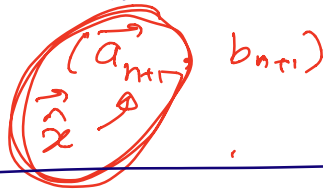Today
- Machine Learning terminology.
- Classification continued.

---

Terminology

Training data :  Data to help you learn (your classifier,
predictor etc.)  e.g $\cdot A\vec{x} = \vec{b}$
$$\Downarrow \;\; (\vec{a_i}, b_i)$$

Test data :  Data to check/test if what you
learned is any good.
$$(\vec{a}_{n+1}, b_{n+1})$$
$\vec{x}$

---

Classification

Initial data :   $\vec{x_1}, \vec{x_2} \cdots \vec{x_m}$   (m data points)

→ pixels of an image
→ observations of a planet.

Associated label:   $l_1, l_2 \cdots l_m$
Binary classification      $l_i \in \{+1, -1\}$

$\{+1\}, \quad \{-1\}$
cat          dog
Neuron 1      Neuron 2.

$$\left\{ (\vec{x_1}, \ell_1) \; , \; (\vec{x_2}, \ell_2) \; \cdots \; (\vec{x_m}, \ell_m) \right\}$$

Find a classifier. In particular, we want to find a **linear classifier**.
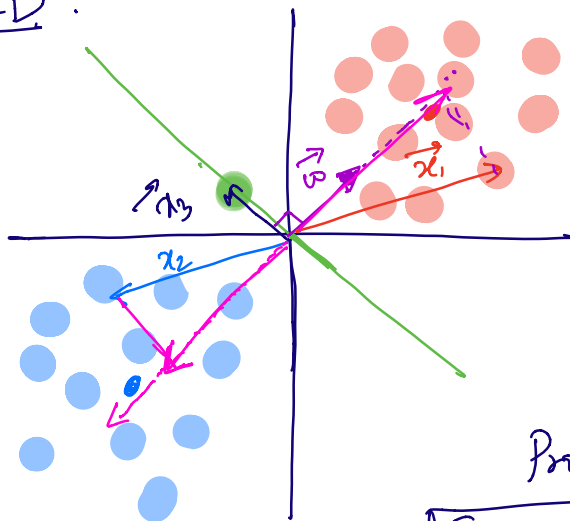
$$f(\vec{x_i}) \longrightarrow \ell_i$$

We want to find a vector $\vec{w}$ such that:

$$\text{sign}\left( \vec{x}^T \vec{w} \right)$$

$$\text{sign}(x) = +1 \qquad \text{if } x > 0$$
$$\text{sign}(x) = -1 \qquad \text{if } x < 0$$

In 2D :



$\|\vec{w}\| = 1$
Find $\vec{w}$ such that
$\vec{x_i}^T \vec{w} > 0$
if $\vec{x_i} \in \{ \text{Red} \}$

$\vec{x_i}^T \vec{w} < 0$ if
$\vec{x_i} \in \{ \text{Blue} \}$

Proj of $\vec{x}$ onto $\vec{w}$ : $\boxed{\vec{x}^T \vec{w}} \times$
$\dfrac{}{\|\vec{w}\|^2} = 1$

Consider $\vec{x_i}^T \vec{w} > 0$ : [Signed Magnitude] of projection of $\vec{x_i}$ onto $\vec{w}$

Consider $\vec{x_2}^T \vec{w} < 0$ : true for blue points

---

Goal: To find $\vec{w}$. How?

"Cost-function" → penalty for being wrong.

$$\vec{x_i}^T \vec{w} \longrightarrow l_i$$

One possible cost function:

$$\underset{\vec{w}}{\arg\min} \sum_{i=1}^{m} (\underbrace{\vec{x_i}^T \vec{w}}_{} - l_i)^2$$

→ if $\vec{x_i}^T \vec{w} = l_i$ = good.

$\vec{x_i}^T \vec{w} \neq l_i$ = bad.

$$= \underset{\vec{w}}{\arg\min} \left\| \begin{bmatrix} -\vec{x_1}^T- \\ -\vec{x_2}^T- \\ \vdots \\ -\vec{x_m}^T- \end{bmatrix} \begin{bmatrix} \vec{w} \end{bmatrix} - \begin{bmatrix} l_1 \\ l_2 \\ \vdots \\ l_m \end{bmatrix} \right\|^2$$

Least squares.

Nothing is special about $(\vec{x}_i^T \vec{w} - l_i)^2$.

## General cost function:

$$C\left(\vec{x}_i^T \vec{w}, l_i\right)$$

$\vec{x}_i^T \vec{w}$ has
the same sign
as $l_i$
→ good.
wont $C(\cdot)$ to
be small.

$\vec{x}_i^T \vec{w}$ has opposite
sign as $l_i$
This is bad.
$C(\cdot)$ to be large.

## We chose: $C(\vec{x}_i^T \vec{w}, l_i) = \exp\left(-l_i \, \vec{x}_i^T \vec{w}\right)$

When $\text{sign}(\vec{x}_i^T \vec{w}) = l_i$

$\quad \exp(\text{negative}) \longrightarrow$ small.

$\quad \text{sign}(\vec{x}_i^T \vec{w}) \neq l_i \quad \leftarrow$ error
$\quad \exp(\text{positive}) \longrightarrow$ big $\qquad$ High cost function!

$$\underset{\vec{w}}{\text{argmin}} \sum_{i=1}^{m} \exp\left(-l_i \, \vec{x}_i^T \vec{w}\right)$$

Our strategy: Make this cost function look like a quadratic. $\vec{w} \in \mathbb{R}^n$

Taylor approximation:

$$f(\vec{w}) \approx f(\vec{w_*}) + \underbrace{\frac{df}{d\vec{w}}\bigg|_{\vec{w} = \vec{w_*}}}_{\text{row vector Derivative}} (\vec{w} - \vec{w_*})$$

$$+ \frac{1}{2}(\vec{w} - \vec{w_*})^T \underbrace{\frac{d^2 f}{d\vec{w}^2}\bigg|_{\vec{w} = \vec{w_*}}}_{\substack{\text{matrix} \\ \text{Hessian}}} (\vec{w} - \vec{w_*})$$

$$= f(\vec{w_*}) + \left[ \frac{\partial f}{\partial w_1} \cdots \frac{\partial f}{\partial w_n} \right]\bigg|_{\vec{w} = \vec{w_*}} (\vec{w} - \vec{w_*})$$

$$+ \frac{1}{2}(\vec{w} - \vec{w_*})^T \begin{bmatrix} \frac{\partial^2 f}{\partial w_1^2} & \cdots & \frac{\partial^2 f}{\partial w_1 \partial w_n} \\ & \vdots & \\ \frac{\partial^2 f}{\partial w_n \partial w_1} & \cdots & \frac{\partial^2 f}{\partial w_n^2} \end{bmatrix} (\vec{w} - \vec{w_*})$$

To find our quadratic form, we need an operating point $\vec{w_*}$.

But to find our operating point, we need a quadratic form!?

Solution: Consider an iterative algorithm.
$$\left(\text{Newton's method.}\right) \longrightarrow \text{roots of a polynomial.}$$

## Algorithm:

① Arbitrarily choose an operating point
$$\vec{w_*} = \vec{w}[0] = \vec{0}$$

② Quadraticize around $\vec{w_*}$
$$f(\vec{w}) = \sum_{i=1}^{m} \underbrace{C(\vec{x_i}^T \vec{w}, l_i)}_{\text{cost function}} \quad \text{around } \vec{w_*}$$

$$f(\vec{w}) \approx \underset{\uparrow}{\vec{w}^T A \vec{w}} + \underset{\uparrow}{\vec{b}^T \vec{w}} + \underset{\uparrow}{d} \quad \left(\text{Generic form} \atop \text{of quadratic}\right)$$
$$\qquad\qquad \text{matrix} \qquad\quad \text{vector} \qquad \text{scalar}$$

③ Find the minimizer of the quadratic.
Call this $\vec{w}[1]$

④ Set $\vec{w_*} = \vec{w}[1]$, and go back to ②.

Quadratic approx

$$f(\vec{w}) \approx$$
$$f(\vec{w_*}) + \left[\frac{\partial f}{\partial w_1} \cdots \frac{\partial f}{\partial w_n}\right]\Big|_{\vec{w}=\vec{w_*}} (\vec{w}-\vec{w_*})$$

$$+ \frac{1}{2}(\vec{w}-\vec{w_*})^T \begin{bmatrix} \frac{\partial^2 f}{\partial w_1^2} & \cdots & \frac{\partial^2 f}{\partial w_1 \partial w_n} \\ \vdots & & \\ \frac{\partial^2 f}{\partial w_n \partial w_1} & \cdots & \frac{\partial^2 f}{\partial w_n^2} \end{bmatrix} (\vec{w}-\vec{w_*})$$

$$f(\vec{w}) = \sum_{j=1}^{m} \exp(-l_i \vec{x_i}^T \vec{w})$$

Consider: Partial of
$$\frac{\partial (\exp(-l_i \vec{x_i}^T \vec{w}))}{\partial w_1}$$

$$= \frac{\partial (\exp(-l_i (x_1 w_1 + x_2 w_2 \cdots + x_n w_n)))}{\partial w_1}$$

$$= -l_i x_i \exp(-l_i \vec{x_i}^T \vec{w})$$

---

## Quadratic approximation:

$$\sum_{i=1}^{m} c(\vec{x_i}^T \vec{w}, l_i) \approx$$

$$\sum_{i=1}^{n} \left[ c(\vec{x_i}^T \vec{w_*}, l_i) \underset{\nearrow}{\overset{\text{constant}}{}} - l_i \exp(-l_i \vec{x_i}^T \vec{w_*}) \vec{x_i}^T (\vec{w} - \vec{w_*}) \right.$$

$$\left. + \frac{1}{2} \exp(-l_i \vec{x_i}^T \vec{w_*}) \langle \vec{x_i}, \vec{w} - \vec{w_*} \rangle^2 \right]$$