_/_/_

# MACHINE LEARNING ALGORITHMS.

No need to explicitly program, the machine learns by itself.

- Supervised learning: Huge, labelled data.
- Unsupervised learning: Clustering; unlabelled data.
- Semi-Supervised learning: Some are labelled, some not.
- Reinforcement learning: Carrot + stick method ?.
- Deep (Generative) learning: Language model, generate something, layers of Neural Nets.
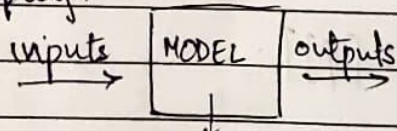
① Supervised Learning:

Regression problem: Outputs a real number.
Classification problem: Binary or Multiclass.

$\mathcal{H}$ - Curly H: All possible hypothesis $(g)$.
$h \in \mathcal{H}$, $h$ is the of needed hypothesis.

31 Dec  Mapping:

inputs → | MODEL | outputs →

Mathematical model:
family of math equations.

Training a model: finding parameters that predict outputs 'well' from inputs for a training dataset of input/output pairs.

**Parameters:** $y = mx + b$.

Here, $m$ & $b$ influence the relationship b/w $y$ & $x$.

**Regression:** Finding the best possible relationship b/w $x$ & $y$.

**Need for a generalized model:** a model that isn't made specifically for the dataset alone, but works for more data with good results.

* Capital bold $X$ - matrix.
* Bold variable - ~~scaler value~~ vector.
* normal variable - scalar value; single value.

* Parameters: $\phi$ ; Model: $y = f[x_i, \phi]$.

* **Loss Function:**
Measures how bad the model is.

Find the parameter that minimize the loss.

$$\hat{\phi} = \underset{\phi}{\text{argmin}} \, [L[\phi]].$$

**Testing a model:** Separate test data with input/output pairs.
↳ checks how well it generalizes.

**Supervised:** mapping from one input to multiple output [ classification model].
mapping from multiple input to one output [ regression model].

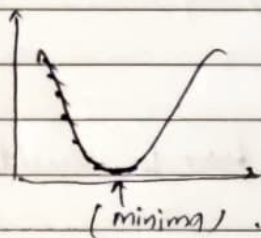$$\bullet \; L[\phi] = \sum_{i=1}^{I} (\rho_0 + \theta_1 x_i - y_i)^2.$$
$\quad\hookrightarrow$ least square loss function.

Gradient Descent Algorithm:
Take / Compute the slope at a point;
Equate to 0.
Check next point, if it is greater than
current point, current point is a minimum.


(minima).

16 jan.

* Hypothesis:

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 \dots$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad - \text{parameters / weights.}$$

$$h(x) = \sum_{i=0}^{n} \theta_i x_i \quad .$$

* Cost Function:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^{n} (h_\theta(x_i) - y_i)^2.$$

Take house price dataset and write code for
Gradient Descent with only 2 features, in
python [Find convergence rate, experiment with]
different learning rate.

☆ Gradient Descent Algorithm :

$$\theta_j := \theta_j - \alpha \nabla_\theta J(\theta).$$

where $\nabla = \dfrac{\partial}{\partial \theta_j}$.

$\alpha$ = learning rate.

$$\dfrac{\partial}{\partial \theta_j} J(\theta) = \dfrac{\partial}{\partial \theta_j} \dfrac{1}{2}\left(h_\theta(x) - y\right)^2.$$

[Partial deriv cause $\theta_j$ has it multivariate].

$$\dfrac{\partial}{\partial \theta_j} J(\theta) = 2 \cdot \dfrac{1}{2} \left(h_\theta(x) - y\right) \dfrac{\partial}{\partial \theta_j} \left(h_\theta(x) - y\right).$$

$$= \left(h_\theta(x) - y\right) \dfrac{\partial}{\partial \theta_j} \left(\sum_{i=0}^{n} \theta_i x_i - y\right).$$

$$\nabla_\theta J(\theta) = \left(h_\theta(x) - y\right) x_j.$$

in case → Multiple variables, then it is Hessian matrix.

$$\Rightarrow \theta_j := \theta_j + \alpha \left(y_i - h_\theta(x_i)\right) x_j.$$

↳ LHS - update rule: Least Mean Square Update Rule.
Also called : Widrow - Hoff Learning Rule.

Repeat until convergence; when the error term
$(y_i - h_\theta(x_i)) x_j$ does not change, hence $\theta_j$ does
not change.

\* Logistic Regression:

Email classifier [Binary]:
$$y = \{Spam, Ham\}.$$

$$h_\theta(x) = g(\theta^T x).$$

$$= \left[ \frac{1}{1 + e^{-\theta^T x}} \right]. \quad \text{(sigmoid)}.$$

$$g(z) = \frac{1}{1 + e^{-z}}.$$

$g(z)$ is always bounded between 0 & 1.

$$g'(z) = \frac{d}{dz} \frac{1}{1 + e^{-z}}. \qquad z \text{ vs } z^{\wedge\wedge?}$$

$$= \frac{1}{(1 + e^{-z})^2} e^{-z}.$$

$$= \left[ \frac{1}{(1 + e^{-z})} \right] \left[ 1 - \frac{1}{(1 + e^{-z})} \right].$$

$$g'(z) = g(z)(1 - g(z)).$$

Assume that:
$$P(y=1 \mid x;\theta) = h_\theta(x).$$
$$P(y=0 \mid x;\theta) = (1 - h_\theta(x)).$$

$$\boxed{P(y \mid x;\theta) = [h_\theta(x)]^y (1 - h_\theta(x))^{1-y}}$$

Likelihood:
$$L(\theta) = P(y \mid x;\theta)$$
$$= \prod_{i=1}^{n} P(y^i \mid x^i;\theta).$$

HW: Compare the standard stochastic descent rule of linear regression and stochastic ascent rule of logistic regression. What similarities do you observe? Is it a coincidence? ___/___/___

$$L(\theta) = P(y \mid x;\theta)$$

$$= \prod_{i=1}^{n} P(y_i \mid x_i ; \theta).$$

$$= \prod_{i=1}^{n} (h_\theta(x_i))^{y_i} (1 - h_\theta(x_i))^{1-y_i}$$

$$=$$

## log likelihood:

$$\ell(\theta) = \log L(\theta).$$

$$\ell(\theta) = \sum_{i=1}^{n} y_i \log(h_\theta(x_i)) + (1-y_i) \log(1 - h_\theta(x_i)).$$

## Maximize the likelihood:

$$\frac{\partial \ell(\theta)}{\partial \theta_j} = \left[ y \frac{1}{g(\theta^T x)} - (1-y)\frac{1}{1-g(\theta^T x)} \right] \frac{\partial}{\partial \theta_j} g(\theta^T x).$$

$$y \in \{0,1\}. \Rightarrow h_\theta(x) = g(x) = g(\theta^T x).$$

$$\frac{\partial \ell(\theta)}{\partial \theta_j} = \left[ \quad \right] g(\theta^T x)(1-g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x .$$

$$= (y(1-g(\theta^T x)) - (1-y) g(\theta^T x)) x_j$$

$$= (y - h_\theta(x)) x_j.$$

## Stochastic Gradient Ascent

$$\theta_j := \theta_j + \alpha (y_i - h_\theta(x_i)) x_j$$

$$\left[ \text{from } \theta = \theta + \alpha \nabla_\theta \ell(\theta) \right].$$

\* Bayesian Decision Theory:

Classification : class $\omega_1$ , class $\omega_2$ .

Accept product.      Reject product.

Decision Rule:

REJECTED

$P(W_1) > P(\omega_2) \Rightarrow \omega_1$.     sometimes good, but

$P(W_1) < P(\omega_2) \Rightarrow W_2$.     not always.

$P(\omega_1), P(\omega_2)$ are called Apriori probabilities.
Decision is based on some feature $x$.

Probability density function of variable $x$ given $\omega_1$ or $\omega_2$:

$P(x|\omega_1)$ ; $p(x|\omega_2)$. $\rightarrow$ Class conditional probabilities

Decision is based on: $p(\omega_1|x)$ ; $p(\omega_2|x)$.

Joint probability: $P(\omega_i, x) = P(\omega_i|x)P(x)$. or $\begin{bmatrix} \text{Product} \\ \text{rule} \end{bmatrix}$.

$P(x|\omega_i)P(\omega_i)$.

$$\Rightarrow \boxed{P(\omega_i|x) = \frac{P(x|\omega_i)P(\omega_i)}{P(x)}}$$

where $p(x) = \sum_i P(x|\omega_i)P(\omega_i)$.

Decision Rule:

$p(\omega_1|x) > p(W_2|x) \Rightarrow \omega_1$.

$P(\omega_1|x) < P(\omega_2|x) \Rightarrow \omega_2$.

Exam portion: Moodle - First 4 chapters of textbook.

Exam time 9-11?? Why not 10-12?          ¿Adat se mjbo?

                                          __/__/__

$$P(x|w_i) P(w_i) > P(x|w_2) P(w_2) \Rightarrow w_1.$$

Here, the decision is made based on apriori probabilities and class conditional probabilities.

6 Feb

* Parametric methods:

- Maximum Likelihood estimation (MLE):

Baye's Theorem:   $P(\theta|D) \propto P(D|\theta) P(\theta)$
                    posterior      likelihood    → joint?

$P(X|\theta)$ ← likelihood

$$l(\theta|x) = P(x|\theta) = \prod_{t=1}^{N} P(x_t|\theta).$$

$X = \{x_t\}_{t=1}^{N}$  ← Draw some (N) data from a known distribution $P(x|\theta)$.

$x_t \sim P(x|\theta)$.   ;  $\theta = \begin{bmatrix} \mu \\ \sigma \end{bmatrix}$.  $\theta_1 = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_n \end{bmatrix}$ and $\theta_2 = \begin{bmatrix} \sigma_1 \\ \vdots \\ \sigma_n \end{bmatrix}$.

log likelihood function:

$$L(\theta|X) = \log l(\theta|X)$$
$$= \sum_{t=1}^{N} \log P(x_t|\theta).$$

Bernoulli Distribution:          Ber(P).        ; Binomial(n, p).
$$P(x) = p^x (1-p)^{1-x}$$     ;  $x \in \{0,1\}$.

$$L(P|X) = \log \prod_{t=1}^{N} P(x_t^i)(1-p)^{(1-x_t^i)}.$$

$$= \sum_{t=1}^{N} x_t \log p + (N - \sum_t x_t) \log(1-p).$$

Maximize : Estimate $\hat{P} = $ ... &

$$\hat{P} = \frac{\sum_t x^b_t}{N}$$

It is called estimate.

* Bias and Variance :

$\quad X \quad f \quad \theta.$ $\qquad$ Bias of a graph distribution, not neural net.

Estimator of $\theta$ ← $d = f(x) \cdot \longrightarrow \hat{P}(\theta)$.

Error $= (d(x) - \theta)^2$.

$$V(d, \theta) \qquad = E\left[(d(x) - \theta)^2\right].$$
Mean Square Error.

Bias of an estimator : $b_\theta(d) = \underline{E\left[d(x)\right] - \theta}$.

What does unbiased estimator mean ?
For all $\theta$, $b_\theta = 0$, we call it unbiased estimator.
$\quad\longrightarrow$ number of samples.

$$E[m] = E\left[\frac{\sum_t x_t}{N}\right] = \frac{1}{N} \frac{\sum_t x_t}{N} = \frac{N\mu}{N} = \mu.$$

$\quad\longrightarrow$ concept used similar to.

Law of large numbers : When sample size ↑, sample
$\qquad\qquad\qquad\qquad$ mean $\leadsto$ actual mean.

Variance : $Var(m) = Var\left[\frac{\sum_t x_t}{N}\right] = \frac{\sigma^2}{N}$.

$Var = E[x^2] - (E[x])^2$.

Mean Square Error:

$$v(d,\theta) = E\{[d(x)-\theta]\}^2$$

$$v(d,\theta) = E[(d(x)-\theta)^2].$$

$$= \underbrace{E[(d-E[d])^2]}_{Variance} + \underbrace{(E(d)-\theta)^2}_{bias}.$$

$$r(d,\theta) = Var(d) + b_\theta(d)^2.$$

$$Error = Variance + bias^2$$