

Reducing Supervisor Burden in DAgger by one class SVM

Huan Ling

2019-01-25

(Ross & Gordon & Bagnell, 2011): DAgger, or Dataset Aggregation

- Imitation learning as interactive supervision
- Aggregate training data from expert with test data from execution

Algorithm 1 DAgger

- 1: $D = \{(s, a)\}$ initial expert demonstrations
 - 2: $\theta_1 \leftarrow$ train learner's policy parameters on D
 - 3: **for** $i = 1 \dots N$ **do**
 - 4: Execute learner's policy π_{θ_i} , get visited states $S_{\theta_i} = \{s_0, \dots, s_T\}$
 - 5: Query the expert at those states to get actions $A = \{a_0, \dots, a_T\}$
 - 6: Aggregate dataset $D = D \cup \{(s, a) \mid s \in S_{\theta_i}, a \in A\}$
 - 7: Train learner's policy $\pi_{\theta_{i+1}}$ on dataset D
 - 8: Return one of the policies π_{θ_i} that performs best on validation set
-

(Ross & Gordon & Bagnell, 2011): DAgger, or Dataset Aggregation

- Imitation learning as interactive supervision
- Aggregate training data from expert with test data from execution

Algorithm 1 DAgger

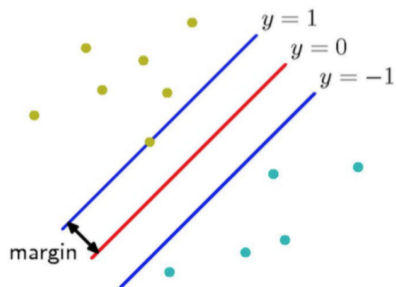
- 1: $D = \{(s, a)\}$ initial expert demonstrations
 - 2: $\theta_1 \leftarrow$ train learner's policy parameters on D
 - 3: **for** $i = 1 \dots N$ **do**
 - 4: Execute learner's policy π_{θ_i} , get visited states $S_{\theta_i} = \{s_0, \dots, s_T\}$
 - 5: Query the expert at those states to get actions $A = \{a_0, \dots, a_T\}$
 - 6: Aggregate dataset $D = D \cup \{(s, a) \mid s \in S_{\theta_i}, a \in A\}$
 - 7: Train learner's policy $\pi_{\theta_{i+1}}$ on dataset D
 - 8: Return one of the policies π_{θ_i} that performs best on validation set
-

- ▶ Dagger is time-consuming and it requires human to annotate data for each state.
- ▶ Can we save some human annotations by algorithm?

Two Class SVM Review

- ▶ We are given training data $(x^{(i)}, t^{(i)})$.
- ▶ We look at classification, so $t^{(i)}$ will represent the class label.
- ▶ We use $t = 1$ for the positive and $t = -1$ for the negative class.

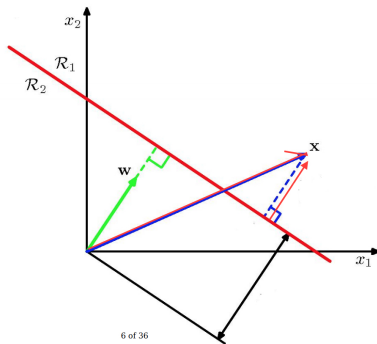
Two Class SVM Review



- ▶ **Goal:** Find classification boundary that leads to the largest margin (buffer) from points on both sides.
- ▶ Subset of vectors that support (determine boundary) are called the **support vectors**.

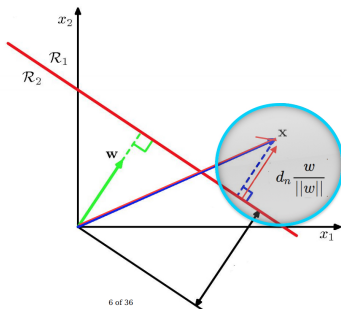
Two Class SVM Review

- ▶ Let decision boundary to be $\omega^T \mathbf{x} + b = 0$, so that ω perpendicular to decision boundary.
- ▶ Let \mathbf{x} to be support vector, d_n is the distance from support vector to boundary.



Two Class SVM Review

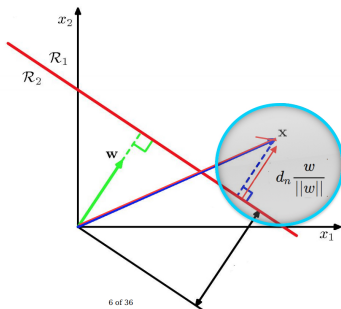
- ▶ $d_n \frac{\omega}{\|\omega\|}$ is the vector perpendicular to decision boundary and has length d_n .



6 of 36

Two Class SVM Review

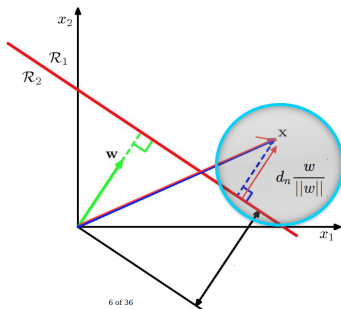
- ▶ $d_n \frac{\omega}{\|\omega\|}$ is the vector perpendicular to decision boundary and has length d_n .
- ▶ Vector $x + d_n \frac{\omega}{\|\omega\|}$ locates on decision boundary.



6 of 36

Two Class SVM Review

- ▶ $d_n \frac{\omega}{\|\omega\|}$ is the vector perpendicular to decision boundary and has length d_n .
- ▶ Vector $x + d_n \frac{\omega}{\|\omega\|}$ locates on decision boundary.
- ▶ So $\omega^T(x + d_n \frac{\omega}{\|\omega\|}) + b = 0$



6 of 36

Two Class SVM Review

- ▶ $d_n \frac{\omega}{\|\omega\|}$ is the vector perpendicular to decision boundary and has length d_n .

Two Class SVM Review

- ▶ $d_n \frac{\omega}{\|\omega\|}$ is the vector perpendicular to decision boundary and has length d_n .
- ▶ Vector $x + d_n \frac{\omega}{\|\omega\|}$ locates on decision boundary.

Two Class SVM Review

- ▶ $d_n \frac{\omega}{\|\omega\|}$ is the vector perpendicular to decision boundary and has length d_n .
- ▶ Vector $x + d_n \frac{\omega}{\|\omega\|}$ locates on decision boundary.
- ▶ So $\omega^T (x + d_n \frac{\omega}{\|\omega\|}) + b = 0$

Two Class SVM Review

- ▶ $d_n \frac{\omega}{\|\omega\|}$ is the vector perpendicular to decision boundary and has length d_n .
- ▶ Vector $x + d_n \frac{\omega}{\|\omega\|}$ locates on decision boundary.
- ▶ So $\omega^T (x + d_n \frac{\omega}{\|\omega\|}) + b = 0$
- ▶ Derive d_n we have, $d_n = \frac{\omega^T x + b}{\|\omega\|}$

Two Class SVM Review

- ▶ $d_n \frac{\omega}{\|\omega\|}$ is the vector perpendicular to decision boundary and has length d_n .
- ▶ Vector $x + d_n \frac{\omega}{\|\omega\|}$ locates on decision boundary.
- ▶ So $\omega^T(x + d_n \frac{\omega}{\|\omega\|}) + b = 0$
- ▶ Derive d_n we have, $d_n = \frac{\omega^T x + b}{\|\omega\|}$
- ▶ We let $t_n \in \{-1, 1\}$ stands for label for current data x , then
$$d_n = \frac{t_n(\omega^T x + b)}{\|\omega\|}$$

Two Class SVM Review

We can set $d_n = \frac{1}{\|w\|}$ for the point x_n closest to the decision boundary, leading to the problem:

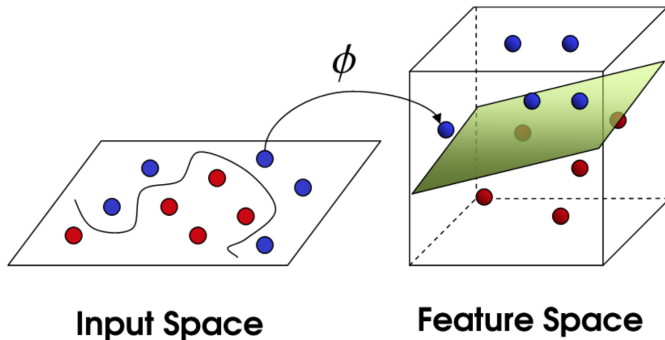
$$\begin{aligned} & \max \frac{1}{\|w\|} \\ \text{s.t. } & t_n(w^T x_n + b) \geq 1, \text{ for } n = 1 \dots N \end{aligned}$$

► Or equivalently:

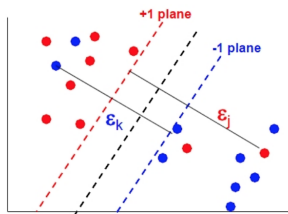
$$\begin{aligned} & \min \frac{1}{2} \|w\|^2 \\ \text{s.t. } & t_n(w^T x_n + b) \geq 1, \text{ for } n = 1 \dots N \end{aligned}$$

Non-linear SVM

- ▶ We change $y(x) = \omega^T(x) + b$ to $y(x) = \omega^T \phi((x)) + b$
- ▶ $\phi()$ is called kernel function.



Two Class SVM Review, slack variables



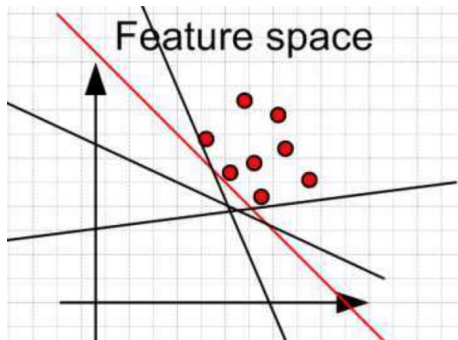
- Introduce slack variables ξ_i

$$\min \frac{1}{2} \|\mathbf{w}\|^2 + \lambda \sum_{i=1}^N \xi_i$$

$$\text{s.t. } \xi_i \geq 0; \quad \forall i \quad t^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \xi_i$$

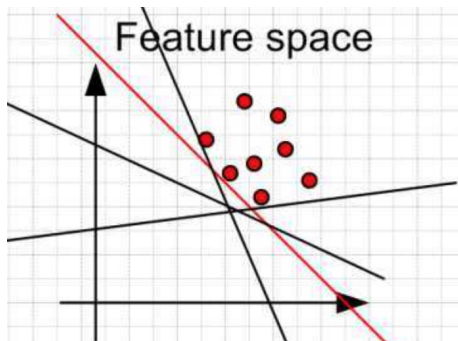
- ▶ Example x_i lies on wrong side of hyperplane then $\xi_i > 1$
- ▶ So $\sum_i \xi_i$ is the upper bounds of the number of training error.

One Class SVM



- ▶ We want to separate all data from **origin**.
- ▶ We want to the decision boundary to be as far as possible to **origin**.
- ▶ The red line above is the ideal decision boundary we want.

One Class SVM



- ▶ Why do we want one class SVM?
- ▶ One class SVM means all data in training set are in the same class.
- ▶ When a new data comes in, we know whether the data is worth labeling.

One Class SVM

$$\begin{aligned} \min_{\mathbf{w}, b, \xi, \rho} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} - \rho + \frac{1}{\nu l} \sum_{i=1}^l \xi_i \\ \text{subject to} \quad & \mathbf{w}^T \phi(\mathbf{x}_i) - \rho \geq -\xi_i \\ & \xi_i \geq 0, \quad i = 1, \dots, l \end{aligned}$$

- ▶ Same as two class SVM, we add slack term ν .
- ▶ The parameter ν controls the penalty or slack term and is an upper bound on the fraction of outliers in training set.
- ▶ In another word, when $\nu = 0.1$ then only 0.1 of training data will be classified to wrong class.

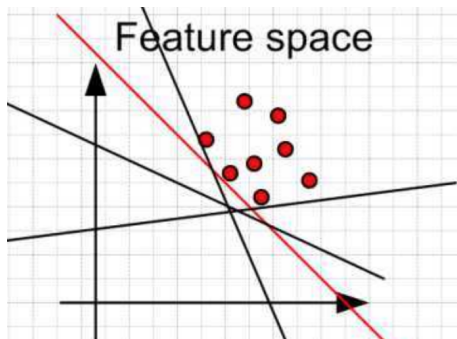
$$\begin{aligned} \min_{\mathbf{w}, b, \xi, \rho} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} - \rho + \frac{1}{\nu l} \sum_{i=1}^l \xi_i \\ \text{subject to} \quad & \mathbf{w}^T \phi(\mathbf{x}_i) - \rho \geq -\xi_i \\ & \xi_i \geq 0, \quad i = 1, \dots, l \end{aligned}$$

- ▶ Same as two class SVM, we add slack term ν .
- ▶ The parameter ν controls the penalty or slack term and is an upper bound on the fraction of outliers in training set.
- ▶ In another word, when $\nu = 0.1$ then only 0.1 of training data will be classified to wrong class.

$$\begin{aligned} \min_{\mathbf{w}, b, \xi, \rho} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} - \rho + \frac{1}{\nu l} \sum_{i=1}^l \xi_i \\ \text{subject to} \quad & \mathbf{w}^T \phi(\mathbf{x}_i) - \rho \geq -\xi_i \\ & \xi_i \geq 0, \quad i = 1, \dots, l \end{aligned}$$

- ▶ Same as two class SVM, we add slack term ν .
- ▶ The parameter ν controls the penalty or slack term and is an upper bound on the fraction of outliers in training set.
- ▶ In another word, when $\nu = 0.1$ then only 0.1 of training data will be classified to wrong class.

One Class SVM



- ▶ One class SVM means all data in training set are in the same class.
- ▶ When a new data comes in, we know whether the data is worth labeling.
- ▶ However, standard one class SVM doesn't have sense of different classes.

- ▶ We separate our training data in two category, success and failure based on y :

$$y_i = \begin{cases} 1 & : l(\pi_\theta(\mathbf{x}_i), \mathbf{u}_i) \leq \varepsilon \\ -1 & : l(\pi_\theta(\mathbf{x}_i), \mathbf{u}_i) > \varepsilon \end{cases}$$

Where as l is loss function, π_θ is current policy, ε is human defined threshold.

One Class SVM

- ▶ We separate our training data in two category, success and failure based on y :

$$y_i = \begin{cases} 1 & : l(\pi_\theta(\mathbf{x}_i), \mathbf{u}_i) \leq \varepsilon \\ -1 & : l(\pi_\theta(\mathbf{x}_i), \mathbf{u}_i) > \varepsilon \end{cases}$$

Where as l is loss function, π_θ is current policy, ε is human defined threshold.

- ▶ $D_s = \{\{x_i, u_i\} \in D_k, y_i = 1\}$, $D_r = \{\{x_i, u_i\} \in D_k, y_i = -1\}$

One Class SVM

- ▶ We separate our training data in two category, success and failure based on y :

$$y_i = \begin{cases} 1 & : l(\pi_\theta(\mathbf{x}_i), \mathbf{u}_i) \leq \varepsilon \\ -1 & : l(\pi_\theta(\mathbf{x}_i), \mathbf{u}_i) > \varepsilon \end{cases}$$

Where as l is loss function, π_θ is current policy, ε is human defined threshold.

- ▶ $D_s = \{\{x_i, u_i\} \in D_k, y_i = 1\}$, $D_r = \{\{x_i, u_i\} \in D_k, y_i = -1\}$
- ▶ A separate One-Class SVM is then trained on each set of states, and providing measures of both set, g_s and g_r .

One Class SVM

- ▶ We separate our training data in two category, success and failure based on y :

$$y_i = \begin{cases} 1 & : l(\pi_\theta(\mathbf{x}_i), \mathbf{u}_i) \leq \varepsilon \\ -1 & : l(\pi_\theta(\mathbf{x}_i), \mathbf{u}_i) > \varepsilon \end{cases}$$

Where as l is loss function, π_θ is current policy, ε is human defined threshold.

- ▶ $D_s = \{\{x_i, u_i\} \in D_k, y_i = 1\}$, $D_r = \{\{x_i, u_i\} \in D_k, y_i = -1\}$
- ▶ A separate One-Class SVM is then trained on each set of states, and providing measures of both set, g_s and g_r .
- ▶ So our decision function becomes:

$$g_\sigma(\mathbf{x}) = \begin{cases} 0 & : g_s(\mathbf{x}) == 1 \text{ and } g_r(\mathbf{x}) == -1 \\ -1 & : \text{otherwise} \end{cases}$$

SHIV and DAgger

- ▶ Both SHIV and DAgger solve the minimization loss iterating two steps.

SHIV and DAgger

- ▶ Both SHIV and DAgger solve the minimization loss iterating two steps.
- ▶ **Step 1:** Compute a θ using the training data so far

SHIV and DAgger

- ▶ Both SHIV and DAgger solve the minimization loss iterating two steps.
- ▶ **Step 1:** Compute a θ using the training data so far
- ▶ **Step 2:** Execute the policy induced by the current θ , and ask for labels for the encountered states.

SHIV and DAgger

- ▶ Both SHIV and DAgger solve the minimization loss iterating two steps.
- ▶ **Step 1:** Compute a θ using the training data so far
- ▶ **Step 2:** Execute the policy induced by the current θ , and ask for labels for the encountered states.
- ▶ **DAgger:** Query supervisor for every new state.

$$D_{k+1} = \mathcal{D}_k \cup \{(\mathbf{x}_t, \tilde{\mathbf{u}}_t) \mid t \in \{0, \dots, T\}\}$$

SHIV and DAgger

- ▶ Both SHIV and DAgger solve the minimization loss iterating two steps.
- ▶ **Step 1:** Compute a θ using the training data so far
- ▶ **Step 2:** Execute the policy induced by the current θ , and ask for labels for the encountered states.
- ▶ **DAgger:** Query supervisor for every new state.

$$D_{k+1} = \mathcal{D}_k \cup \{(\mathbf{x}_t, \tilde{\mathbf{u}}_t) \mid t \in \{0, \dots, T\}\}$$

- ▶ **SHIV:** Using One Class SVM to decide whether to ask a query.

$$D_{k+1} = \mathcal{D}_k \cup \{(\mathbf{x}_t, \tilde{\mathbf{u}}_t) \mid t \in \{0, \dots, T\}, g(\mathbf{x}_t) = -1\}$$

Results

