

Least squares

- So far we have looked at binary prediction problems with zero-one loss

$$l(h, (x, y)) = \mathbb{1}[h(x) \neq y]$$

classifier

labeled
example

- Finding classifiers with optimal zero-one loss is often a hard optimization problem (e.g. NP-hard for halfspaces)
- There are ML tasks where the label y is not binary.
- We consider the setting

where the set of labels

$$Y = \mathbb{R} \quad \text{and} \quad l(h, (x, y)) = \underbrace{(h(x) - y)^2}_{\text{square loss}}$$

↑ space of labels ↑ predictor $h: X \rightarrow \mathbb{R}$ ↑ real-valued label

- For a labeled sample

$$S = ((x_1, y_1), \dots, (x_m, y_m))$$

where $x_1, x_2, \dots, x_m \in X$ and

$$y_1, y_2, \dots, y_m \in \mathbb{R}$$

we define empirical loss

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m l(h, (x_i, y_i))$$

- For a distribution \mathcal{D} over $X \times Y$
... define generalization error

we view

(risk, population loss) as

$$L_{\mathcal{D}}(h) = \mathbb{E} \left[l(h, (x, y)) \right]$$

where (x, y) is a random sample from \mathcal{D} .

- In general, we can use any loss function

$$l: \mathcal{Y}^X \times X \times \mathcal{Y} \rightarrow \mathbb{R}^d$$

- Zero-one loss and square loss are just two examples

Linear least squares
(Ordinary least squares)

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

$$H_d = \{ h_w : w \in \mathbb{R}^d \} \subseteq Y^X$$

where $h_w: \mathbb{R}^d \rightarrow \mathbb{R}$,

$$h_w(x) = w^T x$$

and as before

$$l(h, (x, y)) = (h(x) - y)^2$$

ERM for linear least squares

• We are given a labeled sample

$$S = ((x_1, y_1), \dots, (x_m, y_m)) \text{ where}$$

$x_1, x_2, \dots, x_m \in \mathbb{R}^d$ and

$y_1, y_2, \dots, y_m \in \mathbb{R}$

$y_1, y_2, \dots, y_m \in \mathbb{R}$.

- We want find the minimizer of the empirical loss:

$$\hat{h} = \operatorname{argmin}_{h \in H_d} L_S(h)$$

- This is equivalent to finding $\hat{w} \in \mathbb{R}^d$ such that

$$\hat{w} = \operatorname{argmin}_{w \in \mathbb{R}^d} L_S(h_w)$$

$$= \operatorname{argmin}_{w \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m (w^T x_i - y_i)^2$$

- Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$

$$f(w) = \frac{1}{m} \sum_{i=1}^m (w^T x_i - y_i)^2$$

- Then, \hat{w} is the minimizer of function f .
- Since f is differentiable, any global (or local) minimizer is a stationary point. That is, we can find \hat{w} by solving

$$\nabla f(w) = 0$$

where ∇f is the gradient of f :

$$\nabla f(w) = \left(\frac{\partial f(w)}{\partial w_1}, \frac{\partial f(w)}{\partial w_2}, \dots, \frac{\partial f(w)}{\partial w_n} \right)$$

vector of partial derivatives

- The gradient is

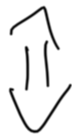
$$\nabla f(w) = \frac{1}{m} \sum_{i=1}^m 2 (w^T x_i - \gamma_i) x_i$$

- The equation $\nabla f(w) = 0$ is equivalent to

$$\sum_{i=1}^m (w^T x_i - \gamma_i) x_i = 0$$

- Equivalently

$$\sum_{i=1}^m (w^T x_i) x_i = \sum_{i=1}^m \gamma_i x_i$$



$$\left(\sum_{i=1}^m x_i x_i^T \right) w = \sum_{i=1}^m \gamma_i x_i$$

$X_i X_i^T$ is
a $d \times d$ matrix
(outer product of X_i
with itself)

• Let $A = \sum_{i=1}^m X_i X_i^T$ and $b = \sum_{i=1}^m \gamma_i X_i$

• The last equation is equivalent
to

$$A w = b$$

• Assuming A is invertible

$$w = A^{-1} b$$

• It turns out that $A w = b$

always has a solution, even if A is not invertible. The reason is that A and b are related.

- If A is not invertible there are infinitely many solutions.
- The solution of $Ax = b$ that minimizes $\|x\|_2$ is unique and can be computed using pseudo-inverse

$$x = A^+ b$$

↑
pseudo-inverse

- Since A is symmetric, it has eigen decomposition

$$A = V^T D V$$

where D is $d \times d$ diagonal matrix and V is $d \times d$ orthogonal matrix, i.e.

$$V^T = V^{-1}$$

• If $D = \begin{pmatrix} \alpha_1 & & 0 \\ & \ddots & \\ 0 & & \alpha_d \end{pmatrix}$ then the pseudo-inverse of D is

$$D^+ = \begin{pmatrix} \beta_1 & & 0 \\ & \ddots & \\ 0 & & \beta_d \end{pmatrix} \text{ where}$$

$$\beta_i = \begin{cases} 1/\alpha_i & \text{if } \alpha_i \neq 0 \\ 0 & \text{if } \alpha_i = 0 \end{cases}$$

- If M is an invertible matrix then

$$M^{-1} = M^{\dagger}$$

- Pseudo-inverse of A is

$$A^{\dagger} = (V^T D V)^{\dagger}$$

$$= V^{\dagger} D^{\dagger} (V^T)^{\dagger}$$

$$= V^T D^{\dagger} V$$

- This the same rule as for standard inverse.

$$V^+ = V^{-1} = V^T \quad \checkmark$$

$$(V^T)^+ = (V^T)^{-1} = V$$

Polynomial Regression

- Let $Y = \mathbb{R}$, $X = \mathbb{R}$ and

$$H_d = \left\{ p : \begin{array}{l} p(x) \text{ is a polynomial,} \\ \deg(p) \leq d \end{array} \right\}$$

$$\ell(h, (x, y)) = (h(x) - y)^2$$

- The problem can be reduced to linear least squares.
- Suppose

$$p(x) = w_0 + w_1 x + \dots + w_d x^d$$

$$= w^T \psi(x)$$

where $w = (w_0, w_1, \dots, w_d) \in \mathbb{R}^{d+1}$

and $\psi: \mathbb{R}^d \rightarrow \mathbb{R}^{d+1}$

$$\psi(x) = (1, x, x^2, \dots, x^d)$$

- ψ is called a feature mapping

(More specifically, ψ is called polynomial feature mapping.)

- Suppose $S = ((x_1, y_1) \dots (x_m, y_m))$

is a labeled sample where

$$x_1, x_2, \dots, x_m \in \mathbb{R}^d$$

and $y_1, y_2, \dots, y_m \in \mathbb{R}$

• The ERM problem is

$$\hat{p} = \underset{p: \deg(p) \leq d}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m (p(x_i) - y_i)^2$$

• Equivalently,

$$\hat{w} = \underset{w \in \mathbb{R}^{d+1}}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m (w^T \psi(x_i) - y_i)^2$$

• We can find \hat{w} using linear least squares

1000 - 1000