

Perceptron

Binary Classification

$(x_1, y_1), (x_2, y_2), \dots$

feature
vec
 $x_i \in \mathbb{R}^d$

Label

$y_i \in \{\pm 1\}$

$$h(x) = y$$

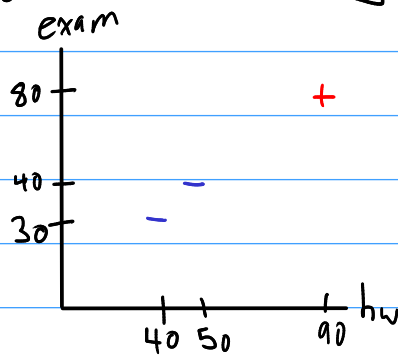
$$x_i = (\text{hw}, \text{exam})$$

$$y_i = \text{Passed?}$$

$$x_1 = (90, 80), y_1 = 1$$

$$x_2 = (40, 30), y_2 = -1$$

$$x_3 = (50, 40), y_3 = -1$$



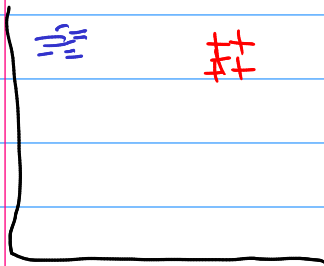
$$x = (50, 60) \quad y = ?$$

Statistical Learning

$(x_1, y_1), \dots, (x_n, y_n) \stackrel{\text{i.i.d.}}{\sim} P$

Goal: Learn $h(\cdot): \mathbb{R}^d \rightarrow \{\pm 1\}$ s.t. $\Pr_{(x,y) \sim P}[h(x) = y]$ is large

Bayes classifier



Online Learning.

At each time $i = 1, 2, \dots$.

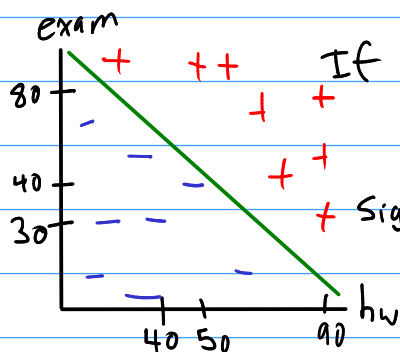
1. Receive x_i .
2. Choose h_i , predict $\hat{y}_i = h_i(x_i)$.
3. View y_i . Suffer mistake if $y_i \neq \hat{y}_i$.

Perceptron

$$x_1 = (90, 80), y_1 = 1$$

$$x_2 = (40, 30), y_2 = -1$$

$$x_3 = (50, 40), y_3 = -1$$



If average of hw + exams > 0.5 then pass

If $0.5 \cdot hw + 0.5 \cdot exam - 0.5 > 0$, then pass

$$\text{sign}(\langle (0.5, 0.5), (hw, exam) \rangle - 0.5)$$

$$x_i = (hw, exam)$$

$$y_i = \text{sign}(\langle (0.5, 0.5), x_i \rangle - 0.5)$$

Algorithm: The Perceptron (Rosenblatt 1958)

Input: Dataset $\mathcal{D} = \{(x_i, y_i) \in \mathbb{R}^d \times \{\pm 1\} : i = 1, \dots, n\}$, initialization $w \in \mathbb{R}^d$ and $b \in \mathbb{R}$, threshold $\delta \geq 0$

Output: approximate solution w and b

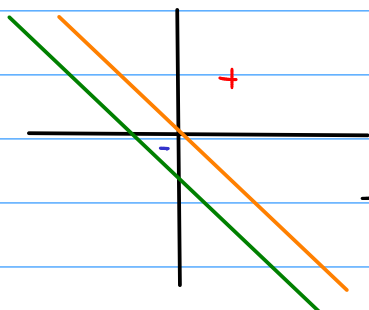
```

1 for  $k = 1, 2, \dots$  do
2   receive training example index  $I_k \in \{1, \dots, n\}$  // the index  $I_k$  can be random
3   if  $y_{I_k}(w^\top x_{I_k} + b) \leq \delta$  then
4      $w \leftarrow w + y_{I_k} x_{I_k}$  // update only after making a 'mistake'
5      $b \leftarrow b + y_{I_k}$ 

```

$$x_1 = (1, 1), y_1 = 1$$

$$x_2 = (-\frac{1}{4}, -\frac{1}{4}), y_2 = -1$$



$$w = \vec{0}, b = 0, \delta = 0$$

$$w \leftarrow w + (1, 1) = (1, 1)$$

$$b \leftarrow b + 1 = 1$$

$$+ (-\frac{1}{2} + 1) = -\frac{1}{2} < 0$$

$$w \leftarrow w + (-\frac{1}{4}, -\frac{1}{4})(-1) = (\frac{5}{4}, \frac{5}{4})$$

$$b \leftarrow b + (-1) = 1 - 1 = 0$$

Padding + Pre-multiplication

$$y_i = \text{sign}(\langle w, x_i \rangle + b) \quad \forall i \in [n]$$

$$= \text{sign}(\langle (w \ b), (x_i \ 1) \rangle)$$

$$\Leftrightarrow y_i \langle (w \ b), (x_i \ 1) \rangle > 0$$

$$(x_i, y_i) \Rightarrow a_i \triangleq y_i (x_i \ 1)$$

$$\langle a_i, w \rangle > 0 \quad \forall i$$

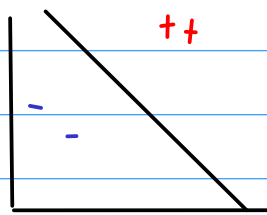
$$Aw > \vec{0} \text{ (entrywise)}$$

$$A = \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix} \in \mathbb{R}^{n \times (d+1)}$$

Linear Separability

$$\exists w \text{ s.t. } \langle a_i, w \rangle \geq s > 0 \quad \forall i$$

$$\Leftrightarrow Aw \geq s \vec{1}, \text{ where } s > 0$$



Error Bound + Margin

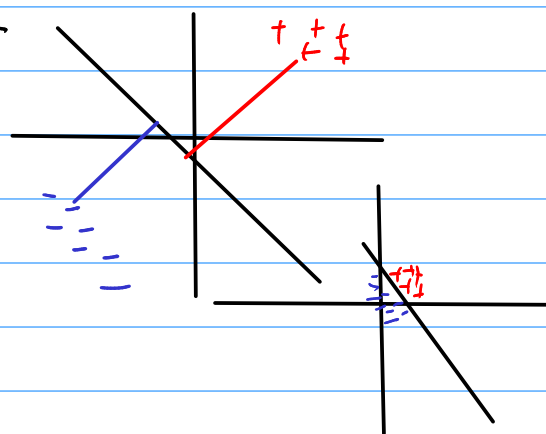
Suppose $\exists w, s > 0$ s.t. $Aw \geq s \vec{1}$.

Then perceptron converges in $\leq R^2 \cdot \left(\frac{\|w\|_2^2}{s^2}\right)$ steps, $R = \max_i \|a_i\|_2$.

$$\min_{(w,s): Aw \geq s \vec{1}} \frac{\|w\|_2^2}{s^2} = \min_{(w,s): \|w\|_2=1, Aw \geq s \vec{1}} \frac{1}{s^2}$$

$$= \frac{1}{\left(\max_{(w,s): \|w\|_2=1, Aw \geq s \vec{1}} s\right)^2} = \frac{1}{\left(\max_{\|w\|_2=1} \min_i \langle a_i, w \rangle\right)^2}$$

margin γ



Uniqueness?
No. SVM

Non-Separable?

- Perceptron cycles
- Not the right algo.

When to stop?

- all pts are correct
- error stops decreasing
- weights converge
- fixed # of iters

Multiclass $(k=2 \rightarrow \text{binary})$
- one-vs-all - Train k classifiers \leftarrow # of classes
- one-vs-one - k^2 classifiers \leftarrow Pick largest $\langle w_i, x \rangle + b$