Context

Set-up Consider data $\mathcal{X} \subset \mathbb{R}^d$, labels $\mathcal{Y} \subset \mathbb{R}$ and a classifier $f: \mathcal{X} \mapsto \mathcal{Y}$. We approximate f by \hat{f} using n labelled data $(\boldsymbol{x}_1, y_1), \dots, (\boldsymbol{x}_n, y_n) \in \mathcal{X}_0 \times \mathcal{Y}, \mathcal{X}_0 \subset \mathcal{X}$.

Bias-variance tradeoff For $y = f(\boldsymbol{x}) + \epsilon$, $\epsilon \sim \mathrm{N}(0, \sigma^2)$, $\mathrm{E}([y - \hat{f}(\boldsymbol{x})]^2) = [\mathrm{Bias}(\hat{f}(\boldsymbol{x}))]^2 + \mathrm{var}(\hat{f}(\boldsymbol{x})) + \sigma^2,$ $\mathrm{Bias}(\hat{f}(\boldsymbol{x})) = \mathrm{E}(\hat{f}(\boldsymbol{x})) - f(\boldsymbol{x}),$ $\mathrm{var}(\hat{f}(\boldsymbol{x})) = \mathrm{E}([\hat{f}(\boldsymbol{x}) - \mathrm{E}(\hat{f}(\boldsymbol{x}))]^2).$

1 Bayes classifier

$$\hat{f}(\boldsymbol{x}) = \arg\max_{y \in \mathcal{Y}} P(\boldsymbol{X} = \boldsymbol{x}|Y = y)P(Y = y)$$

Training $\mathcal{O}(nd^2 + d^3)$

- Compute the class prior P(Y = y) from \mathcal{X}_0 .
- Choose a probabilistic model P(X = x | Y = y).
- Find parameters for P(X = x|Y = y) that minimize the likelihood (MVN: compute \overline{x} , S and S^{-1}).

Testing $\mathcal{O}(d^2)$ Evaluate $\hat{f}(\boldsymbol{x})$ (MVN: products $\boldsymbol{S}^{-1}\boldsymbol{x}$).

2 k-NNs

$$\hat{f}(\boldsymbol{x}) = \text{label}\left(\arg\min_{1 \leq i \leq n} d_{\boldsymbol{W}}(\boldsymbol{x}, \boldsymbol{x}_i)\right) \quad (k = 1)$$

Training $\mathcal{O}(nd^2s + n\log(n)d)$

- Find distance $d_{\mathbf{W}}(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i \mathbf{x}_j)^T \mathbf{W}(\mathbf{x}_i \mathbf{x}_j)$, with s steps of BFGS, that best separates the data.
- Choose k. Find medians to produce k-d trees.

Testing $\mathcal{O}(\log(n)d)$ Evaluate $\hat{f}(\boldsymbol{x})$.

3 Perceptron $\mathcal{Y} = \{-1, 1\}$

$$\hat{f}(\boldsymbol{x}) = \operatorname{sign}\left(\boldsymbol{w}^T \boldsymbol{x} + w_0\right)$$

Training $\mathcal{O}(ndP)$, with $P = \max_{1 \le i \le n} ||x_i||^2 / \gamma^2$ (margin γ)

- Initialize $\mathbf{w} = \mathbf{0}$ and $w_0 = 0$.
- For k = 1, 2, ...,
 - if there is $x_i \in \mathcal{X}_0$ such that

$$\operatorname{sign}\left(\boldsymbol{w}^T\boldsymbol{x}_i + w_0\right) \neq y_i,$$

set $\mathbf{w} = \mathbf{w} + y_i \mathbf{x}_i$ and $w_0 = w_0 + y_i$.

Testing $\mathcal{O}(nd)$ Evaluate $\hat{f}(\boldsymbol{x})$.

4 SVMs $\mathcal{Y} = \{-1, 1\}$

$$\hat{f}(\boldsymbol{x}) = \operatorname{sign}\left(\boldsymbol{w}^T \boldsymbol{x} + w_0\right)$$

Training $\mathcal{O}(nd^2s)$ Minimize, with s steps of PGD,

$$J(\mathbf{w}) = \frac{1}{2} ||\mathbf{w}||^2$$
, such that $y_i(\mathbf{w}^T \mathbf{x}_i + w_0) \ge 1$, $1 \le i \le n$.

Testing $\mathcal{O}(d)$ Evaluate $\hat{f}(\boldsymbol{x})$.

5 Logistic regression $\mathcal{Y} = \{0, 1\}$

$$\hat{f}(\boldsymbol{x}) = \frac{1}{1 + e^{-(\boldsymbol{w}^T \boldsymbol{x} + w_0)}}$$

Training $\mathcal{O}(nd^3s)$ Minimize the likelihood with s steps of Newton's method.

Testing $\mathcal{O}(d)$ Evaluate $\hat{f}(\boldsymbol{x})$.

6 Neural networks

$$\hat{f}(\boldsymbol{x}) = \sum_{i=1}^{N} \alpha_i \sigma(\boldsymbol{w}_i^T \boldsymbol{x} + w_{i,0})$$
 (one layer)

Training $\mathcal{O}(ndsN^2)$ Minimize, with s steps of SGD,

$$J(\boldsymbol{\alpha}, \boldsymbol{w}, \boldsymbol{w}_0) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{1}[\hat{f}(\boldsymbol{x}_i) \neq y_i].$$

Testing $\mathcal{O}(dN)$ Evaluate $\hat{f}(\boldsymbol{x})$.

7 Decision trees

$$\hat{f}(\boldsymbol{x}) = \text{label}(C(\boldsymbol{x}))$$

Training Find feature and threshold that maximize

$$u(C) - (p_L u(C_L) + p_R u(C_R), \ u(C) = -\sum_{y \in \mathcal{V}} p_y \log p_y.$$

Testing $\mathcal{O}(\log(n)d)$ Evaluate $\hat{f}(x)$.

8 Random forest

$$\hat{f}(\boldsymbol{x}) = \frac{1}{T} \sum_{t=1}^{T} \text{label}(C_t(\boldsymbol{x}))$$

Training Sample T times, with replacement, m < n training examples from \mathcal{X}_0 to construct T decision trees.

Testing $\mathcal{O}(\log(n)dT)$ Evaluate $\hat{f}(\boldsymbol{x})$.