

Statistical Foundations of Learning

Debarghya Ghoshdastidar

School of Computation, Information and Technology
Technical University of Munich

k-Nearest Neighbour Classification

Outline

- k -nearest neighbour classification
 - Generalisation error of k -NN for finite k as $m \rightarrow \infty$
- Consistency / Universal consistency: Asymptotically achieving Bayes risk
- Plug-in classifiers
 - Stone's theorem: Universal consistency of plug-in classifiers
 - Universal consistency of k -NN rule
 - Proof of Stone's theorem

Nearest neighbour rule

- Assume $\mathcal{X} \subset \mathbb{R}^p$ and we use Euclidean distance $\|x - x'\|$
- Given $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$
- Nearest neighbour and k -nearest neighbour rules:
 - For test data $x \in \mathcal{X}$, sort x_1, \dots, x_m according to $\|x - x_i\|$
 - $\pi_k(x) \in [m]$ such that $x_{\pi_k(x)}$ is k -th nearest neighbour of x

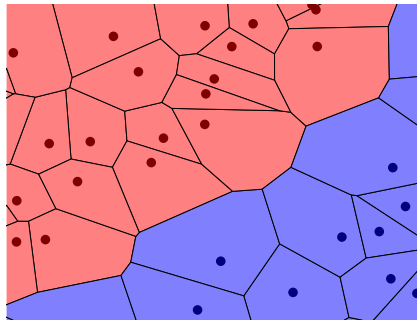
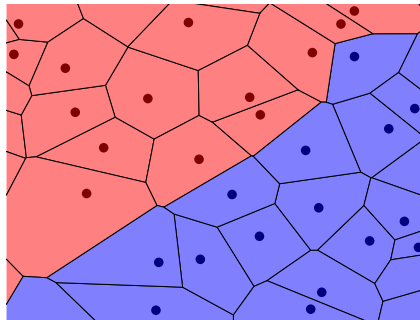
$$\|x - x_{\pi_1(x)}\| \leq \|x - x_{\pi_2(x)}\| \leq \dots \leq \|x - x_{\pi_m(x)}\|$$

- $h_S^{NN}(x) = y_{\pi_1(x)}$ and $h_S^{kNN}(x) = \text{majority vote of } y_{\pi_1(x)}, \dots, y_{\pi_k(x)}$

There can be ties (not discussed here)

Nearest neighbour rule

- Finite m : Decision boundary depends significantly on S
- Large m : Can learn very complex decision boundaries (more complex for $k > 1$)



Recap: Bayes risk

Learning problem characterised by:

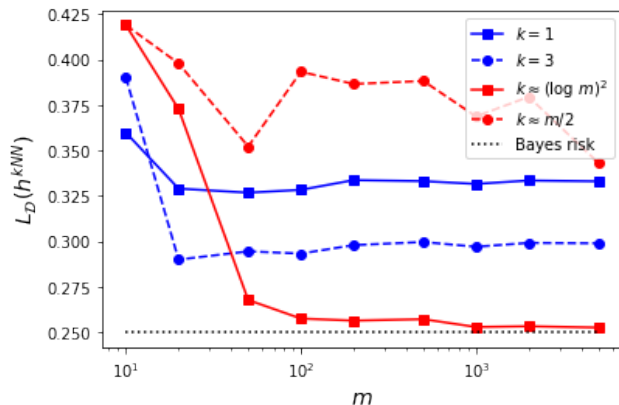
$$\mathcal{D} = \underbrace{\mathcal{D}_{\mathcal{X}}}_{\text{marginal of features}} \times \underbrace{\mathbb{P}_{\mathcal{Y}|\mathcal{X}}(y|x)}_{\text{conditional probability of label}}$$

- $\eta(x) = \mathbb{P}_{\mathcal{Y}|\mathcal{X}}(y = 1|x)$
- Bayes risk $L_{\mathcal{D}}^* = \mathbb{E}_{x \sim \mathcal{D}_{\mathcal{X}}} [\min\{\eta(x), 1 - \eta(x)\}]$
- Bayes risk is smallest possible risk / generalisation error for a learning problem
- Bayes risk achieved by Bayes classifier (needs knowledge of η)

Example: Performance of k -NN rule

- Predicting software crash:

$$\mathcal{D}_{\mathcal{X}} = \text{Uniform}[0, 1] \quad \text{and} \quad \eta(x) = |1 - 2x|$$



- Test error of k -NN for different values of m and k
 - Test data has 5000 samples
 - Errors averaged over 25 trials
- We will try to explain these results mathematically

Expected generalisation error for NN rule

Theorem kNN.1 (Asymptotic expected risk of NN rule)

- Define $L_{\mathcal{D}}^{NN} = \mathbb{E}_{x \sim \mathcal{D}_X} [2\eta(x)(1 - \eta(x))]$
As $m \rightarrow \infty$, $\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S^{NN})] \rightarrow L_{\mathcal{D}}^{NN}$
- Comparison with Bayes risk: $L_{\mathcal{D}}^* \leq L_{\mathcal{D}}^{NN} \leq 2L_{\mathcal{D}}^*$
- Easy to verify 2nd statement (exercise)
- We will prove the 1st statement (under some assumptions)

Expected generalisation error of k -NN

Theorem kNN.2 (Asymptotic expected risk of k -NN rule)

- Assume k is fixed

- Define $L_{\mathcal{D}}^{kNN} = \lim_{m \rightarrow \infty} \mathbb{E}_{S \sim \mathcal{D}^m} \left[L_{\mathcal{D}} \left(h_S^{kNN} \right) \right]$ *(limit exists)*

$$L_{\mathcal{D}}^* \leq L_{\mathcal{D}}^{kNN} \leq \left(1 + \frac{2}{\sqrt{k}} \right) L_{\mathcal{D}}^*$$

- Proof skipped. Similar to proof for $k = 1$ (more involved)
- Result suggests that we need $k \rightarrow \infty$ to get close to Bayes risk

Intuition why NN (or k -NN) works

- Let $x_1, x_2, \dots, x_m \sim_{iid} \mathcal{D}_{\mathcal{X}}$. Consider some x .
- Intuition: For large m , $x_{\pi_1(x)}, \dots, x_{\pi_k(x)}$ are arbitrarily close to x
Hence, the label of x is likely to be same as $y_{\pi_1(x)}, \dots, y_{\pi_k(x)}$
- Bit more formal:
If $\frac{k}{m} \rightarrow 0$ as $m \rightarrow \infty$, then $x_{\pi_1(x)}, \dots, x_{\pi_k(x)} \rightarrow x$ in probability
 - Convergence in probability: A sequence of random variables z_1, z_2, \dots is said to converge to random variable x in probability if, for every $\epsilon > 0$, $\lim_{m \rightarrow \infty} \mathbb{P}(\|z_m - x\| > \epsilon) = 0$.

Formal proof of Theorem kNN.1 (under assumptions)

- Goal: Show $\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S^{\text{NN}})] \rightarrow \mathbb{E}_{x \sim \mathcal{D}_{\mathcal{X}}} [2\eta(x)(1 - \eta(x))]$ as $m \rightarrow \infty$
- $\ell = 0\text{-}1$ loss $\implies L_{\mathcal{D}}(h_S^{\text{NN}}) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathbf{1}\{y \neq y_{\pi_1(x)}\}]$
- $S = \{(x_1, y_1), \dots, (x_m, y_m)\} \sim \mathcal{D}^m$, and test data $(x, y) \sim \mathcal{D}$
 - View as $x, x_1, \dots, x_m \sim_{iid} \mathcal{D}_{\mathcal{X}}$ generated first
 - Then labels generated according to $\eta(\cdot)$

$$\begin{aligned}\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S^{\text{NN}})] &= \mathbb{E}_{S, (x,y) \sim \mathcal{D}^{m+1}} [\mathbf{1}\{y \neq y_{\pi_1(x)}\}] \\ &= \mathbb{E}_{x, x_1, \dots, x_m} [\mathbb{E}_{y, y_1, \dots, y_m} [\mathbf{1}\{y \neq y_{\pi_1(x)}\} \mid x, x_1, \dots, x_m]]\end{aligned}$$

Proof of Theorem kNN.1 (continued)

- Conditioned on $x, x_{\pi_1(x)}$, labels y and $y_{\pi_1(x)}$ are independent

$$\mathbb{E}_{y, y_1, \dots, y_m} [\mathbf{1} \{y \neq y_{\pi_1(x)}\} \mid x, x_1, \dots, x_m] = \eta(x)(1 - \eta(x_{\pi_1(x)})) + (1 - \eta(x))\eta(x_{\pi_1(x)})$$

- Need to show: As $m \rightarrow \infty$,

$$\mathbb{E}_{x, x_1, \dots, x_m} [\eta(x)(1 - \eta(x_{\pi_1(x)})) + (1 - \eta(x))\eta(x_{\pi_1(x)})] - \mathbb{E}_x [2\eta(x)(1 - \eta(x))] \rightarrow 0$$

Equivalently,

$$\left| \mathbb{E}_{x, x_1, \dots, x_m} [\eta(x)(1 - \eta(x_{\pi_1(x)})) + (1 - \eta(x))\eta(x_{\pi_1(x)}) - 2\eta(x)(1 - \eta(x))] \right| \rightarrow 0$$

- Jensen's inequality: If $f(z)$ is a convex function, then $f(\mathbb{E}[z]) \leq \mathbb{E}[f(z)]$

Example: $|\mathbb{E}[z]| \leq \mathbb{E}[|z|]$

Proof of Theorem kNN.1 (continued)

$$\begin{aligned} & \left| \eta(x)(1 - \eta(x_{\pi_1(x)})) + (1 - \eta(x))\eta(x_{\pi_1(x)}) - 2\eta(x)(1 - \eta(x)) \right| \\ &= \left| (1 - 2\eta(x))(\eta(x_{\pi_1(x)}) - \eta(x)) \right| \\ &= |1 - 2\eta(x)| \cdot |\eta(x_{\pi_1(x)}) - \eta(x)| \\ &\leq |\eta(x_{\pi_1(x)}) - \eta(x)| \quad \text{since } |1 - 2\eta(x)| \leq 1 \text{ for all } x \end{aligned}$$

- Hence, suffices to show that: $\mathbb{E}_{x, x_1, \dots, x_m} [|\eta(x_{\pi_1(x)}) - \eta(x)|] \rightarrow 0$

Proof of Theorem kNN.1 (adding assumptions on $\mathcal{D}_{\mathcal{X}}, \eta$)

(A1) $\eta : \mathcal{X} \rightarrow [0, 1]$ is uniformly continuous

- Uniform continuity:

η is uniformly continuous if for every $\delta > 0$, there exists $\epsilon_{\delta} > 0$ such that for every $x, x' \in \mathcal{X}$ with $\|x - x'\| \leq \epsilon_{\delta}$, $|\eta(x) - \eta(x')| \leq \delta$

- Equivalently, there exists $\epsilon_{\delta} > 0$ such that $|\eta(x) - \eta(x')| > \delta \implies \|x - x'\| > \epsilon_{\delta}$

(A2) $\text{support}(\mathcal{D}_{\mathcal{X}}) = \mathcal{X}$

- Define $\mathcal{D}_{\mathcal{X}}(x; \epsilon) = \mathbb{P}_{x' \sim \mathcal{D}_{\mathcal{X}}}(\|x' - x\| \leq \epsilon)$, probability mass in ϵ -neighbourhood of x

- $\text{support}(\mathcal{D}_{\mathcal{X}}) = \text{set of all } x \in \mathcal{X} \text{ for which } \mathcal{D}_{\mathcal{X}}(x; \epsilon) > 0 \text{ for every } \epsilon > 0$

- These assumptions are not necessary, but lead to a simpler proof

Proof of Theorem kNN.1 (continued)

- Choose any $\delta \in (0, 1)$. We can write

$$\begin{aligned}\mathbb{E} [|\eta(x_{\pi_1(x)}) - \eta(x)|] &= \mathbb{E} [|\eta(x_{\pi_1(x)}) - \eta(x)| \cdot \mathbf{1} \{|\eta(x_{\pi_1(x)}) - \eta(x)| > \delta\}] \\ &\quad + \mathbb{E} [|\eta(x_{\pi_1(x)}) - \eta(x)| \cdot \mathbf{1} \{|\eta(x_{\pi_1(x)}) - \eta(x)| \leq \delta\}] \\ &\leq \mathbb{P} (|\eta(x_{\pi_1(x)}) - \eta(x)| > \delta) + \delta\end{aligned}$$

- From uniform continuity of η : $|\eta(x_{\pi_1(x)}) - \eta(x)| > \delta \implies \|x_{\pi_1(x)} - x\| > \epsilon_\delta$

$$\begin{aligned}\mathbb{P}_{x, x_1, \dots, x_m} (|\eta(x_{\pi_1(x)}) - \eta(x)| > \delta) &\leq \mathbb{P}_{x, x_1, \dots, x_m} (\|x_{\pi_1(x)} - x\| > \epsilon_\delta) \\ &= \mathbb{P}_{x, x_1, \dots, x_m} \left(\min_{i \in \{1, \dots, m\}} \|x_i - x\| > \epsilon_\delta \right) \\ &= \mathbb{E}_x \left[\mathbb{P}_{x_1, \dots, x_m} \left(\min_{i \in \{1, \dots, m\}} \|x_i - x\| > \epsilon_\delta \mid x \right) \right]\end{aligned}$$

Proof of Theorem kNN.1 (continued)

- Recall x_1, \dots, x_m are independent

$$\begin{aligned}\mathbb{P}_{x_1, \dots, x_m} \left(\min_{i \in \{1, \dots, m\}} \|x_i - x\| > \epsilon_\delta \mid x \right) &= \prod_{i=1}^m \mathbb{P}_{x_i} \left(\|x_i - x\| > \epsilon_\delta \mid x \right) \\ &= \prod_{i=1}^m \left(1 - \mathbb{P}_{x_i} \left(\|x_i - x\| \leq \epsilon_\delta \mid x \right) \right) \\ &= (1 - \mathcal{D}_{\mathcal{X}}(x; \epsilon_\delta))^m\end{aligned}$$

- For every $x \in \text{support}(\mathcal{D}_{\mathcal{X}})$, $\mathcal{D}_{\mathcal{X}}(x; \epsilon_\delta) > 0$ and so

$$(1 - \mathcal{D}_{\mathcal{X}}(x; \epsilon_\delta))^m \rightarrow 0 \text{ as } m \rightarrow \infty$$

- By assumption (A2), this is true for every $x \in \mathcal{X}$. So for every $\epsilon_\delta > 0$

$$\mathbb{E}_x \left[\mathbb{P}_{x_1, \dots, x_m} \left(\min_{i \in \{1, \dots, m\}} \|x_i - x\| > \epsilon_\delta \mid x \right) \right] \rightarrow 0 \text{ as } m \rightarrow \infty$$

Proof of Theorem kNN.1 (conclusion)

Combining everything

$$\begin{aligned} & \left| \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}} (h_S^{\text{NN}})] - L_{\mathcal{D}}^* \right| \\ &= \left| \mathbb{E}_{x, x_1, \dots, x_m} [\eta(x)(1 - \eta(x_{\pi_1(x)})) + (1 - \eta(x))\eta(x_{\pi_1(x)}) - 2\eta(x)(1 - \eta(x))] \right| \\ &\leq \mathbb{E}_{x, x_1, \dots, x_m} [|\eta(x_{\pi_1(x)}) - \eta(x)|] \\ &\leq \underbrace{\mathbb{P}_{x, x_1, \dots, x_m} (|\eta(x_{\pi_1(x)}) - \eta(x)| > \delta)}_{\rightarrow 0 \text{ as } m \rightarrow \infty} + \delta \quad \text{for any chosen } \delta > 0 \end{aligned}$$

For any $\delta > 0$, $\lim_{m \rightarrow \infty} \left| \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}} (h_S^{\text{NN}})] - L_{\mathcal{D}}^* \right| \leq \delta$

Hence, limit must be zero. Concludes the proof.

Proof idea of Theorem kNN.1, kNN.2 (without assumptions)

- Use above proof till you get $\mathbb{E}_{x, x_1, \dots, x_m} [|\eta(x_{\pi_1(x)}) - \eta(x)|]$
- Then one applies 2nd statement below with $k = 1$, which has no assumption on $\eta, \mathcal{D}_{\mathcal{X}}$
- Proof skipped. If interested, see Lemmas 5.3-5.4 in Devroye's book

Lemma kNN.3 (Convergence of function computed on k-NN)

Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be an integrable function. If $\frac{k}{m} \rightarrow 0$, then

$$(i) \quad \frac{1}{k} \sum_{i=1}^k \mathbb{E} [|f(x_{\pi_i(x)})|] \leq \left(\left(1 + \frac{2}{\sqrt{2-\sqrt{3}}} \right)^p - 1 \right) \mathbb{E}[|f(x)|] \quad (\text{Stone's lemma})$$

$$(ii) \quad \mathbb{E}_{x, x_1, \dots, x_m \sim \mathcal{D}_{\mathcal{X}}^{m+1}} \left[\frac{1}{k} \sum_{i=1}^k |f(x_{\pi_i(x)}) - f(x)| \right] \rightarrow 0 \quad \text{as } m \rightarrow \infty$$

Consistency and Universal consistency

- \mathcal{D} = distribution on $\mathcal{X} \times \mathcal{Y}$
- h_S = predictor learned by algorithm \mathcal{A} given sample $S \sim \mathcal{D}^m$
- \mathcal{A} is **consistent** with respect to \mathcal{D} and specified loss if

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)] \rightarrow L_{\mathcal{D}}^* \quad \text{as } m \rightarrow \infty$$

- \mathcal{A} is **universally consistent** if it is consistent for every \mathcal{D}

Practical approach to Bayes classification

- Bayes binary classifier

- $h^*(x) = \mathbf{1} \{ \eta(x) \geq \frac{1}{2} \}$

or, $h^*(x) = \text{sign} \left(\eta(x) - \frac{1}{2} \right) \in \{-1, +1\}$

- Main challenge: $\eta(\cdot)$ not known

- Plug-in classifier:

- $\hat{\eta}(\cdot)$ = estimate $\eta(\cdot)$ from labelled examples S

- Predictor

$$\hat{h}(x) = \begin{cases} 1 & \text{if } \hat{\eta}(x) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases} \quad \text{OR} \quad \hat{h}(x) = \begin{cases} 1 & \text{if } \hat{\eta}(x) \geq \frac{1}{2} \\ -1 & \text{otherwise} \end{cases}$$

Example of plug-in classifier: Naïve Bayes

- From Bayes theorem:

$$\eta(x) = \mathbb{P}(y = 1|x) = \frac{\mathbb{P}(y = 1)\mathbb{P}(x|y = 1)}{\mathbb{P}(x)} \quad 1 - \eta(x) = \frac{\mathbb{P}(y = 0)\mathbb{P}(x|y = 0)}{\mathbb{P}(x)}$$

- Rewriting Bayes classifier: $h^*(x) = \mathbf{1} \{ \mathbb{P}(y = 1)\mathbb{P}(x|y = 1) > \mathbb{P}(y = 0)\mathbb{P}(x|y = 0) \}$
- A plug-in classifier: $\hat{h}(x) = \mathbf{1} \{ \hat{p}_1 \hat{b}_1(x) > \hat{p}_0 \hat{b}_0(x) \}$
 - Easy to estimate $\mathbb{P}(y = i)$; \hat{p}_i = fraction of training data with label- i
 - Difficult to estimate class-conditional probability/density $\mathbb{P}(x|y)$ if x is high dimensional
 - Naïve Bayes: For $x = (x^{(1)}, \dots, x^{(p)}) \in \mathbb{R}^p$, assume $\mathbb{P}(x|y = i) = \underbrace{\prod_{j=1}^p \mathbb{P}(x^{(j)}|y = i)}_{\hat{b}_i(x) \text{ estimates this}}$

Example of plug-in classifier: NN rule

- $\hat{\eta}(\cdot)$ as weighted average:
 - Given $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$
 - For test data x , define weights $w_1(x), \dots, w_m(x) \in [0, 1]$ with $\sum_{i=1}^m w_i(x) = 1$

$$\hat{\eta}(x) = \sum_{i=1}^m \mathbf{1}\{y_i = 1\} w_i(x)$$

- NN rule: $w_i(x) = 1$ for $i = \pi_1(x)$, and 0 otherwise

$$\hat{\eta}(x) = \mathbf{1}\{y_{\pi_1(x)} = 1\} \implies \hat{h}(x) = y_{\pi_1(x)}$$

- Questions: What are $w_1(\cdot), \dots, w_m(\cdot)$ for kNN?

Can we write Naïve Bayes in terms of weighted average?

Universal consistency of plug-in classifiers

Theorem kNN.4 (Stone's consistency theorem)

\hat{h} is universally consistent if weights for estimating $\hat{\eta}$ satisfy:

(i) $\exists c$ such that, for every non-negative integrable function f with $\mathbb{E}_{x \sim \mathcal{D}_{\mathcal{X}}}[f(x)] < \infty$,

$$\mathbb{E}_{x, x_1, \dots, x_m \sim \mathcal{D}^{m+1}} \left[\sum_{i=1}^m w_i(x) \cdot f(x_i) \right] \leq c \mathbb{E}_{x \sim \mathcal{D}_{\mathcal{X}}}[f(x)]$$

(ii) For all $a > 0$, $\lim_{m \rightarrow \infty} \mathbb{E}_{x, x_1, \dots, x_m \sim \mathcal{D}^{m+1}} \left[\sum_{i=1}^m w_i(x) \cdot \mathbf{1} \{ \|x_i - x\| > a \} \right] = 0$

(iii) $\lim_{m \rightarrow \infty} \mathbb{E}_{x, x_1, \dots, x_m \sim \mathcal{D}^{m+1}} \left[\max_{i \in [m]} w_i(x) \right] = 0$

k -nearest neighbour rule

- Assume $\mathcal{X} \subset \mathbb{R}^p$ and we use Euclidean distance $\|x - x'\|$
- Given $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$
- k -nearest neighbour rule:
 - For test data $x \in \mathcal{X}$, sort x_1, \dots, x_m according to $\|x - x_i\|$
 - $\pi_k(x)$ = index for is k -th nearest neighbour of x
 - Predict $h_S^{kNN}(x)$ = majority vote of $y_{\pi_1(x)}, \dots, y_{\pi_k(x)}$

OR for ± 1 labels,
$$h_S^{kNN}(x) = \text{sign} \left(\frac{1}{k} \sum_{i=1}^k y_{\pi_i(x)} \right)$$

Universal consistency of k -NN

Theorem kNN.5 (Universal consistency of k -NN)

If $k \rightarrow \infty$ and $\frac{k}{m} \rightarrow 0$ as $m \rightarrow \infty$, then for all distributions \mathcal{D} ,

$$\mathbb{E}_{S \sim \mathcal{D}^m} \left[L_{\mathcal{D}} \left(h_S^{kNN} \right) \right] \rightarrow L_{\mathcal{D}}^* \quad \text{as } m \rightarrow \infty$$

Proved by verifying conditions of Stone's theorem

Proof: Universal consistency of k -NN

- k -NN as plug-in classifier

$$\hat{\eta}(x) = \frac{1}{k} \sum_{i=1}^k \mathbf{1}\{y_{\pi_i(x)} = 1\} = \sum_{i=1}^m \mathbf{1}\{y_i = 1\} \underbrace{\frac{\mathbf{1}\{\pi_i(x) \leq k\}}{k}}_{=w_i(x)}$$

- $\sum_{i=1}^m w_i(x) \cdot f(x) = \frac{1}{k} \sum_{i=1}^k f(x_{\pi_i(x)}) \implies$ Condition (i) holds due to Stone's lemma
- Condition (ii) holds since $x_{\pi_k(x)} \rightarrow x$ in probability if $\frac{k}{m} \rightarrow 0$
- Condition (iii) holds since $\max_{i \in [m]} w_i(x) = \frac{1}{k} \rightarrow 0$ as $k \rightarrow \infty$

Recap Stone's theorem

Theorem kNN.6 (Stone's consistency theorem)

Let $\hat{\eta}(x) = \sum_{i=1}^m \mathbf{1}\{y_i = 1\} w_i(x)$, and $\hat{h}(x) = \mathbf{1}\{\hat{\eta}(x) \geq \frac{1}{2}\}$.

\hat{h} is universally consistent if weights $w_i(x)$ satisfy:

(i) $\exists c$ such that, for every non-negative integrable function f with $\mathbb{E}_{x \sim \mathcal{D}_{\mathcal{X}}}[f(x)] < \infty$,

$$\mathbb{E}_{x, x_1, \dots, x_m \sim \mathcal{D}^{m+1}} \left[\sum_{i=1}^m w_i(x) \cdot f(x_i) \right] \leq c \mathbb{E}_{x \sim \mathcal{D}_{\mathcal{X}}}[f(x)]$$

(ii) For all $a > 0$, $\lim_{m \rightarrow \infty} \mathbb{E}_{x, x_1, \dots, x_m \sim \mathcal{D}^{m+1}} \left[\sum_{i=1}^m w_i(x) \cdot \mathbf{1}\{\|x_i - x\| > a\} \right] = 0$

(iii) $\lim_{m \rightarrow \infty} \mathbb{E}_{x, x_1, \dots, x_m \sim \mathcal{D}^{m+1}} \left[\max_{i \in [m]} w_i(x) \right] = 0$

Proof of Stone's theorem: Main idea

Lemma kNN.7 (Risk bound for plug-in classifier)

Consider 0-1 loss, and let $\hat{h}(x) = \mathbf{1} \left\{ \hat{\eta}(x) \geq \frac{1}{2} \right\}$ be a plug-in classifier. Then

$$L_{\mathcal{D}}(\hat{h}) - L_{\mathcal{D}}^* \leq 2\mathbb{E}_{x \sim \mathcal{D}_{\mathcal{X}}} \left[|\hat{\eta}(x) - \eta(x)| \right] \leq 2\sqrt{\mathbb{E}_{x \sim \mathcal{D}_{\mathcal{X}}} \left[(\hat{\eta}(x) - \eta(x))^2 \right]}$$

Proof idea: (complete the steps)

- 2nd inequality uses Jensen's inequality: For a convex function f , $f(\mathbb{E}[z]) \leq \mathbb{E}[f(z)]$
- 1st inequality: First show that $L_{\mathcal{D}}(\hat{h}) - L_{\mathcal{D}}^* = 2\mathbb{E}_x \left[\left| \eta(x) - \frac{1}{2} \right| \cdot \mathbf{1} \left\{ h^*(x) \neq \hat{h}(x) \right\} \right]$

Then find upper bound for the term inside observing that, whenever $h^*(x) \neq \hat{h}(x)$,
 $\left| \eta(x) - \frac{1}{2} \right| \leq |\eta(x) - \hat{\eta}(x)|$

Proof of Stone's theorem: Main idea

- Taking expectation w.r.t S , we can write

$$\mathbb{E}_S[L_{\mathcal{D}}(\hat{h}) - L_{\mathcal{D}}^*] \leq 2\mathbb{E}_{S,x} \left[|\hat{\eta}(x) - \eta(x)| \right] \leq 2\sqrt{\mathbb{E}_{S,x} \left[(\hat{\eta}(x) - \eta(x))^2 \right]}$$

- Due to above, suffices to show

$$\mathbb{E}_{S,x} \left[(\hat{\eta}(x) - \eta(x))^2 \right] \rightarrow 0 \quad \text{as } m \rightarrow \infty$$

- Assumption: $\eta : \mathcal{X} \rightarrow [0, 1]$ is uniformly continuous
 - The assumption is not necessary, but simplifies parts of the proof

Proof of Stone's theorem: Main idea

- Recall $\hat{\eta}(x) = \sum_{i=1}^m \mathbf{1}\{y_i = 1\} w_i(x)$ and define $\tilde{\eta}(x) = \sum_{i=1}^m \eta(x_i) w_i(x)$

$$\begin{aligned}\text{Then } (\hat{\eta}(x) - \eta(x))^2 &= (\hat{\eta}(x) - \tilde{\eta}(x) + \tilde{\eta}(x) - \eta(x))^2 \\ &\leq 2\left((\hat{\eta}(x) - \tilde{\eta}(x))^2 + (\tilde{\eta}(x) - \eta(x))^2\right)\end{aligned}$$

- Separately show expectation of each squared term goes to 0

- Note: $\hat{\eta}(x) - \tilde{\eta}(x) = \sum_{i=1}^m (\mathbf{1}\{y_i = 1\} - \eta(x_i)) w_i(x)$
and $\tilde{\eta}(x) - \eta(x) = \sum_{i=1}^m (\eta(x_i) - \eta(x)) w_i(x)$

Proof: $(\hat{\eta}(x) - \tilde{\eta}(x))^2 \rightarrow 0$ in expectation

$$\begin{aligned}\mathbb{E}_{S,x} \left[(\hat{\eta}(x) - \tilde{\eta}(x))^2 \right] &= \mathbb{E} \left[\sum_{i=1}^m \sum_{j=1}^m (\mathbf{1} \{y_i = 1\} - \eta(x_i)) w_i(x) \cdot (\mathbf{1} \{y_j = 1\} - \eta(x_j)) w_j(x) \right] \\&= \mathbb{E} \left[\sum_{i=1}^m \underbrace{(\mathbf{1} \{y_i = 1\} - \eta(x_i))^2}_{\leq 1} (w_i(x))^2 \right] && \begin{array}{l} \text{given } x, x_i, x_j, \\ y_i \text{ independent of } y_j \\ \text{Expectation} = 0 \text{ for } i \neq j \end{array} \\&\leq \mathbb{E} \left[\sum_{i=1}^m (w_i(x))^2 \right] \\&\leq \mathbb{E} \left[\max_{i \in [m]} w_i(x) \cdot \underbrace{\sum_{i=1}^m w_i(x)}_{=1} \right] = \mathbb{E} \left[\max_{i \in [m]} w_i(x) \right] \rightarrow 0 \quad (\text{condition (iii)})\end{aligned}$$

Proof: $(\tilde{\eta}(x) - \eta(x))^2 \rightarrow 0$ in expectation

$$\begin{aligned}\mathbb{E}_{S,x} \left[(\tilde{\eta}(x) - \eta(x))^2 \right] &= \mathbb{E} \left[\left(\sum_{i=1}^m w_i(x) \cdot (\eta(x_i) - \eta(x)) \right)^2 \right] \\ &\leq \mathbb{E} \left[\sum_{i=1}^m w_i(x) \cdot (\eta(x_i) - \eta(x))^2 \right] \quad (\text{Jensen's inequality})\end{aligned}$$

- Jensen's inequality: For convex f and weights w_1, \dots, w_m such that $\sum_i w_i = 1$,

$$f \left(\sum_i w_i z_i \right) \leq \sum_i w_i f(z_i)$$

- Fix some $\epsilon > 0$. Since η is uniformly continuous,

$$\exists a_\epsilon > 0 \text{ such that } \|x_i - x\| \leq a_\epsilon \implies |\eta(x_i) - \eta(x)| \leq \epsilon$$

Proof: $(\tilde{\eta}(x) - \eta(x))^2 \rightarrow 0$ in expectation

$$\begin{aligned} & \mathbb{E} \left[\sum_{i=1}^m w_i(x) \cdot (\eta(x_i) - \eta(x))^2 \right] \\ & \leq \underbrace{\mathbb{E} \left[\sum_{i=1}^m w_i(x) \cdot \epsilon^2 \cdot \mathbf{1} \{ \|x_i - x\| \leq a_\epsilon \} \right]}_{\leq \epsilon^2} + \underbrace{\mathbb{E} \left[\sum_{i=1}^m w_i(x) \cdot 1 \cdot \mathbf{1} \{ \|x_i - x\| > a_\epsilon \} \right]}_{\rightarrow 0 \text{ due to condition (ii)}} \end{aligned}$$

- From above, for any $\epsilon > 0$,

$$\lim_{m \rightarrow \infty} \mathbb{E}_{S,x} \left[(\tilde{\eta}(x) - \eta(x))^2 \right] \leq \epsilon^2$$

- Hence, limit is 0

Proof: Need for condition (i)

- Above proof, assuming η is uniformly continuous, does not need condition (i)
- If uniform continuity is not assumed
 - $\eta(\cdot)$ is bounded \implies Can be approximated by a uniformly continuous function η^*

for any $\epsilon > 0$, \exists unif. cont. η^* such that $\mathbb{E}_x \left[(\eta(x) - \eta^*(x))^2 \right] < \epsilon$

- Use previous slide to prove $\mathbb{E}_{S,x} \left[(\tilde{\eta}^*(x) - \eta^*(x))^2 \right] \rightarrow 0$
- Need condition (i) to bound $\mathbb{E}_{S,x} \left[(\tilde{\eta}(x) - \tilde{\eta}^*(x))^2 \right]$ and $\mathbb{E}_{S,x} \left[(\eta^*(x) - \eta(x))^2 \right]$

Conclusion / Up Next

- Given infinite training data ($m \rightarrow \infty$), NN (or kNN) is good compared with optimal (Bayes) predictor, $L_{\mathcal{D}}^{NN} \leq 2L_{\mathcal{D}}^*$
- With $k \rightarrow \infty$ and $k/m \rightarrow 0$, kNN is universally consistent
... optimal for any \mathcal{D} if it has infinite training data, $m \rightarrow \infty$
- Next: What happens for finite m ?
 - For finite m , complex models can easily overfit
 - We restrict ERM to certain classes of models (say linear classifier)
 - For ERM solution \hat{h} , we bound $L_{\mathcal{D}}(\hat{h})$ in terms of complexity of model class