Resume:-

1. LSTMs, Bi-directional LSTMs
   - why LSTM > RNNs
   - cell state vs hidden state
   - Difference b/w LSTM & Bi-LSTM

L

2. NLP
   - Transformer architecture
   - self-attention
   - Glove, ELMO, FastText!?

3. *Information retrieval*
   - metrics :- MRR, NDCG
   - Traditional methods

4. Classical ML
   - Mathematics : Logistic, Linear Regression
              SVMs, PCA, SVD,
                      ✓   /   /

   - Linear algebra revision ✓

   -

                        Chi-square test
   ~ Hypothesis testing (linear regression)
     [Confidence interval]     t-test, f-test et.
        p-value          non-parametric testing

# Coding practice

Deep learning

- Optimisation
Local minima, saddle points, local maxima

→ Effect of sample sizes on p-value
and confidence interval.

→ standard error

→ A/B testing case studies   } resources
think about business metrics } around this

→ Using statistics to evaluate choices

Boosting algorithms
    XgBoost.
Bagging   Random Forest

$VU^T = S$
$(USU^T)^{-1} = ...$
$U^T = U^{-1}$
$= \boxed{V S^{-1} U^T}$
$(V^{-T})^{-1} = V$  $UU^T = S$
$U^T = U^{-1}$

## Matrix Rank

→ number of linearly independent rows : row rank
→ number of linearly independent columns : column rank

$\Rightarrow$ row rank = column rank = Rank

$\boxed{Rank(m,n) \leq min(m,n)}$

For a $\boxed{\text{square matrix}}$ if determinant is non-zero, square matrix is a full rank matrix.

## SVD

General setting

$m \times n$   $m \times m$  $m \times n$  $n \times n$

diagonal matrix
$S_i = \sigma_i (i = 1 \sim k, ~) $ else $=0$

$A = USU^T$ ,  $B = USU^T$ is the best k-value approximation to A

Compact SVD $(m \times n)$  $m \times n$  $n \times n$  $r \times r$

$U, V$ : orthogonal matrices $UU^T = I$, $VV^T = S$
$U, V$ : rotate , $S$ : stretching (diagonal matrices)

$\vec{y} = A\vec{x} = (US U^T)(\vec{x})$
$S$ : singular matrix with non-zero diagonal elements
→ if $S$ has 0. element calculate Pseudo-inverse
→ Condition number $= \dfrac{S_1}{S_n}$  largest / smallest

for all the points

If the second derivative of function exists, it is convex if $\partial^2 f(x) \geq 0$ i.e. whether all the eigenvalues of hessian are non negative

→ Local minima is global minima for convex func
→ There can be multiple minima though

Positive definite matrix : Eigenvalues $>$ 0

**\*** If the hessian is everywhere
$\uparrow$ positive semi-definite then function is
convex.

$$M = \begin{bmatrix} \dfrac{\partial^2 f}{\partial x_1^2} & \dfrac{\partial^2 f}{\partial x_1 \partial x_2} & - - & \dfrac{\partial^2 f}{\partial x_1 \partial x_n} \\ & & & \\ \vdots & & & \\ & & & \\ \dfrac{\partial^2 f}{\partial x_n \partial x_1} & \dfrac{\partial^2 f}{\partial x_n \partial x_2} & - - - & \dfrac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

$\checkmark$ All +ve eigen values $\rightarrow$ local minima
All -ve eigenvalues $\rightarrow$ local maxima
mix $\rightarrow$ saddle point

Momentum method
additional $\leftarrow V_t = \beta V_{t-1} + g_{t, t-1}$
rate vector $x_t = x_{t-1} - \eta \cdot V_t$
$\quad\quad\quad\quad\quad\quad\quad - \rightarrow$ learning rate

# Information Theory

Entropy: $H(x) = -\sum_x P(x) \log P(x)$

Mutual information:

$$I(x, y) = H(X) + H(Y) - H(X, Y)$$

<u>How much does X inform Y</u>

KL Divergence $\rightarrow$ model    Amount of information loss

$$D_{KL}(P \| q) = \sum_x P(x) \log \frac{P(x)}{q(x)}$$

when 'q' is used to approximate 'P'.

$\downarrow$ True distribution

\* Maximizing log-likelihood of observing data X with respect to model parameters $\theta$ is equivalent to minimizing KL Divergence between the likelihood and true source distribution of the data.

# Principle of maximum entropy

$\rightarrow$ Max entropy distribution 'agrees with what is known, but expresses minimum uncertainty with respect to all other matters'.

Maximize $-\sum_{P(x)} P(x) \log P(x)$

constraint:- $\sum P(x) - 1 = 0$    $P(x)$ must integrate (sum) = 1

Max entropy distribution when:-

① over a finite discrete range $\{0, 1, \ldots N\}$

  Uniform distribution

② continuous r.v. $X$ with mean by $\mu$

$$\int_X dx \, (x \cdot P(x)) - \mu = 0$$

$\Rightarrow$ Exponential distribution!

③ Maximum entropy distribution with a variance
  $\sigma$

→ variance constraint: $\int dx \, (x-\mu)^2 \cdot P(x) - \sigma^2$

  of
  Also implicit |mean| constraint

$\Rightarrow$ Normal distribution!

KL divergence measured the difference between two probability distributions over the same variable $x$.

# Logistic regression

MLE estimate:
Bernouli
$$P(Y=1 \mid x) = h_\theta(x) = P$$
$$P(Y=0 \mid x) = 1 - h_\theta(x) = 1-P$$
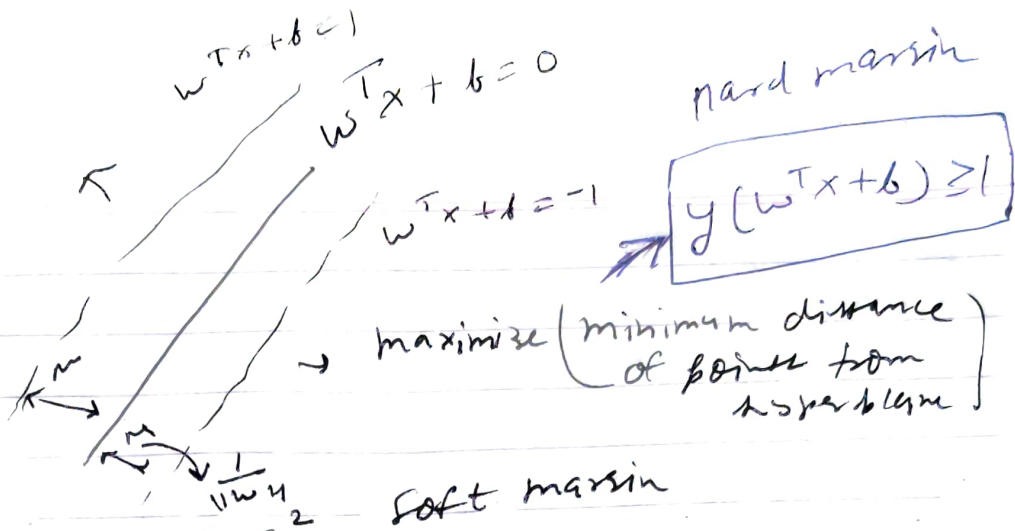
$$Likelihood = (P)^y (1-P)^{1-y}$$

$$\frac{(W-k+2P)}{S} + 1 : output \; shape$$

$$Pr\left(|x-\mu| \geq k\sigma\right) \leq \frac{1}{k^2}$$

or equivalently

$$Pr\left(|x-\mu| \geq k\right) \leq \frac{\sigma^2}{k^2}$$

SVM

$w^T x + b = 1$

$w^T x + b = 0$

$w^T x + b = -1$

hard margin

$$y(w^T x + b) \geq 1$$

maximize $\left(\begin{array}{c}\text{minimum distance}\\ \text{of points from}\\ \text{hyperplane}\end{array}\right)$
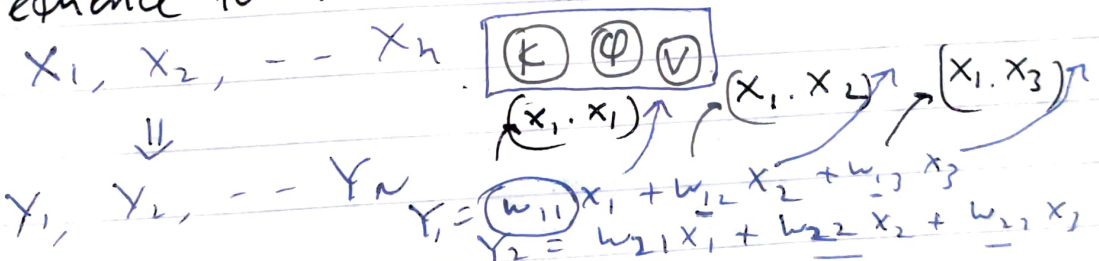
$\frac{1}{\|w\|_2}$   soft margin

min. $\frac{1}{2}\|w\|^2 + C \xi$

Such that   $y_i(w_i x_i + b_i) \geq 1 - \xi_i$

Self - attention

sequence to sequence operation

$X_1, X_2, --- X_n$   $\boxed{K}$ $\boxed{\varphi}$ $\boxed{V}$   $(X_1 . X_2)$   $(X_1 . X_3)$

$(X_1 . X_1)$

$\Downarrow$

$Y_1, Y_2, --- Y_n$

$Y_1 = w_{11} X_1 + w_{12} X_2 + w_{13} X_3$

$Y_2 = w_{21} X_1 + w_{22} X_2 + w_{23} X_3$

Basic operation: $Y_i = \sum_{j} (W_{ij}) X_j$   weighted average of inputs

sum over $j$

$\rightarrow$ softmax operation to normalize

$W_{ij} = (X_i . X_j)$   $(X \times X^T) X$

key, query, value

$K = w_K^T X$, $\varphi = w_\varphi^T X$, $V = w_v^T X$

$\boxed{\dfrac{Softmax(K \cdot \varphi)(V)}{\sqrt{dim}}}$   self - attention

## SGD or GD

$$\theta = \theta - \alpha \nabla L_\theta$$

## Add Momentum

$$v = \beta v_{t-1} + \nabla L_\theta$$
$$\theta = \theta - \alpha v$$

## Adagrad

$$r = r + \nabla L_\theta \odot \nabla L_\theta \qquad g : \text{gradient}$$
$$\theta = \theta - \frac{\alpha}{\epsilon + \sqrt{r}} \nabla L_\theta$$

## RMS prop with momentum

decay factor

$$r = (1-\rho) \nabla L_\theta \odot \nabla L_\theta + \rho (r) \quad \text{exponential average}$$
$$v = \beta v + \alpha \frac{1}{\rho} \odot g$$
$$\theta = \theta - \underline{\alpha} v \qquad \text{decay factor}$$

## Adam

$$s = \rho_1 s + (1 - \rho_1) g$$
$$r = \rho_2 r + (1 - \rho_2) g \odot g$$

$$s = \frac{s}{1 - \rho_1}$$

$$r = \frac{r}{1 - \rho_2}$$

$$\theta = \theta - \alpha \frac{s}{r}$$

# SVD relation to Eigen-decomposition

SVD:
$$A = U \Sigma V^T$$

Eigen decomposition

$$A = X \Lambda X^T \qquad *$$

> A needs to be symmetric
> U, V, X are Orthonormal
> $\Lambda, \Sigma$ are diagonal

$$AA^T = (U\Sigma V^T)(U\Sigma V^T)^T = (U\Sigma V^T)(V \Sigma^T U^T)$$
$$= U\Sigma \Sigma^T U^T$$
$$\qquad\qquad \downarrow \quad \downarrow \quad \downarrow$$
$$\qquad\qquad X \quad \Lambda \quad X^T$$

$$A^T A = V\Sigma \Sigma^T V^T$$
$$\qquad\quad \uparrow \quad \uparrow \quad \uparrow$$
$$\qquad\quad X \quad \Lambda \quad X^T$$

Shows how to SVD using Eigenvalue decomposition

$$\lambda_i = \sigma_i^2$$

SVD:-

1. Optimal low-rank approximation
2. Interpretability problem
3. Lack of sparsity

# LSTM

3 gates: Input, forget and Output

$$I = \sigma(\qquad)$$
$$f = \sigma(\qquad)$$
$$O = \sigma(\qquad)$$

Additional memory cell $C_t$:   using $H_{t-1}$

candidate   $\bar{C}_t = \tanh(X_t, H_{t-1}, W)$ ✓

$C_t =$ $\underbrace{\text{memory cell}}$ → forget from the to review cell store

$$F_t \odot C_{t-1} + I_t \odot C_t$$
↳ read from the new
cell state

$$\boxed{H_t} = O_t \odot \tanh(C_t)$$

# Gated Recurrent Units:-
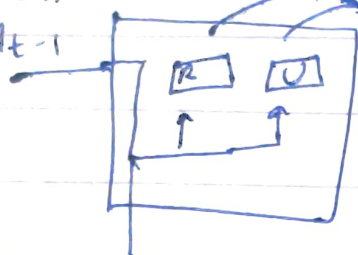
RNNs:

$$H_t = f(W_h, X_t, H_{t-1}) \quad \checkmark$$
$$O_t = g(H_t, W_o) \quad \checkmark$$

GRUs introduce two designs (gates):
1. Reset gate  2. Update gate

hidden state
$H_{t-1}$



Sigmoid gates

$$\checkmark \; R_t = \sigma(X_t W_{xr} + H_{t-1} W_{hr} + b_r)$$
$$\checkmark \; Z_t = \sigma(X_t W_{xz} + H_{t-1} W_{hz} + b_z)$$

Input $X_t$

Reset gate controls how much previous state we might want to remember

Update gate controls how much of the new state is just a copy of the old state

Candidate hidden states:-

① $\tilde{H}_t = \tanh(X_t W_{xh} + (R_t \odot H_{t-1}) W_{hh} + b_h)$

element wise multiplication

② $H_t = Z_t \odot H_{t-1} + (1 - Z_t) \odot \tilde{H}_t$

Helps with vanishing gradient problem in RNNs and better captures the dependencies for sequences with large time step distances