# Supervised Learning Review

Assume access to training data $(x, y)$

$$X = \{x_1, x_2, \ldots, x_n\}$$

$$= \begin{array}{c} \\ n \end{array} \overset{d}{\begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ \vdots & & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nd} \end{bmatrix}} \qquad X \in \mathbb{R}^{n \times d}$$

$$y = \{y_1, y_2, \ldots, y_n\} \qquad y \in \mathbb{R}^{n \times c}$$

In the simplest case, $c = 1$.

$$y = \begin{array}{c} \\ n \end{array} \overset{c}{\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}} \qquad \text{e.g.} \quad \begin{bmatrix} -1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \qquad \begin{array}{l} x \text{ are our instances} \\ y \text{ are our labels} \end{array}$$

The goal is to find a model with parameters $\Theta$, $f_\Theta$, s.t. $f_\Theta(x_i) = y_i$.

The real hope is that $f_\theta$ generalizes to unseen or held-out data.

We estimate this via a TEST SET $(x', y')$ which might contain $\ell$ examples.

In practice we usually use three datasets



Train — estimate $\theta$

dev — explore models & hyperparams

Test — eval. final performance.

!!! Do not touch your test data !!!

The idea behind learning is to find $\hat{\theta}$ s.t. the loss on the Train data is small. formally

Training objective

$$\underset{\hat{\theta}}{\arg\min} \sum_{i=1}^{N} \mathcal{L}_{\hat{\theta}}(x_i, y_i)$$

$$= \mathcal{L}(\underbrace{f_{\hat{\theta}}(x_i)}_{\hat{y}_i}, y_i)$$

This requires (at very least) specifying $f_\theta$

A classic model is the Perceptron



$W_0 + W_1 \cdot x_{i1} \ldots$
$+ W_n \cdot x_{in}$
$= \sum_j^n W_j x_{ij}$
$= w \cdot x_i$
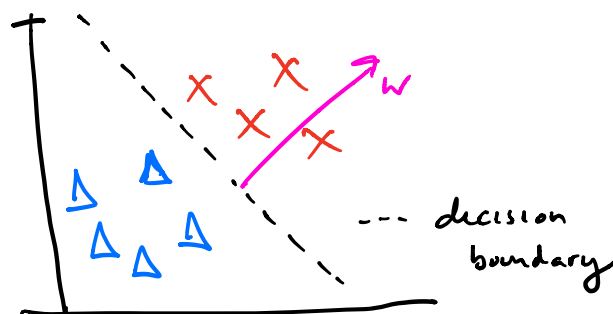
Here $\theta = W$ and a prediction is made:

$$\hat{y}_i = \begin{cases} 1 & \text{if } W \cdot x_i > \tau \\ -1 & \text{otherwise} \end{cases}$$

$W$ defines a decision boundary (line or plane) that delineates examples from $y_i = 1 / y_i = -1$

This is perpendicular to $W$; $u \cdot v = 0$ iff $u$ is perp. to $v$.



--- decision boundary

So here the natural loss is Zero/one loss

$$\mathcal{L}(\hat{y}_i, y_i) = \begin{cases} 0 & \text{if } y_i == \hat{y}_i \\ 1 & \text{otherwise} \end{cases}$$

We need an estimation procedure to find a good $\hat{\theta}$.

```
fitPerceptron (X, Y, MAX_ITERS)
    W ← ∅⃗        // init weights to Zero-vector
    for max_iters
        for (x_i, y_i) ∈ X, Y
            a ← W · x_i        // activation for x_i
            if y_i · a ≤ 0     // we made a mistake
                W ← W + y_i · x_i
    return W
```

Let's see this in Colab.