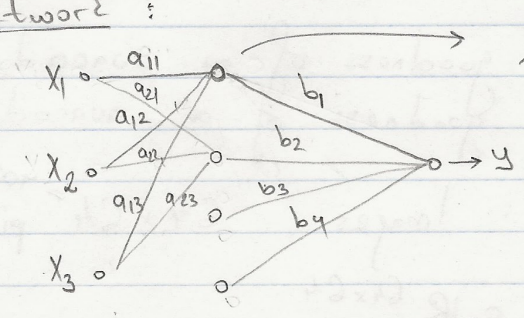


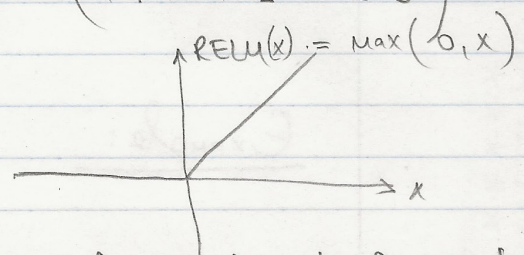
Example: $\mathcal{H} = \{ h_{\mathcal{Q}}(x) = \mathcal{Q}_1 x^n + \mathcal{Q}_2 x^{n-1} + \dots + \mathcal{Q}_n x + \mathcal{Q}_{n+1} \}$
 ↪ set of polynomial predictors of degree $\leq n$.

Example $\mathcal{H} = \{ h_{\mathcal{Q}} : h_{\mathcal{Q}}$ is the set of functions we can implement with a neural network with a given structure and coefficients $\mathcal{Q} \in \mathbb{R}^d \}$

Neural network:



$$h_1 = \text{ReLU}(a_{11}x_1 + a_{12}x_2 + a_{13}x_3)$$



$$y = b_1 h_1 + b_2 h_2 + b_3 h_3 + b_4 h_4$$

$$f: \mathbb{R}^3 \rightarrow \mathbb{R}$$

Loss function: $l: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ non-negative

$l(\hat{y}, y)$: takes two labels and quantifies how different the two labels are.

E.g. $\mathcal{Y} \subseteq \mathbb{R}$ regression problems.

$$l(\hat{y}, y) = (\hat{y} - y)^2 \quad \text{squared loss.}$$

E.g. $\mathcal{Y} \in \{1, 2, \dots, k\}$ discrete set classification problems with k classes.

$$l(\hat{y}, y) = \begin{cases} 1 & \text{if } y \neq \hat{y} \\ 0 & \text{o/w} \end{cases} = \mathbb{1}_{\{\hat{y} \neq y\}}(\hat{y}, y)$$

0-1 loss.

The loss of a predictor $h \in \mathcal{H}$ on a given sample (x, y) will be given by

$$l(h(x), y) \quad \begin{array}{l} \nearrow \text{true label of } x \\ \downarrow \text{the label} \\ \text{predicted by } h \text{ for } x \end{array}$$

Empirical Loss of a predictor h : (the average loss of a predictor h on our sample set $\{(x_i, y_i)\}_{i=1}^n$)

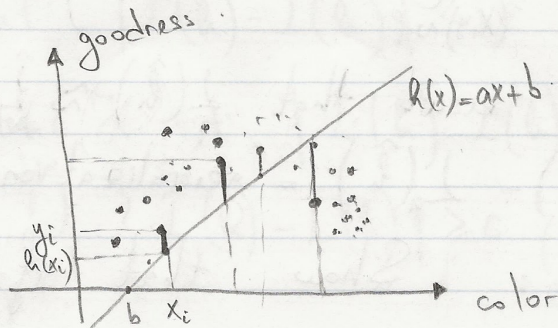
$$L_n(h) = \frac{1}{n} \sum_{i=1}^n l(h(x_i), y_i)$$

Empirical Risk Minimization (ERM):

Find $h \in \mathcal{H}$ that has the smallest $L_n(h)$

$$\hat{h} = \underset{h \in \mathcal{H}}{\operatorname{argmin}} L_n(h) = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n l(h(x_i), y_i)$$

Ex:



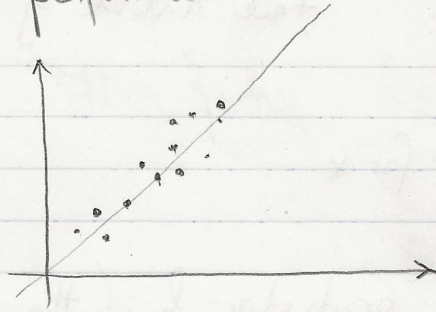
$$\mathcal{H} = \{h_{(a,b)} : h_{(a,b)} = ax + b\}$$

$$l(h(x_i), y_i) = (h(x_i) - y_i)^2 = (ax_i + b - y_i)^2$$

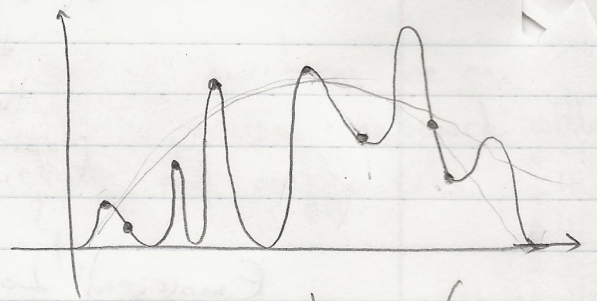
$$\min_{h \in \mathcal{H}} L_n(h) = \min_{a, b} \frac{1}{n} \sum_{i=1}^n (ax_i + b - y_i)^2$$

Linear Least Squares Regression.

How do we know that the performance of a chosen predictor \hat{h} on the data set is representative of its performance in the real world/population?



The data may not be representative



Memorization /
Poor generalization
Overfitting

Data generation model:

$\{(x_i, y_i)\}_{i=1}^n$ " i.i.d $\sim P$ (P is the true ^{unknown} distribution from which nature generates (x, y))

Population loss of a predictor h :

$$L(h) = \mathbb{E}_{(x,y) \sim P} [l(h(x), y)] = \int l(h(x), y) f_{x,y}(x, y) dx$$

How can we ensure that $L(\hat{h}) \approx L_n(\hat{h})$:

$$L(\hat{h}) - L_n(\hat{h}) = \text{generalization error}$$

Uniform convergence: Show that if we

take n i.i.d. samples from P and evaluate $L_n(h)$ for all $h \in \mathcal{H}$, then $L_n(h) \approx L(h)$ for all $h \in \mathcal{H}$ (regardless of P).

Then no matter what $\hat{h} \in \mathcal{H}$ is chosen by the ERM

algorithm, $L_n(\hat{h}) \approx L(\hat{h})$, i.e. the ERM solution will generalize well to the real world.