

Empirical Risk Minimization.

Available dataset (X, Y) = collection of samples from $P(X, Y)$

Underlying joint probability distribution unknown to us.

Goal of model development — iteratively train fw using (X, Y)

Optimality obtained when model trained on $(X, Y)_{\text{train}}$ results in minimal cost $Q(w)$ on $(X, Y)_{\text{test}}$

identify optimal function f_w^* from all possible hypotheses f_w .

Proxy to future unknown data.

$$Q(w) = L(Y, f_w(x))$$

One realization for $P(X, Y)$ — we want model that produces minimal expected cost across all variations of $(X, Y)_{\text{test}}$

cost is a function of random input data pair (x, y)

Measures generalization error of model.

RISK

Risk for model f_w :

$$R(f_w) = E[Q(w)] = E[L(Y, f_w(x))] = \int L(Y, f_w(x)) dP(X, Y)$$

Since joint distribution $P(X, Y)$ is unknown, true risk approx. via empirical risk

Realization = observed value of random variable.

Integral over all possible realizations of data based on joint distribution $P(X, Y)$

MSE on $(X, Y)_{\text{test}}$.

$$R(f_w)_{\text{empirical}} = \frac{1}{n} \sum_{i=1}^n (y_i - f_w(x_i))^2$$

one that minimizes empirical risk based on test data.

$$f_w^* = \arg \min_{f_w \in F_w} R_{\text{empirical}}(f_w)$$

for regression:

$$R(f_w) = E[y - f_w(x)]^2$$

In linear regression: f is a linear model parameterized by w .

→ If real mapping function is g : $Y = g(X) + \epsilon$

additive random noise
↓
Gaussian $\begin{cases} \mu=0 \\ \sigma=\sigma \end{cases}$

⇒ We can build a model such that $E[f(X)] = E[g(X) + \epsilon | x=X]$

Now, bias = lack of capacity in model to perfectly fit g (underfitting)

Variance = error due to fitting a spurious relation by random noise (overfitting)
→ seem related, but really not.

Consider an unseen case X_{new} (approx)

let - $f(X_{\text{new}}) = f$

(TRUE) $g(X_{\text{new}}) = g$ → deterministic + non-random.

$$R(f) = E[(Y - f(X_{\text{new}}))^2 | x = X_{\text{new}}]$$

$$= E[(y - f)^2]$$

$$= E[(y - g + g - f)^2]$$

$$= E[(y - g)^2] + E[(g - f)^2] + 2E[(y - g)(g - f)]$$

$$= E[(y-g)^2] + E[(g-f)^2] + 2E[(y-g)(g-f)]$$

$y = g + \varepsilon$

math later

$$= E[\varepsilon^2] + E(g-f)^2$$

$$= \sigma^2 + E[g - E(f)]^2 + E[f - E(f)]^2$$

$$R(f) = \underbrace{\sigma^2}_{\text{cannot be reduced}} + \underbrace{\text{Bias}^2(f)}_{\text{controllable by training}} + \text{Var}(f)$$

Simple model

- high bias
- low variance

Complex model

- low bias
- high variance