

1.1 Introduction to Machine Learning :-

Questions: * Is AI/ML just a fad?

- New innovation wave
- Automated pattern recog.
- Automated decision making.

Input + labels \rightarrow Model | Input + Model \rightarrow labels

AI \geq ML \geq Neural Networks \geq DL

* Why ML now? : - Lots of data - Lots of computatⁿ - New frameworks & alg.

When ML not useful : - Unsupervised learning - Too little data - Model is too simple

Sweet spot for ML : - Lots of structure - Explainability not critical - Prediction accuracy \gg but stability - Complex but stationary

* What is model complexity : Linear + non linear + different inputs + fn of fn.

* Other considerations: Memory Transferability Computation Speed & parallelization

* Life stages : Need identification \rightarrow Data gathering \rightarrow Model choice \rightarrow Train valid \rightarrow Deployment \rightarrow Monitoring

Directions in ML research : - Train with unlabelled data - High level labelling
- Train with fewer labelled data - One task to another

- Explaining decisions - More cautious in recognising new scenarios.

* Myth or Reality? - M, M, M, M, R, M, M, M, M

Takeaways: - Not Silver Bullet
- Still useful - Data is currency of ML
- Critically design, question & then accept ML

1.2 ML for smart monkeys :-

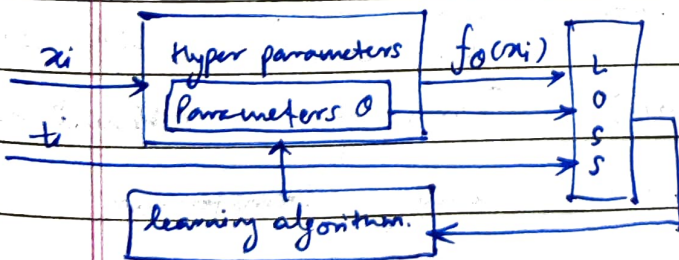
[High level recipe for a machine learning concept or model.]

Elements of a model : Input x_i and function $f_\theta(x_i)$

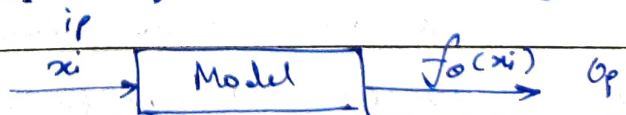
Utility of a model : Target output $t_i \rightarrow$ Bring $f_\theta(x_i) \rightarrow t_i$

Model

Aim to do this by minimizing $L(t_i, f_\theta(x_i), \theta)$
(Optimization)



Training phase.



(Testing phase)

Examples of hyper parameters : Random forest - No. of nodes & tree depth
Neural Network - No. of nodes in layers

Recipe for ML training:

- ① Decide on type of ML problem
- ② Prepare data
- ③ Shortlist ML frameworks
- ④ Prepare training, test & validation sets
- ⑤ Train, validate, repeat
- ⑥ Test ONCE

① Broad types of ML problems:

Supervised : x_i, t_i both given (Most advances are here)

Unsupervised : x_i only given.

Supervised ML using classification / ordinal / Continuous (Ranking) / Regression or range

Unsupervised classification / Dimensional redⁿ (clustering) (maybe PCA?)

② Preparing the data:

- Remove useless data (No variance / Not available) ↘ Maybe already depends on values
 - Example (How many days spent in hospital already depends on given patient)
- Reduce Redundancy (Get rid of copies of highly correlated variables)
- Handle missing data (Impute if sporadic, drop if too frequent)
- Transform variables (Convert discrete to one hot bit and normalise continuous variables)

③ Popular ML frameworks

Vector :	Logistic regression	Linear regression	Kmeans	PCA x-PCA
(Basic data)	SVM, RF, NN		Fuzzy Cmeans DBSCAN	tLE, ISOMAP
	(Classification)	(Regression)	(Clustering)	(Dimensional reduction)
Series, text :	Recurrent NN, long short-term memory, Transformer,			
	1D CNN, Hidden Markov Models			
Images :	2D CNN, Markov Random fields	Probabilistic graphical models		
Video :	3D CNN, CNN + LSTM, MRF			

2.4.) Pseudo-inverse of a Matrix ($X \rightarrow X^+$)

$$(XX^+)X = X \quad [\text{like identity but not identity}]$$

PAGE No.	
DATE	/ /

and $X^+ = (X^H X)^{-1} X^H$ (Conjugate transpose = X^H)

Eigen decomposition \rightarrow Square matrix

$$AV = \lambda V, \quad V \text{ is a unit vector (we assume this)} \quad \lambda_i \text{ is a scalar.}$$

$N \times N \quad N \times 1$

Then we can get $A = Q \Lambda Q^{-1}$

$$\lambda_1 > \lambda_2 > \lambda_3 \dots \text{ for } \Lambda \quad (Q = [v_1 \ v_2 \ v_3 \dots v_N])$$

Tensors $\alpha \in \mathbb{R}$ and $x \in \mathbb{R}^N$, $A \in \mathbb{R}^{N \times M}$ Tensor: $T \in \mathbb{R}^{M \times N \times P \times Q}$

Tensor transpose needs order of permutation of dimensions.

2.5) Functions $f: X \rightarrow Y \quad x \in X \quad f(x) \rightarrow y$

\hookrightarrow May be many to one but CANNOT be one to many. (obviously)

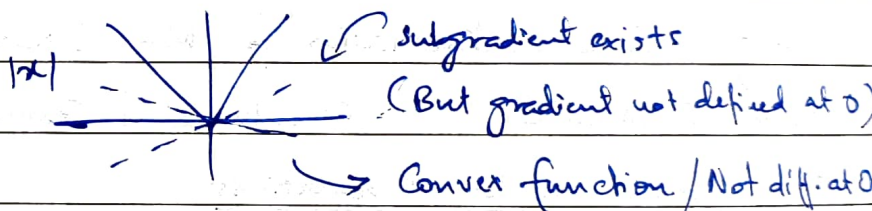
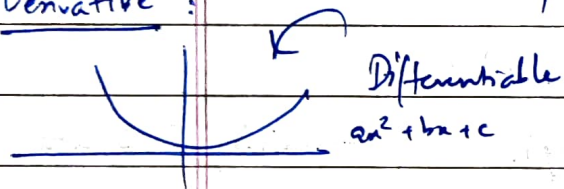
Continuity: x is continuous: $\lim_{\Delta x \rightarrow 0} f(x + \Delta x) = f(x)$

\hookrightarrow Lipschitz continuity

\hookrightarrow Eg. x^2 - Not Lipschitz cont.

Derivative:

$$|f(x_1) - f(x_2)| \leq K |x_1 - x_2| \quad \exists K \text{ such that}$$

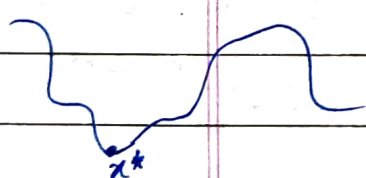


$$\lambda f(x_1) + (1-\lambda)f(x_2) \geq f(\lambda x_1 + (1-\lambda)x_2) \quad \triangleq \text{Convex function.}$$

Critical points: $f'(x) = 0$ The x which satisfy this. \hookrightarrow Maxima / Minimal Inflection point.

\hookrightarrow Maxima: $f''(x) < 0$ | $f''(x) > 0$: Minima | $f'''(x) = 0$ Inflection point.

2.6) Gradient Descent and ascent:



Reach local minima x^* We can compute $f(x)$ & $f'(x)$

If $f'(x) > 0$ - Move left / $f'(x) < 0$ - Move right

$$x_1 = x_0 - \eta f'(x_0) \quad \text{or} \quad x_n = x_{n-1} - \eta \nabla f(x_n)$$

If η too small \rightarrow Too many steps needed | η too large \rightarrow overshoots.

$$x_{n+1} = x_n - \frac{f'(x_n)}{f''(x_n)} \quad \hookrightarrow \text{If second derivative pretty deep / constant over long range.}$$

2.7) Multivariate functions

eg. $z = ax^2 + by^2$ Contour plots
Gradient (Derivative of Multivariate)
 Partial derivatives taken as a vector



$$\nabla = \left[\frac{\partial}{\partial x_1} \frac{\partial}{\partial x_2} \dots \frac{\partial}{\partial x_n} \right]$$

or $f(x_1, x_2, \dots, x_n) = y$ so

$$\begin{bmatrix} \vec{x} \end{bmatrix} = \begin{bmatrix} \vec{x} \end{bmatrix} - \eta \begin{bmatrix} \nabla f \end{bmatrix}$$

gradient descent

Newton's method for multi variate fn:

of $f(x_1, x_2)$ and $\nabla f = \left[\frac{\partial f}{\partial x_1} \frac{\partial f}{\partial x_2} \right]$

The 2nd derivative: H (Hessian) =
$$\begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix}$$
 off diagonal terms are same.

for $f(x_1, \dots, x_n)$ $H_{ij} = \left[\frac{\partial^2 f}{\partial x_i \partial x_j} \right]$ - So an $n \times n$ matrix.

for eg. $f = ax^2 + bx^2 \rightarrow H = \begin{bmatrix} 2a & 0 \\ 0 & 2b \end{bmatrix}$ - Eigen vectors tell
 are Eigen value - There is a minime
 are Eigen value - There is a maxime (Relate to 2nd derivative test)

Update: $\vec{x} = \vec{x} - H^{-1} \nabla f$ order N^3 (Newton method)

Constrained Optimisation Using Lagrange Multiplier:

eg. Maximize $f(x)$, subject to $g(x) = 0$ (constraint)

How to find the tightest? \Rightarrow When the contours of f are tangential to $g(x)$.

$\therefore L(x) \equiv f(x) + \lambda g(x)$

Maximize this subject to $\lambda \neq 0$ $\nabla f(x) + \lambda \nabla g(x) = 0$

\rightarrow The derivative (gradient) of $f(x)$ & $g(x)$ are parallel to each other
 Precise condition for having same tangent or same normal equiv.

2.8) Probability Theory (Basics)

1. Random Variable ($X=x$) x is a value that X can take

eg. $X \in \{0, 1\}$ (heads or tails toss)

Probability Mass function $P(X=x) = P_X(x)$ - fn called prob. mass fn

$\sum_x P_X(x) = 1$ Depicted using a bar chart. (Discrete RV)

Common PMFs

① Bernoulli $x \in \{0, 1\}$ $E[X] = p$ ($E[X] = \sum x P_X(x)$)

Bern($x|p$) = $p^x(1-p)^{1-x}$ (This is write in a succinct form)

$$E[f(x)] = \sum_x f(x) P_x(x) \quad \boxed{\text{Var}(X) = E((X - E(X))^2)}$$

$$\text{Var}(X) = E[X^2] - E[X]^2 \quad (\text{Trick}) \quad (y - y^2)$$

(2) Binomial distribution $\sim (y, N)$ [N Bernoulli iid variables]

$$\text{Distribution: } P_X(x) = {}^N C_x p^x (1-p)^{N-x}$$

$x \in \{0, \dots, N\}$

$$E[X] = np \quad \text{Var}(X) = np(1-p) \quad (E22)$$

29.) Entropy of a Random Variable

$H(X)$ - High when variable more random / Less when less random

$$\therefore H(X) \triangleq - E[\log(P_X(x))]$$

$$= - \sum_x P_X(x) \log P_X(x)$$

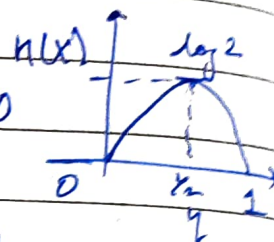
For Bernoulli RV:

$$H(X) = \log 2$$

fair coin

$$H(X) = 0$$

completely biased



Multiple Random Variables: (Joint / Conditional / Marginal)

x - Dice outcome | $y \in \{0, 1\} \Rightarrow \text{Even}$ | $z \in \{0, 1\}$ ^{not} less than 3

$$\therefore P(y, z) =$$

	$z=0$	$z=1$
$y=0$	$1/6$	$2/6$
$y=1$	$1/6$	$2/6$

Joint PDF

Conditional: $P(y, z) / P(z) = P(y|z)$ (Obviously)

Marginal distribution: of single variable $P(z) = \sum_y P(y, z)$

2.10.) Continuous Random Variable and pdfs (Prob. DENSITY fn) [26.18]

$$f_X(x) \Rightarrow P(X=x) = 0 = f_X(x) dx$$

$$\int_{-\infty}^{\infty} f_X(x) dx = 1$$

$$f_X(x) \geq 0 \quad \forall x$$

$f_X(x)$ can be > 1 also

$$F_X(x) = \int_{-\infty}^x f_X(t) dt \quad (\text{Cumulative Density function})$$

Examples: (1) Gaussian $\sim N(x; \mu, \sigma)$ = $\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

$$H(x) = \int_{-\infty}^{\infty} \log(f_X(x)) f_X(x) dx$$

In general.

$$E[f(x)] = \int_{-\infty}^{\infty} f(x) f_X(x) dx$$

② Beta distribution ($f: [0,1] \rightarrow \mathbb{R}^+$)

$$p(x|a,b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}$$

PAGE No.	
DATE	/ /

③ Uniform RV

Normal $x \in [0,1]$

$a=b$ $f(x) = 1/(b-a)$ for $f: [a,b] \rightarrow \mathbb{R}$

Multivariate PDFs

$f(x_1, x_2)$ (Probability density)

or

$$\int_{-\infty}^{\infty} f(\vec{x}) d\vec{x} = 1$$

In general.

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x_1, x_2) dx_1 dx_2 = 1$$

$$p(x_2|x_1=a) = \frac{p(x_1=a, x_2)}{p(x_1=a)} \quad (\text{Conditional density})$$

$$p(x_1) = \int_{x_2} p(x_1, x_2) dx_2 \quad (\text{Marginal density})$$

① Multivariate Gaussian Distribution:

$$\bar{y} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \text{ or } [y_1, \dots, y_n]$$

$$\Sigma = [\sigma_{ij}] \text{ symmetric matrix}$$

$$f(\vec{x}) = \frac{1}{\sqrt{2\pi}^K \sqrt{\det(\Sigma)}} e^{-\frac{1}{2} (\vec{x} - \bar{y})^T \Sigma^{-1} (\vec{x} - \bar{y})}$$

(done.)