC learning problem requires no prior assumptions about the data-generating process g . It requires that the learning algorithm \mathcal{A} will find a predictor $\mathcal{A}(\mathcal{S})$ whose error is not much larger than the best possible error of a predictor in some given benchmark hypothesis class. So essentially, in practice, the researcher's effort falls on the hypothesis class \mathcal{H} .

Definition 4. (Agnostic PAC Learnability for General Loss Functions) A hypothesis class \mathcal{H} is agnostic PAC learnable with respect to a domain \mathcal{Z} of size m and a loss function $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}_+$, if there exist a function $m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$ and a learning algorithm \mathcal{A} with the following property:

- for every $\epsilon \in (0, 1)$, $\delta \in (0, 1)$, and distribution g over \mathcal{Z} , when running algorithm \mathcal{A} given training set $\mathcal{S}_m = \{z_1, \dots, z_m\}$ with $z_i \stackrel{\text{iid}}{\sim} g$ for $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ then \mathcal{A} returns $\mathcal{A}(\mathcal{S}) \in \mathcal{H}$ such as

$$(1.1) \quad \Pr_{\mathcal{S} \sim g} \left(R_g(\mathcal{A}(\mathcal{S}_m)) - \min_{h' \in \mathcal{H}} (R_g(h')) \leq \epsilon \right) \geq 1 - \delta$$

Note 5. It may be easier to work with expectations: by using Markov inequality (1.1) becomes

$$(1.2) \quad \Pr \left(R_g(\mathfrak{A}(\mathcal{S}_m)) - \min_{h' \in \mathcal{H}} (R_g(h')) \leq \epsilon \right) \geq 1 - \underbrace{\frac{1}{\epsilon} \mathbb{E}_{\mathcal{S} \sim g} \left(R_g(\mathfrak{A}(\mathcal{S}_m)) - \min_{h' \in \mathcal{H}} (R_g(h')) \right)}_{\delta}$$

and hence we need to work with expectations and bounded above.

Hint: Markov inequality $\Pr(X \geq a) \leq \frac{1}{a} \mathbb{E}(X)$ for $X > 0$.

Remark 6. The accuracy parameter ϵ determines how far the output rule h can be from the optimal one (this corresponds to the ‘approximately correct’). The confidence parameter δ indicates how likely the classifier is to meet that accuracy requirement (corresponds to the “probably” part of “PAC”).

2. ANALYSIS OF THE RISK BASED ON THE TRADE-OFF FITTING VS STABILITY

Note 7. Let $R^* = \min_{h \in \mathcal{H}} (R(h))$ be an ideal/optimal (hence minimum) Risk, \mathfrak{A} be a learning algorithm, and $\mathfrak{A}(\mathcal{S})$ the learnign rule from \mathfrak{A} under training dataset \mathcal{S} . The Risk of a learning algorithm \mathfrak{A} can be decomposed as

$$(2.1) \quad R_g(\mathfrak{A}(\mathcal{S})) - R^* = \underbrace{\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) - R^*}_{(I)} + \underbrace{R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))}_{(II)}$$

Note 8. Over-fitting can be represented by $R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))$. However, $\hat{R}_{\mathcal{S}}(\cdot)$ is a random variable and for our computational convenience here we focus on its expectation w.r.t. $\mathcal{S} \sim g$. Hence, we provide the following (arguable) definition of over-fitting on which we base our analysis .

Definition 9. For a learning algorithm \mathfrak{A} , as a measure of over-fitting we consider the expected difference between true Risk and empirical Risk

$$(2.2) \quad \mathbb{E}_{\mathcal{S} \sim g} \left(R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right)$$

Definition 10. We say that learning algorithm \mathfrak{A} suffers from over-fitting when (2.2) is ‘too’ large.

Note 11. The expected Risk of a learning algorithm $\mathfrak{A}(\mathcal{S})$ can be decomposed as

$$(2.3) \quad \mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S}))) = \underbrace{\mathbb{E}_{\mathcal{S} \sim g} \left(\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right)}_{(I)} + \underbrace{\mathbb{E}_{\mathcal{S} \sim g} \left(R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right)}_{(II)}$$

(by applying expectations in (2.1) and ignoring R^*), where (I) indicates how well $\mathfrak{A}(\mathcal{S})$ fits the training set \mathcal{S} , and (II) indicates the discrepancy between the true and empirical risks of $\mathfrak{A}(\mathcal{S})$. In Section 3, we argue that (II) measures the over-fitting and a certain type of stability of $\mathfrak{A}(\mathcal{S})$.

Note 12. The following result connects the task of upper bounding (and minimizing this bound) the Expected Risk in (2.3) with PAC learning (4).

Theorem 13. *Let \mathfrak{A} be a learning algorithm that guarantees the following:*

- If $m \geq m_{\mathcal{H}}(\epsilon)$ then for every distribution g , it is

$$E_{\mathcal{S} \sim g} \left(R_g(\mathfrak{A}(\mathcal{S}_m)) - \min_{h' \in \mathcal{H}} (R_g(h')) \right) \leq \epsilon$$

Then \mathfrak{A} it satisfied Definition 4:

- for every $\delta \in (0, 1)$, if $m \geq m_{\mathcal{H}}(\epsilon\delta)$ then

$$\Pr_{\mathcal{S} \sim g} \left(R_g(\mathfrak{A}(\mathcal{S}_m)) - \min_{h' \in \mathcal{H}} (R_g(h')) \leq \epsilon \right) \geq 1 - \delta$$

Proof. Let $\xi = R_g(\mathfrak{A}(\mathcal{S}_m)) - \min_{h' \in \mathcal{H}} (R_g(h'))$. From Markov's inequality $\Pr(\xi \geq E(\xi)/\delta) \leq \frac{1}{E(\xi)/\delta} E(\xi) = \delta$. Namely, $\Pr(\xi \leq E(\xi)/\delta) = 1 - \delta$. But it is given that if $m \geq m_{\mathcal{H}}(\epsilon\delta)$ for every distribution g it is $E(\xi) \leq (\epsilon\delta)$. So by substitution $\Pr(\xi \leq (\epsilon\delta)/\delta) = 1 - \delta$ implies $\Pr(\xi \leq \epsilon) = 1 - \delta$. Now substitute back ξ and we conclude the proof. \square

Note 14. We aim to design a learning algorithm $\mathfrak{A}(\mathcal{S})$ that both fits the training set and is stable; i.e. minimizing both terms in (2.3). As seen later, there may be a trade-off between (I) and (II); expected empirical risk term and stability term. Hence, we aim to upper bound (2.3) by upper bounding (I) and (II) individually.

3. STABILITY AND OVER-FITTING

Notation 15. Let $\mathcal{S} = \{z_1, \dots, z_m\}$ be a training sample, and \mathfrak{A} be a learning algorithm with output $\mathfrak{A}(\mathcal{S})$.

Note 16. A learning algorithm can be stable if a small change of the input to the algorithm does not change the output of the algorithm much. Formalizing this in maths, we can say that if $\mathcal{S}^{(i)} = \{z_1, \dots, z_{i-1}, z', z_{i+1}, \dots, z_m\}$ is another training dataset equal to \mathcal{S} but the i th element which is replaced by another $z' \sim g$, then a good learning algorithm \mathfrak{A} would produce a small value of

$$\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \geq 0$$

Definition 17. We say that a learning algorithm \mathfrak{A} is **on-average-replace-one-stable** with rate $\epsilon : \mathbb{N} \rightarrow \mathbb{R}$ if for every distribution g

$$E_{\substack{\mathcal{S} \sim g \\ z' \sim g, i \sim U\{1, \dots, m\}}} \left(\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \right) \leq \epsilon(m)$$

where $\epsilon(\cdot)$ has to be a decreasing function.

Note 18. Following we discuss the association of stability and over-fitting, based on the over-fitting Definition 10.

Theorem 19. For any learning algorithm \mathfrak{A}

$$E_{\mathcal{S} \sim g} \left(R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{D}}(\mathfrak{A}(\mathcal{S})) \right) = E_{\substack{\mathcal{S} \sim g, z' \sim g \\ i \sim U\{1, \dots, m\}}} \left(\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \right)$$

where g is a distribution, $\mathcal{S} = \{z_1, \dots, z_m\}$ and $\mathcal{S}^{(i)} = \{z_1, \dots, z_{i-1}, z', z_{i+1}, \dots, z_m\}$ are training datasets with $z', z_1, \dots, z_m \stackrel{\text{iid}}{\sim} g$.

Proof. As $z', z_1, \dots, z_m \stackrel{\text{iid}}{\sim} g$, then for every i

$$\mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S}))) = \mathbb{E}_{\substack{\mathcal{S} \sim g \\ z' \sim g}} (\ell(\mathfrak{A}(\mathcal{S}), z')) = \mathbb{E}_{\substack{\mathcal{S} \sim g \\ z' \sim g}} (\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i))$$

and

$$\mathbb{E}_{\mathcal{S} \sim g} (\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))) = \mathbb{E}_{\substack{\mathcal{S} \sim g \\ i \sim U\{1, \dots, m\}}} (\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i))$$

□

4. IMPLEMENTATION IN REGULARIZED LOSS LEARNING PROBLEMS

Definition 20. Assume $(\mathcal{H}, \mathcal{Z}, \ell)$. Regularized Loss Minimization (RLM) learning rule is the one that results as the output of jointly minimizing the empirical risk $\hat{R}_{\mathcal{Z}}(h)$ and a regularization function $J : \mathcal{H} \rightarrow \mathbb{R}$ that is

$$(4.1) \quad h^* = \underbrace{\arg \min}_{h \in \mathcal{H}} (\hat{R}_{\mathcal{S}}(h) + J(h))$$

Remark 21. The motivation for considering the regularization function J in (4.1) is to: (1.) control complexity and (2.) improve stability; as we will see later.

Note 22. We make our example more specific and narrow it to the Ridge RLM learning problem (could be LASSO, Elastic Net, etc.).

Definition 23. The Ridge RLM learning problem $(\mathcal{H}, \mathcal{Z}, \ell)$, here $\mathcal{H} = \mathcal{W} \subset \mathbb{R}^d$, uses regularization function $J(w; \lambda) = \lambda \|w\|_2^2$ with $\lambda > 0$, $w \in \mathcal{W}$ and produces learning rule

$$(4.2) \quad \mathfrak{A}(\mathcal{S}) = \arg \min_{w \in \mathcal{W}} (\hat{R}_{\mathcal{S}}(w) + \lambda \|w\|_2^2)$$

Note 24. Recall (Term 1) that the regularization function in Ridge RLM learning problem penalizes complexity. Essentially, implies a sequence of hypothesis $\mathcal{H}_1 \subset \mathcal{H}_2 \subset \dots$ with $\mathcal{H}_i = \{w \in \mathbb{R}^d : \|w\|_2 < i\}$.

Note 25. Below, we will try to analyze the behavior of Ridge RLM learning rule (4.2) w.r.t. the Risk decomposition (2.3). In particular, to upper bounded w.r.t. the shrinkage term λ , training sample size m , and other characteristics.

4.1. Bounding the empirical risk (I) in (2.3).

Note 26. From (4.2), we have

$$\begin{aligned} \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) &\leq \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) + \lambda \|\mathfrak{A}(\mathcal{S})\|_2^2 \\ &\leq \hat{R}_{\mathcal{S}}(w') + \lambda \|w'\|_2^2; \quad \forall w' \in \mathcal{W} \end{aligned}$$

and by taking expectations w.r.t. \mathcal{S} , it is

$$(4.3) \quad \mathbb{E}_{\mathcal{S} \sim g} \left(\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right) \leq R_g(w') + \lambda \|w'\|_2^2; \quad \forall w' \in \mathcal{W}$$

because $\mathbb{E}_{\mathcal{S} \sim g} \left(\hat{R}_{\mathcal{S}}(\cdot) \right) = \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\mathcal{S} \sim g} (\ell(\cdot, z_i)) = R_g(\cdot)$.

Note 27. We observe that part (I) in the expected risk decomposition (2.3), (aka the upper bound of the expected empirical risk) increases with the regularization term $\lambda > 0$. (!!!) –Although expected, it did not start well.

4.2. Bounding the empirical risk (II) in (2.3).

Note 28. We do that by constraining the loss function to be convex and Lipschitz.

Assumption 29. *The loss function $\ell(\cdot, z)$ in (4.2) is convex for any $z \in \mathcal{Z}$.*

Note 30. Let $\tilde{R}_{\mathcal{S}}(w) = \hat{R}_{\mathcal{S}}(w) + \lambda \|w\|_2^2$. $\tilde{R}_{\mathcal{S}}(\cdot)$ is 2λ -strongly convex as the sum of a convex function $\hat{R}_{\mathcal{S}}(\cdot)$ (Assumption 29) and a 2λ -strongly convex function $J(\cdot; \lambda) = \lambda \|\cdot\|_2^2$ (results directly from Definition 43).

Note 31. Because $\tilde{R}_{\mathcal{S}}(\cdot)$ is 2λ -strongly convex and $\mathfrak{A}_{\text{Ridge}}(\mathcal{S})$ is its minimizer, according to Lemma 45 in Handout (...), for $\mathfrak{A}(\mathcal{S})$ and any $w \in \mathcal{W}$, it is

$$(4.4) \quad \tilde{R}_{\mathcal{S}}(w) - \tilde{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \geq \lambda \|w - \mathfrak{A}(\mathcal{S})\|^2, \quad \forall w \in \mathcal{W}$$

Note 32. Also, for any $w, u \in \mathcal{W}$, it is

$$\begin{aligned} \tilde{R}_{\mathcal{S}}(w) - \tilde{R}_{\mathcal{S}}(u) &= \left(\hat{R}_{\mathcal{S}}(w) + \lambda \|w\|_2^2 \right) - \left(\hat{R}_{\mathcal{S}}(u) + \lambda \|u\|_2^2 \right) \\ &= \left(\hat{R}_{\mathcal{S}^{(i)}}(w) + \lambda \|w\|_2^2 \right) - \left(\hat{R}_{\mathcal{S}^{(i)}}(u) + \lambda \|u\|_2^2 \right) \\ &\quad + \frac{\ell(w, z_i) - \ell(u, z_i)}{m} + \frac{\ell(w, z') - \ell(u, z')}{m} \\ &= \tilde{R}_{\mathcal{S}^{(i)}}(w) - \tilde{R}_{\mathcal{S}^{(i)}}(u) + \frac{\ell(w, z_i) - \ell(u, z_i)}{m} + \frac{\ell(w, z') - \ell(u, z')}{m} \end{aligned}$$

Choosing $w = \mathfrak{A}(\mathcal{S}^{(i)})$ and $u = \mathfrak{A}(\mathcal{S})$, and the fact that $\tilde{R}_{\mathcal{S}^{(i)}}(\mathfrak{A}(\mathcal{S}^{(i)})) \leq \tilde{R}_{\mathcal{S}^{(i)}}(\mathfrak{A}(\mathcal{S}))$ it is

$$(4.5) \quad \tilde{R}_{\mathcal{S}}(w) - \tilde{R}_{\mathcal{S}}(u) \leq \frac{\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i)}{m} + \frac{\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z') - \ell(\mathfrak{A}(\mathcal{S}), z')}{m}$$

Note 33. Then (4.5) and (4.4) imply

$$(4.6) \quad \lambda \left\| \tilde{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}^{(i)})) - \tilde{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right\| \leq \frac{\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i)}{m} + \frac{\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z') - \ell(\mathfrak{A}(\mathcal{S}), z')}{m}$$

Note 34. Now that we brought it in that form, we can use an additional assumption on the loss to bound it.

Assumption 35. *The loss function $\ell(\cdot, z)$ in (4.2) is convex for any $z \in \mathcal{Z}$ and ρ -Lipschitz.*

Note 36. Given ρ -Lipschitzness in Assumption 35, it is

$$(4.7) \quad \begin{aligned} \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) &\leq \rho \|\mathfrak{A}(\mathcal{S}^{(i)}) - \mathfrak{A}(\mathcal{S})\| \\ \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z') - \ell(\mathfrak{A}(\mathcal{S}), z') &\leq \rho \|\mathfrak{A}(\mathcal{S}^{(i)}) - \mathfrak{A}(\mathcal{S})\| \end{aligned}$$

and hence (4.6) yields

$$(4.8) \quad \|\mathfrak{A}(\mathcal{S}^{(i)}) - \mathfrak{A}(\mathcal{S})\| \leq 2 \frac{\rho}{\lambda m}$$

Note 37. Plugging (4.8) in (4.7) yields

$$\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \leq 2 \frac{\rho^2}{\lambda m}$$

Note 38. Using Theorem 19, we get an upper bound for the stability / over-fitting

$$\mathbb{E}_{\mathcal{S} \sim g} \left(R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{D}}(\mathfrak{A}(\mathcal{S})) \right) = \mathbb{E}_{\substack{\mathcal{S} \sim g, z' \sim g \\ i \sim U\{1, \dots, m\}}} \left(\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \right) \leq 2 \frac{\rho^2}{\lambda m}$$

Note 39. After this saga, the researcher could come to the conclusion that: In a Ridge regularization learning problem with loss function which is convex and ρ -Lipschitz, and the regularizer is $J(\cdot; \lambda) = \lambda \|\cdot\|^2$ with $\lambda > 0$, the learning rule trained against iid sample $\mathcal{S} = \{z_i\}_{i=1}^m$ is on-average-replace-one-stable with rate $\epsilon(m) = 2 \frac{\rho^2}{\lambda m}$; i.e.

$$(4.9) \quad \mathbb{E}_{\mathcal{S} \sim g} \left(R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{D}}(\mathfrak{A}(\mathcal{S})) \right) \leq 2 \frac{\rho^2}{\lambda m}$$

Note 40. From (4.9), we see that stability improves (and over-fitting decreases) as the shrinkage parameter λ increases.

4.3. Bounding the Risk (2.3).

Note 41. Given the bounds (4.3) and (4.9), the decomposition of the expected Risk in (2.3) yields that: In a Ridge regularization learning problem with loss function which is convex and ρ -Lipschitz, and the regularizer is $J(\cdot; \lambda) = \lambda \|\cdot\|^2$ with $\lambda > 0$, the learning rule trained against iid sample $\mathcal{S} = \{z_i\}_{i=1}^m$ has expected Risk bound

$$(4.10) \quad \mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S}))) \leq \underbrace{R_g(w') + \lambda \|w'\|_2^2}_{(I)} + \underbrace{2 \frac{\rho^2}{\lambda m}}_{(II)}; \quad \forall w' \in \mathcal{W}$$

Note 42. From 4.10, we see that there is a trade-off between Empirical Risk (I) and stability (II) with regards the regularization parameter λ . We wish to use the optimal $\lambda > 0$ corresponding to the smallest bound in (4.10); it has to both fit the training data well (but perhaps not too well) and be very stable to different training data from the same g (but perhaps not too stable)!

Definition 43. (Strongly Convex functions) A function f is λ -strongly convex function is for all w, u , and $\alpha \in (0, 1)$ we have

$$(4.11) \quad f(\alpha w + (1 - \alpha)u) \leq \alpha f(w) + (1 - \alpha)f(u) - \frac{\lambda}{2}\alpha(1 - \alpha)\|w - u\|^2$$

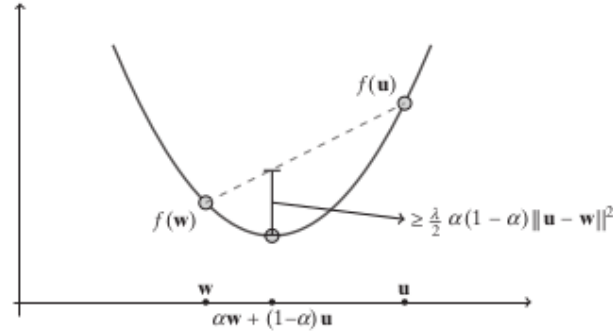


FIGURE 4.1. Strongly convex function

Proposition 44.

- (1) The function $f(w) = \lambda \|w\|_2^2$ is 2λ -strongly convex
- (2) If f is λ -strongly convex and g is convex then $f + g$ is λ -strongly convex

Lemma 45. If f is λ -strongly convex and u is a minimizer of f then for any w

$$f(w) - f(u) \geq \frac{\lambda}{2}\|w - u\|^2$$

Proof. Exercise 8 in the Exercise sheet. □