

What is machine learning



Step 1: Turn X into a vector space
(feature space)

Step 2: Characterize a loss:

$\text{loss}(f(x), y)$ measures how well $f(x)$ predicts y .

(usually $\text{loss}(y, y) = 0$, $\text{loss}(z, y) \geq 0$)

Step 3: pick a family \mathcal{F} of functions that might predict y from x

Step 4: collect data (examples) $(x_1, y_1), \dots, (x_n, y_n)$

Step 5: Minimize $\frac{1}{n} \sum_{i=1}^n \text{loss}(f(x_i), y_i)$
 $f \in \mathcal{F}$

Step 6: ? What do you do with the optimal f ?

How does it perform on new data?

Example :

$$X = \mathbb{R}^d, \quad Y = \mathbb{R}$$

$$\text{loss}(f(x), y) = (f(x) - y)^2$$

$$\mathcal{F} = \{ f(x) = w^T x : w \in \mathbb{R}^d \}$$

$$\underset{w}{\text{Minimize}} \quad -\frac{1}{n} \sum_{i=1}^n \frac{1}{2} (w^T x_i - y_i)^2 \quad \text{LEAST SQUARES}$$

We can rewrite this, defining

$$\bar{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \quad \bar{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

Then cost is $\frac{1}{2n} \|Xw - y\|^2$

$$\|v\| = \left(\sum_{i=1}^n v_i^2 \right)^{1/2} \quad \text{is Euclidean norm}$$

Least-squares solution

I. Orthogonal projection

V an inner product space

$S \subset V$ a subspace

Any $\vec{v} \in V$ can be decomposed as

$$\vec{v} = \vec{v}_S + \vec{v}_L \quad \text{where} \quad \vec{v}_S \in S \quad \langle \vec{v}_L, \vec{u} \rangle = 0 \quad \forall \vec{u} \in S$$

\bar{P}_S is the linear projection onto S

$$\bar{P}_S(\vec{v}) = \vec{v}_S$$

$$\text{FACT: } \|\vec{v} - \bar{P}_S \vec{v}\| \leq \|\vec{v} - \vec{u}\| \quad \forall \vec{u} \in S$$

$$\Rightarrow \min_{\vec{u} \in S} \|\vec{v} - \vec{u}\| = \|\vec{v} - \bar{P}_S \vec{v}\| \quad (\text{definition of min})$$

$$\arg \min_{\vec{u} \in S} \|\vec{v} - \vec{u}\| = \bar{P}_S \vec{v}$$

$\arg \min$ = set of minimizers
"argument of the minimum"

$$\begin{aligned} \text{PROOF OF FACT: } \|\vec{v} - \vec{u}\|^2 &= \|\vec{v} - \bar{P}_S \vec{v} + \bar{P}_S \vec{v} - \vec{u}\|^2 \\ &= \|\vec{v} - \bar{P}_S \vec{v}\|^2 + \|\bar{P}_S \vec{v} - \vec{u}\|^2 \geq \|\vec{v} - \bar{P}_S \vec{v}\|^2 \end{aligned}$$

Projections and Least-Squares

$$\min_{\vec{w}} \|\bar{X} \vec{w} - \vec{y}\| = \min_{\vec{v} \in \text{range}(\bar{X})} \|\vec{v} - \vec{y}\|$$

$$\arg \min_{\vec{v} \in \text{range}(\bar{X})} \|\vec{v} - \vec{y}\| = P_{\text{range}(\bar{X})} \vec{y}$$

Now we also know

$$P_{\text{range}(\bar{X})} \vec{y} - \vec{y} \in \text{range}(\bar{X})^\perp = \text{null}(\bar{X}^T)$$

$$\bar{X}^T (P_{\text{range}(\bar{X})} \vec{y} - \vec{y}) = 0$$

$$\bar{X}^T (\bar{X} \vec{w}_{OLS} - \vec{y}) = 0$$

if $\bar{X}^T \bar{X}$ is invertible:

$$\vec{w}_{OLS} = (\bar{X}^T \bar{X})^{-1} \bar{X}^T \vec{y}$$

only invertible if $d < n \dots$

$\bar{X} \vec{w}_{OLS}$ always unique

\vec{w}_{OLS} not always unique.

Note: $\text{null}(\bar{X}) = \text{range}(\bar{X}^T)^\perp$. Therefore can

assume $\vec{w}_{OLS} = \bar{X}^\perp \alpha \rightarrow \vec{w}_{OLS} = \sum_{i=1}^n \alpha_i \vec{x}_i$

\vec{w}_{OLS} is in span of examples.

Optimization approach

$$\cancel{W} \cdot W_{\star} \in \arg \min \| \bar{X} \bar{w} - \bar{y} \|^2$$

Means

$$\| \bar{X} (\bar{w}_{\star} + \Delta \bar{w}) - \bar{y} \|^2 \geq \| \bar{X} \bar{w}_{\star} - \bar{y} \|^2 \quad \forall \Delta \bar{w}$$

expanding terms, this means

$$\langle \bar{X} \bar{w} - \bar{y}, \bar{X} \Delta \bar{w} \rangle + \| \bar{X} \Delta \bar{w} \|^2 \geq 0$$

$$\Rightarrow \langle \bar{X}^T (\bar{X} \bar{w} - \bar{y}), \Delta \bar{w} \rangle + \| \bar{X} \Delta \bar{w} \|^2 \geq 0$$

$$\Rightarrow \langle \bar{X}^T (\bar{X} \bar{w} - \bar{y}), \frac{\Delta \bar{w}}{\| \Delta \bar{w} \|} \rangle + \frac{\| \bar{X} \Delta \bar{w} \|^2}{\| \Delta \bar{w} \|} \geq 0$$

again implying

$$\bar{X}^T (\bar{X} \bar{w} - \bar{y}) = 0$$