# Stability and Generalization in Pattern Recognition

Fateme Sheikholeslami

Dept. of ECE and DTC, University of Minnesota

*April, 2016*

1. O. Bousquet, and A. Elisseeff. "Stability and generalization," *Journal of Machine Learning Research,* vol. 2, pp. 499-526, Mar 2002.
2. J. Shawe-Taylor, and N. Cristianini, *Kernel methods for pattern analysis,* Cambridge University Press, 2004.

# Road map

- **Overview on pattern recognition**

- **Concentration of a fixed function**
  - McDiarmid inequality
  - Hoeffding's inequality

- **Concentration of a class of functions**
  - Capacity and regularization
    - Rademacher complexity

- **Stable algorithms**

- **Generalization bounds**
  - Polynomial bounds
  - Exponential bounds

- **Stability and generalization of regularized RKHS learning**

# Pattern recognition

❑ Choose a function from a class of functions which achieves a certain objective

- Often interested in $\min_{f \in \mathcal{F}} \mathbb{E}[f(\mathbf{x})]$

❑ Considerations

- $\{\mathbf{x}_i\}_{i=1}^N$ is drawn from an unknown pdf
- $\mathbb{E}[f(\mathbf{x})]$ is approximated by its empirical value $\hat{\mathbb{E}}[f(\mathbf{x})] := \frac{1}{N} \sum_{\mathbf{x}_i \in \mathcal{S}} f(\mathbf{x}_i)$ on a "training" set

$$\mathcal{S} = \{\mathbf{x}_1, ..., \mathbf{x}_N\}$$

❑ Performance

- What conclusion can be made about $\mathbb{E}[f(\mathbf{x})]$ based on its empirical measure?

# Concentration of a fixed function on a finite dataset

## Question 1

How concentrated a <u>fixed</u> function of a finite dataset $h(X_1, ..., X_M) \in \mathbb{R}$ is around its mean?

## McDiarmid's Inequality

Let $X_i \in \mathcal{A}$ denote independent random variables, and assume $h : \mathcal{A}^N \to \mathbb{R}$

If $\sup_{x_1, ..., x_N, \hat{x}_i \in \mathcal{A}} |h(x_1, ..., x_N) - h(x_1, ..., \hat{x}_i, ..., x_N)| \leq c_i, \quad 0 \leq c_i \leq N$

$$\implies \quad \forall \epsilon > 0 \quad \Pr\Big(h(x_1, ..., x_M) - \mathbb{E}[h(x_1, ..., x_N)] \geq \epsilon\Big) \leq \exp(\frac{-2\epsilon^2}{\sum_{i=1}^N c_i^2})$$

## Hoeffdings's Inequality

If $X_1, ..., X_N$ are independent r.v. satisfying $X_i \in [a_i, b_i]$, then for the r.v.

$S_N := \sum_{i=1}^N X_i$, we have $\quad \forall \epsilon > 0 \quad \Pr\Big(S_N - \mathbb{E}[S_N] \geq \epsilon\Big) \leq \exp(\frac{-2\epsilon^2}{\sum_{i=1}^N (b_i - a_i)^2})$

# Example: concentration of the sum of a finite dataset

- **Example: consider** $S_N(\mathbf{x}_1, ..., \mathbf{x}_N) = \frac{1}{N} \sum_{i=1}^{N} x_i = \hat{\mathbb{E}}[X]$ **where** $x_i \in [a, b]$.

$$|S_N(x_1, ..., x_N) - S_N(x_1, ..., \hat{x}_i, ..., x_N)| \leq (b-a)/N \implies \Pr\left(|\hat{\mathbb{E}}[X] - \mathbb{E}[X]| \geq \epsilon\right) \leq 2\exp\left(\frac{-2N\epsilon^2}{(b-a)^2}\right)$$

- **Example: consider the center of mass for the sample set** $\mathcal{S} = \{\mathbf{x}_1, ..., \mathbf{x}_N\}$

$$\phi_\mathcal{S} := \frac{1}{N} \sum_{i=1}^{N} \phi(\mathbf{x}_i)$$

➤ What can be concluded about its concentration?

   Measure of accuracy $g(\mathcal{S}) := \|\phi_\mathcal{S} - \mathbb{E}[\phi(\mathbf{x})]\|$

- Can we apply McDiarmid's Inequality? $\qquad \hat{\mathcal{S}} = \{\mathbf{x}_1, ..., \hat{\mathbf{x}}_i, ..., \mathbf{x}_N\}$

$$|g(\mathcal{S}) - g(\hat{\mathcal{S}})| = \|\phi_\mathcal{S} - \mathbb{E}_\mathbf{x}[\phi(\mathbf{x})]\| - \|\phi_{\hat{\mathcal{S}}} - \mathbb{E}_\mathbf{x}[\phi(\mathbf{x})]\| \leq \|\phi_\mathcal{S} - \phi_{\hat{\mathcal{S}}}\| = \frac{1}{N}\|\phi(\mathbf{x}_i) - \phi(\hat{\mathbf{x}}_i)\| \leq \frac{2R}{N}$$

   Yes! $\implies \Pr\left(g(\mathcal{S}) - \mathbb{E}_\mathcal{S}[g(\mathcal{S})] \geq \epsilon\right) \leq \exp\left(\frac{-2N\epsilon^2}{4R^2}\right) \qquad \boxed{\|\phi(\mathbf{x})\| \leq R}$

# Example: concentration of sample center of mass in feature space

- **Furthermore**

$$\mathbb{E}_{\mathcal{S}}[g(\mathcal{S})] = \mathbb{E}_{\mathcal{S}}[\|\phi_{\mathcal{S}} - \mathbb{E}[\phi(\mathbf{x})]\|] = \mathbb{E}_{\mathcal{S}}[\|\phi_{\mathcal{S}} - \mathbb{E}_{\tilde{\mathcal{S}}}[\phi_{\tilde{\mathcal{S}}}]\|]$$

$$= \mathbb{E}_{\mathcal{S}}[\|\mathbb{E}_{\tilde{\mathcal{S}}}[\phi_{\mathcal{S}} - \phi_{\tilde{\mathcal{S}}}]\|] \le \mathbb{E}_{\mathcal{S}\tilde{\mathcal{S}}}[\|\phi_{\mathcal{S}} - \phi_{\tilde{\mathcal{S}}}\|]$$

$$= \mathbb{E}_{\boldsymbol{\sigma}\mathcal{S}\tilde{\mathcal{S}}}\Big[\frac{1}{N}\Big\|\sum_{i=1}^{N}\sigma_i(\phi(\mathbf{x}_i) - \phi(\tilde{\mathbf{x}}_i))\Big\|\Big] \qquad \text{where } \sigma_i \in \{\pm 1\} \text{ with equal probability}$$

$$= \mathbb{E}_{\boldsymbol{\sigma}\mathcal{S}\tilde{\mathcal{S}}}\Big[\frac{1}{N}\Big\|\sum_{i=1}^{N}(\sigma_i\phi(\mathbf{x}_i) - \sigma_i\phi(\tilde{\mathbf{x}}_i))\Big\|\Big] \le 2\mathbb{E}_{\boldsymbol{\sigma}\mathcal{S}}\Big[\frac{1}{N}\Big\|\sum_{i=1}^{N}\sigma_i\phi(\mathbf{x}_i)\Big\|\Big]$$

$$= \frac{2}{N}\mathbb{E}_{\boldsymbol{\sigma}\mathcal{S}}\Big[\Big(\Big\langle\sum_{i=1}^{N}\sigma_i\phi(\mathbf{x}_i), \sum_{j=1}^{N}\sigma_j\phi(\mathbf{x}_j)\Big\rangle\Big)^{1/2}\Big]$$

$$\le \frac{2}{N}\Big(\mathbb{E}_{\boldsymbol{\sigma}\mathcal{S}}\Big[\sum_{i,j=1}^{N}\sigma_i\sigma_j\kappa(\mathbf{x}_i,\mathbf{x}_j)\Big]\Big)^{1/2} = \frac{2}{N}\Big(\mathbb{E}_{\boldsymbol{\sigma}\mathcal{S}}\Big[\sum_{i=1}^{N}\kappa(\mathbf{x}_i,\mathbf{x}_j)\Big]\Big)^{1/2} \le \frac{2R}{\sqrt{N}}$$
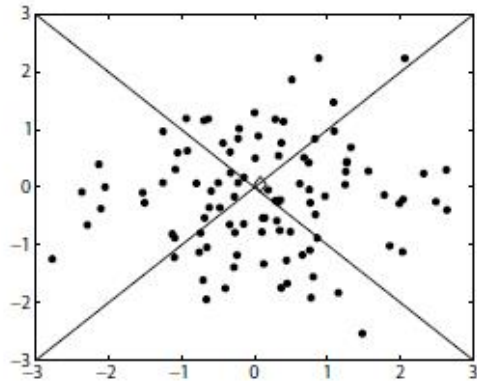
- **Previously, we had**
$$\Pr\Big(g(\mathcal{S}) - \mathbb{E}_{\mathcal{S}}[g(\mathcal{S})] \ge \epsilon\Big) \le \exp\Big(\frac{-2N\epsilon^2}{4R^2}\Big)$$

- **Setting** $\delta := \exp\Big(\dfrac{-2N\epsilon^2}{4R^2}\Big)$ **and after substitution, with probability at least** $1 - \delta$ **we have**
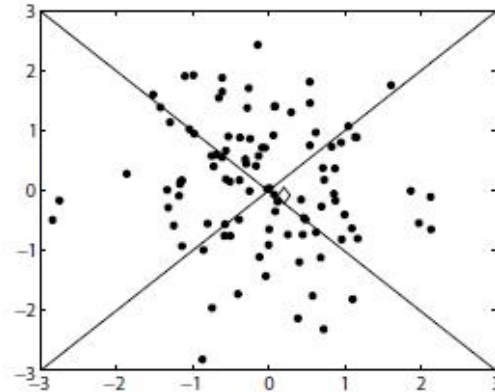
$$g(\mathcal{S}) \le \frac{R}{\sqrt{N}}\Big(2 + \sqrt{2\ln\frac{1}{\delta}}\Big)$$

# Example: concentration of sample mean

■ Sample mean of random draws of 2-dimensional Gaussian variables

The empirical centre of mass based on a random sample

. The empirical centre of mass based on a second random sample.

A random variable that depends (in a "smooth" way) on the influence of many independent variables (but not too much on any of them) is essentially constant.

Talagrand 1996.

# Capacity of a class

- Let us go back to pattern recognition

  - Find the function from a class of functions which achieves a certain objective

    - For instance  $\min_{f \in \mathcal{F}} \hat{\mathbb{E}}[f(x, y)]$

Question 2: How concentrated is empirical mean of the sought pattern to its true mean?

- Example

  - Find a function  $f \in P_{10}$  that maps creditcard numbers to the card holder's phone number

    $P_{10} :=$ Set of polynomials of degree 10

  - Given 10 training pairs, $\exists f \in P_{10}$ such that perfectly maps the training points!

  - Performance on unseen data? Arbitrarily poor!  $\Rightarrow$   overfitting!

- ❖  Performance of a pattern relies on  (1) concentration of the function value
                                        (2) Richness (capacity) of the class

# Rademacher Complexity

- Measures the capacity of a class by its ability to fit random data

- Let $\{\sigma_1, ..., \sigma_N\}$ be independent uniform $\{\pm 1\}$ -valued Rademacher r.v.

- For set $\mathcal{S} = \{\mathbf{x}_1, ..., \mathbf{x}_N\}$, define Empirical Rademacher complexity of class $\mathcal{F}$ as

$$\hat{R}_N(\mathcal{F}) = \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{f \in \mathcal{F}} \left| \frac{2}{N} \sum_{i=1}^{N} \sigma_i f(\mathbf{x}_i) \right| \, \middle| \, \mathbf{x}_1, ..., \mathbf{x}_N \right]$$

- Rademacher complexity of $\mathcal{F}$

$$R_N(\mathcal{F}) = \mathbb{E}_{\mathcal{S}}[\hat{R}_N(\mathcal{F})] = \mathbb{E}_{\mathcal{S}\boldsymbol{\sigma}} \left[ \sup_{f \in \mathcal{F}} \left| \frac{2}{N} \sum_{i=1}^{N} \sigma_i f(\mathbf{x}_i) \right| \right]$$

# Rademacher Complexity of kernel-based functions

■ Consider the class of linear functions in a kernel defined feature space

$$\mathcal{F} := \{f | f : \mathbf{x} \to \sum_{i=1}^{N} \alpha_i \kappa(\mathbf{x}_i, \mathbf{x}), \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} \leq B^2\}$$

■ Consider the class $\quad \mathcal{F}_B := \{f | f : \mathbf{x} \to \langle \mathbf{w}, \boldsymbol{\phi}(\mathbf{x}) \rangle, \|\mathbf{w}\| \leq B\}$ ⬭ Regularization!

If $\kappa : X \times X \to \mathbb{R}$ is a kernel, and $\mathcal{S} = \{\mathbf{x}_1, ..., \mathbf{x}_N\}$ is a sample of points, then the empirical Rademacher complexity of the class $\mathcal{F}_B$ satisfies

$$\hat{R}_N(\mathcal{F}_B) \leq \frac{2B}{N} \sqrt{\sum_{i=1}^{N} \kappa(\mathbf{x}_i, \mathbf{x}_i)} = \frac{2B}{N} \sqrt{\text{tr}(\mathbf{K})}$$

*Proof:*

$$\hat{R}_N(\mathcal{F}_B) = \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{f \in \mathcal{F}_B} \left| \frac{2}{N} \sum_{i=1}^{N} \sigma_i f(\mathbf{x}_i) \right| \right] = \frac{2}{N} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\|\mathbf{w}\| \leq B} \left| \left\langle \mathbf{w}, \sum_{j=1}^{N} \sigma_j \boldsymbol{\phi}(\mathbf{x}_j) \right\rangle \right| \right]$$

$$\leq \frac{2B}{N} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \left\| \sum_{i=1}^{N} \sigma_i \boldsymbol{\phi}(\mathbf{x}_i) \right\| \right] = \frac{2B}{N} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \left( \left\langle \sum_{i=1}^{N} \sigma_i \boldsymbol{\phi}(\mathbf{x}_i), \sum_{j=1}^{N} \sigma_j \boldsymbol{\phi}(\mathbf{x}_j) \right\rangle \right)^{1/2} \right]$$

$$= \frac{2B}{N} \left( \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sum_{i,j=1}^{N} \sigma_i \sigma_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \right] \right)^{1/2} = \frac{2B}{N} \sqrt{\sum_{i=1}^{N} \kappa(\mathbf{x}_i, \mathbf{x}_i)} \quad . \blacksquare$$

# Properties of Rademacher complexity

- Theorem: Let $\mathcal{F}_1, \mathcal{F}_2, ..., \mathcal{F}_m$ and $\mathcal{G}$ be classes of real functions. Then

  a) If $\mathcal{F} \subseteq \mathcal{G}$ then $\hat{R}_N(\mathcal{F}) \leq \hat{R}_N(\mathcal{G})$

  b) $\hat{R}_N(\mathcal{F}) = \hat{R}_N(\mathrm{conv}\mathcal{F})$

  c) For every $c \in \mathbb{R}, \hat{R}_N(c\mathcal{F}) = |c|\hat{R}_N(\mathcal{F})$

  d) If $\mathcal{A}: \mathbb{R} \to \mathbb{R}$ is $L$-Lipschitz and satisfies $\mathcal{A}(0) = 0$, then $\hat{R}_N(\mathcal{A} \circ \mathcal{F}) \leq 2L\hat{R}_N(\mathcal{F})$

  d) For any function $h, \hat{\mathcal{R}}_N(\mathcal{F} + h) \leq \hat{\mathcal{R}}_N(\mathcal{F}) + 2\sqrt{\hat{\mathbb{E}}[h^2]/N}$

  e) For any $1 \leq q < \infty$, let $\mathcal{L}_{\mathcal{F},h,q} = \{|f - h|^q | f \in \mathcal{F}\}$. If $\|f - h\|_\infty \leq 1$

  for every $f \in \mathcal{F}$, then $\hat{R}_N(\mathcal{L}_{\mathcal{F},h,q}) \leq 2q\left(\hat{R}_N(\mathcal{F}) + 2\sqrt{\hat{\mathbb{E}}[h^2]/N}\right)$.

  f) $\hat{R}_N(\sum_{i=1}^m \mathcal{F}_i) \leq \sum_{i=1}^m \hat{R}_N(\mathcal{F}_i)$

# Concentration of a class of functions

Fix $\delta \in (0, 1)$, and let $\mathcal{F} := \{f | f : X \to [0, 1]\}$. Let $\{\mathbf{x}_i\}_{i=1}^N$ be ind. drawn from distribution $\mathcal{D}$. Then, w.p. at least $1 - \delta$ over random draws of sample size $N$

$$\mathbb{E}_{\mathbf{x}}[f(\mathbf{x})] \leq \hat{\mathbb{E}}[f(\mathbf{x})] + R_N(\mathcal{F}) + \sqrt{\frac{\ln(2/\delta)}{2N}} \qquad \forall f \in \mathcal{F}$$

$$\leq \hat{\mathbb{E}}[f(\mathbf{x})] + \hat{R}_N(\mathcal{F}) + 3\sqrt{\frac{\ln(2/\delta)}{2N}}$$

- *Sketch of proof*

  For a fixed f:
  $$\mathbb{E}_{\mathbf{x}}[f(\mathbf{x})] \leq \hat{\mathbb{E}}_{\mathbf{x}}[f(\mathbf{x})] + \sup_{h \in \mathcal{F}} \left( \mathbb{E}_{\mathbf{x}} h - \hat{\mathbb{E}} h \right)$$

  Applying <u>McDiarmid's ineq</u>. on the second term (why?), w.p. at least $1 - \delta/2$

  $$\sup_{h \in \mathcal{F}} \left( \mathbb{E}_{\mathbf{x}} h - \hat{\mathbb{E}} h \right) \leq \mathbb{E}_{\mathcal{S}} \left[ \sup_{h \in \mathcal{F}} \left( \mathbb{E}_{\mathbf{x}} h - \hat{\mathbb{E}} h \right) \right] + \sqrt{\frac{\ln(2/\delta)}{2N}} \qquad \delta/2 := \exp\left( -2N\epsilon^2 \right)$$

  $$\mathbb{E}_{\mathbf{x}}[f(\mathbf{x})] \leq \hat{\mathbb{E}}[f(\mathbf{x})] + \underbrace{\mathbb{E}_{\mathcal{S}} \left[ \sup_{h \in \mathcal{F}} \left( \mathbb{E}_{\mathbf{x}} h - \hat{\mathbb{E}} h \right) \right]} + \sqrt{\frac{\ln(2/\delta)}{2N}}$$
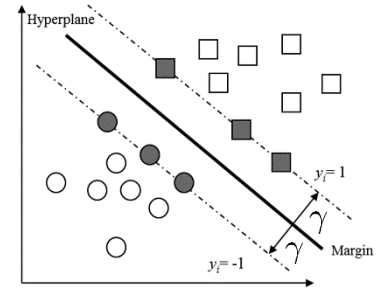
  Can be shown:
  $$\leq R_N(\mathcal{F}) \leq \hat{R}_N(\mathcal{F}) + 2\sqrt{\frac{\ln(2/\delta)}{2N}}$$

  w.p. at least $1 - \delta/2$ . ∎

# Concentration of kernel-based SVM classifier

- Given a function $g(\mathbf{x})$, a dataset $\mathcal{S} = \{(\mathbf{x}_1, y_1), ..., (\mathbf{x}_N, y_N)\}$, a desired margin $\gamma$

- Define slack variable

$$\xi_i := (\gamma - y_i g(\mathbf{x}_i))_+ = \begin{cases} \gamma - y_i g(\mathbf{x}_i) & \gamma - y_i g(\mathbf{x}_i) > 0 \\ 0 & \text{otherwise} \end{cases}$$



## Theorem

Fix $\gamma > 0$, and let $\mathcal{F} := \{f | f(\mathbf{x}, y) = -y g(\mathbf{x}), \ g(\mathbf{x}) = \langle \phi(\mathbf{x}), \mathbf{w} \rangle, \ \|\mathbf{w}\|_{\mathcal{H}} \le 1\}$.

Let $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ be ind. drawn from distribution $\mathcal{D}$. Then, w.p. at least $1 - \delta$ over ind. draws of sample size $N$ we have

$$\mathbb{P}\left(y \ne \operatorname{sign}\left(g(\mathbf{x})\right)\right) \le \frac{1}{N\gamma} \sum_{i=1}^N \xi_i + \frac{4}{N\gamma} \sqrt{\operatorname{tr}(\mathbf{K})} + 3\sqrt{\frac{\ln(2/\delta)}{2N}}$$

# Sketch of proof

- **Define** $\mathcal{P}(a) = \begin{cases} 1, & \text{if} \quad a > 0; \\ 1 + a/\gamma, & \text{if} \quad -\gamma \leq a \leq 0; \\ 0, & \text{otherwise} \end{cases}$ and $\mathcal{H}(a) = \begin{cases} 1, & \text{if} \quad a > 0; \\ 0, & \text{otherwise} \end{cases}$

- **Since** $\mathcal{P}(a)$ dominates $\mathcal{H}(a)$, we have

$$\mathbb{E}_{\mathbf{x}}[\mathcal{H}(f(\mathbf{x}, y)) - 1] \leq \mathbb{E}_{\mathbf{x}}[\mathcal{P}(f(\mathbf{x}, y)) - 1]$$

$$\leq \hat{\mathbb{E}}_{\mathbf{x}}[\mathcal{P}(f(\mathbf{x}, y)) - 1] + \hat{R}_\ell((\mathcal{P} - 1)o\mathcal{F}) + 3\sqrt{\frac{\ln(2/\delta)}{2N}}$$

- $(\mathcal{P} - 1)(a)$ is Lipschitz continuous with $L = 1/\gamma$ and $(\mathcal{P} - 1)(0) = 0$

$$\mathbb{P}\left(y \neq \text{sign}\left(g(\mathbf{x})\right)\right) = \mathbb{E}_{\mathbf{x}}[\mathcal{H}(f(\mathbf{x}, y))] \leq \underbrace{\hat{\mathbb{E}}_{\mathbf{x}}[\mathcal{P}(f(\mathbf{x}, y))]}_{\leq \xi_i/\gamma} + \underbrace{2\hat{R}_\ell(\mathcal{F})/\gamma}_{\hat{R}_\ell(\mathcal{G}) \leq \frac{2}{N}\sqrt{\text{tr}(\mathbf{K})}} + 3\sqrt{\frac{\ln(2/\delta)}{2N}}$$

.■

14

# Algorithms

- Assume data is given as $\mathcal{S} = \{\mathbf{z}_1, ..., \mathbf{z}_N\}$ where $\mathbf{z}_i := (\mathbf{x}_i, y_i)$

- Loss functions are usually of interest $f(\mathbf{z}) = \ell(g, \mathbf{z})$ where, e.g. $g(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} \subset \mathcal{G}$

- Define Algorithm $A : \mathcal{Z}^N \to \mathcal{G}$ that maps dataset $\mathcal{S}$ into a function $A_\mathcal{S} \subset \mathcal{G} : \mathcal{X} \to \mathcal{Y}$

  e.g., $A_\mathcal{S} = \arg\min_{g \in \mathcal{G}} \frac{1}{N} \sum_{i=1}^N \ell(g, \mathbf{z}_i) + \lambda \|g\|_\mathcal{H}^2$

  $$\mathcal{S} \to \boxed{\begin{array}{c}\text{algorithm} \\ A\end{array}} \xrightarrow{A_\mathcal{S}}$$

- Define "Risk" functions

  $$R(A, \mathcal{S}) := \mathbb{E}_\mathbf{z}[\ell(A_\mathcal{S}, \mathbf{z})]$$
  $$R_{emp}(A, \mathcal{S}) := \sum_{i=1}^N \ell(A_\mathcal{S}, \mathbf{z}_i)/N$$
  $$R_{loo}(A, \mathcal{S}) := \sum_{i=1}^N \ell(A_{\mathcal{S} \setminus i}, \mathbf{z}_i)/N$$

- So far, we have studied

  $$\mathbb{P}\left( \sup_{f \in \mathcal{F}} \left| \mathbb{E}[f] - \hat{\mathbb{E}}[f] \right| > \epsilon \right)$$

- Different approaches study

  $$\mathbb{P}\left( \left| R(A, \mathcal{S}) - R_{emp}(A, \mathcal{S}) \right| > \epsilon \right)$$

  while assuming a notion of "stability" for the algorithm $A$.

# Algorithm stability

D1) Algorithm $A$ has *pointwise hypothesis stability* $\beta_N$ if

$$\forall i \in \{1, ..., N\}, \ \mathbb{E}_{\mathcal{S}}\Big[\big|\ell(A_{\mathcal{S}}, \mathbf{z}_i) - \ell(A_{\mathcal{S} \setminus i}, \mathbf{z}_i)\big|\Big] \leq \beta_N$$

D2) Algorithm $A$ has *hypothesis stability* $\beta_N$ if

$$\forall i \in \{1, ..., N\}, \ \mathbb{E}_{\mathcal{S}, \mathbf{z}}\Big[\big|\ell(A_{\mathcal{S}}, \mathbf{z}) - \ell(A_{\mathcal{S} \setminus i}, \mathbf{z})\big|\Big] \leq \beta_N$$

D3) Algorithm $A$ has *uniform stability* $\beta_N$ if

$$\forall \mathcal{S} \in \mathcal{Z}^N, \ \forall i \in \{1, ..., N\}, \ \max_{\mathbf{z} \in \mathrm{supp}(\mathcal{D})} |\ell(A_{\mathcal{S}}, \mathbf{z}) - \ell(A_{\mathcal{S} \setminus i}, \mathbf{z})| \leq \beta_N$$

❖ Algorithm $A$ is considered *stable* if $\beta_N$ decreases as $1/N$.

# Polynomial bounds with hypothesis stability

Theorem

For Algorithm $A$ with hypothesis stability $\beta_1$ and pointwise stability $\beta_2$ w.r.t. a loss function $0 \leq \ell(A_{\mathcal{S}}, \mathbf{z}) \leq M$, w.p. at least $1 - \delta$ we have

$$R(A, \mathcal{S}) \leq R_{emp}(A, \mathcal{S}) + \sqrt{\frac{M^2 + 12MN\beta_2}{2N\delta}}$$

and

$$R(A, \mathcal{S}) \leq R_{loo}(A, \mathcal{S}) + \sqrt{\frac{M^2 + 6MN\beta_1}{2N\delta}}$$

# Exponential bounds with uniform stability

■ Consider a regression task

**Theorem**

For Algorithm $A$ with uniform stability $\beta$ w.r.t. a loss function $0 \leq \ell(A_{\mathcal{S}}, \mathbf{z}) \leq M$,

w.p. at least $1 - \delta$ we have

$$R(A, \mathcal{S}) \leq R_{emp}(A, \mathcal{S}) + 2\beta + (4N\beta + M)\sqrt{\frac{\ln(1/\delta)}{2N}}$$

and

$$R(A, \mathcal{S}) \leq R_{loo}(A, \mathcal{S}) + \beta + (4N\beta + M)\sqrt{\frac{\ln(1/\delta)}{2N}}$$

❖ The bound is tight if $\beta$ scales as $1/N$.

❖ Specialized bounds for classification task is also available.

Question

Are commonly-used learning algorithms stable?

# Uniform stability of regularized RKHS learning

■ Consider the class of linear functions in a kernel defined feature space $\mathcal{G}$

Definition: Loss function $\ell(g, \mathbf{z})$ on $\mathcal{G} \times \mathcal{Y}$ is $\sigma$-admissible w.r.t. $\mathcal{G}$ if the associated cost $\ell(g, \mathbf{z}) = c(g(\mathbf{x}), y)$ is convex w.r.t. its first argument, and

$$\forall y_1, y_2 \in \mathcal{D}, \forall y' \in \mathcal{Y}, |c(y_1, y') - c(y_2, y')| \leq \sigma |y_1 - y_2|$$

where $\mathcal{D} = \{y | \exists g \in \mathcal{G}, \exists \mathbf{x} \in \mathcal{X} : g(\mathbf{x}) = y\}$ .

## Theorem

Assume for given kernel $\kappa(\mathbf{x}, \mathbf{x}) \leq \kappa^2 < \infty$, and let loss $\ell(g, \mathbf{z})$ be $\sigma$-admissible w.r.t. $\mathcal{G}$. Then the learning algorithm $A$ defined by

$$A_{\mathcal{S}} = \arg\min_{g \in \mathcal{G}} \frac{1}{N} \sum_{i=1}^{N} \ell(g, \mathbf{z}_i) + \lambda \|g\|_{\mathcal{H}}^2$$

has uniform stability $\beta \leq \dfrac{\sigma^2 \kappa^2}{2\lambda N}$ .

# Examples on regularized RKHS learning

❖ Regularized RKHS learning

$$A_{\mathcal{S}} = \arg\min_{g \in \mathcal{G}} \frac{1}{N} \sum_{i=1}^{N} \ell(g, \mathbf{z}_i) + \lambda \|g\|_{\mathcal{H}}^2$$

❖ Ex1) Bounded SVM regression

- $\ell(g, \mathbf{z}) = |g(\mathbf{x}) - y|_\epsilon = \begin{cases} 0 & \text{if } |g(\mathbf{x}) - y| \leq \epsilon \\ |g(\mathbf{x}) - y| - \epsilon & \text{otherwise} \end{cases}$ and $\mathcal{Y} = [0, B]$.

- $\ell(g, \mathbf{z})$ is 1-admissible $\Rightarrow$ $\beta \leq \dfrac{\kappa^2}{2\lambda N}$ $\Rightarrow$ $R \leq R_{emp} + \dfrac{\kappa^2}{\lambda N} + \left( \dfrac{2\kappa^2}{\lambda} + \kappa\sqrt{\dfrac{B}{\lambda}} \right) \sqrt{\dfrac{\ln(1/\delta)}{2N}}$

❖ Ex2) Regularized least squares

- $\ell(g, \mathbf{z}) = (g(\mathbf{x}) - y)^2$ and $\mathcal{Y} = [0, B]$.

- $\ell(g, \mathbf{z})$ is $2B$-admissible $\Rightarrow$ $\beta \leq \dfrac{2B^2\kappa^2}{\lambda N}$ $\Rightarrow$ $R \leq R_{emp} + \dfrac{4\kappa^2\beta^2}{\lambda N} + \left( \dfrac{8\kappa^2 B^2}{\lambda} + 2B \right) \sqrt{\dfrac{\ln(1/\delta)}{2N}}$

# Summary

- **Concentration of a fixed function**
  - McDiarmid inequality
  - Hoeffding's inequality

- **Concentration of a class of functions**
  - Capacity and regularization
    - Rademachar complexity

- **Algorithm stability**

- **Generalization bounds**
  - Polynomial bounds
  - Exponential bounds

- **Regularized RKHS learning**

*Thank You!*