

Training loss $L_{\text{train}} = L(f, X_{\text{train}}, Y_{\text{train}})$.

Test loss $L_{\text{test}} = L(f, X_{\text{test}}, Y_{\text{test}})$.

Validation Loss $L_{\text{valid}} = L(f, X_{\text{valid}}, Y_{\text{valid}})$.

These pairs of (X, Y) are sampled from population loss D_x, D_y .

ultimate goal: minimize $\underset{f}{\text{minimize}} \quad L_{\text{population}} = \mathbb{E}_{X, Y \sim D_x, D_y} L(f, X, Y)$.

Overfit vs. Underfit:

Underfitting - your function may not have enough representation power

overfitting - your function has too much representation power.
 - gives $L_{\text{train}} \approx 0$; but generalizes bad.
 - Regularization.

Modern View:

For some neural network, there are implicit regularizations

↳ explicit regularization not necessary.

(perhaps because of SGD algo.).

Unsupervised Learning

- Clustering → Data mining / recommendation system.
- Principal component Analysis (Find the most important components (directions))
- Generative Model ⇒ Describe distribution of data . \hookrightarrow minimize MSE.
 by mapping Gaussian to target distribution
- Anomaly Detection.
- Dimensional Reduction

Semi-Supervised Learning

data ↗ labeled
unlabeled

Can we use unlabeled data to improve prediction?

Assumptions:
points with small distances
are more likely to share same labels.