

DEEP LEARNING



WII DS

Structured Predictn \leftrightarrow Hidden Markov Model (Supervised Learning)
 \rightarrow Generative Model.

Naive Bayes \xleftrightarrow{vs} Logistic Regression

Discriminative Model for Structured Predictn

Nxt wh - Unsupervised Regressn (PCA)

Advanced Topics

$$P(y|x) = \frac{1}{1 + \exp(-y(\vec{\theta} \cdot \vec{x} + \theta_0))}$$

\downarrow
+1 or -1 tagging
 \downarrow

Now what if multiple classes

Logistic Regression

$$P(y|x) = \frac{\exp(\sum_{k=1}^K \lambda_k \cdot f_k(x, y))}{\sum_y \exp(\sum_{k=1}^K \lambda_k \cdot f_k(x, y))}$$

Softmax Regression

Entropy - Degree of uncertainty in data sent

In war, sent 1 bit \rightarrow event occur in war

2 bit \rightarrow 4 possible event occurred

A	0	0
B	0	1
C	1	0
D	1	1

3 bit \rightarrow 8 events

$\log_2 k$ bit $\rightarrow 2^k$ possible event

\rightarrow Expected no of bit : $\sum_x P(x) \cdot \log_2 \frac{1}{P(x)}$

\hookrightarrow Degree of uncertainty high if $P_1, P_2, P_3, P_4 = \frac{1}{4}$

\hookrightarrow Entropy

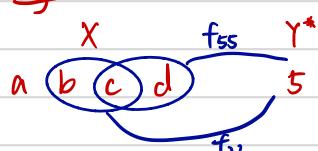
Constraint Optimization Problem

Given (X, Y^*)

$$\max_{\min_{P(Y|X)}} \sum_y P(Y|X) \cdot \log P(Y|X)$$

constraint by ① $\sum_y P(Y|X) \cdot f_k(x, y) = f_k(x, Y^*) \rightarrow \sum_y P(Y|X) \cdot f_k(x, y) - f_k(x, Y^*) = 0$
② $\sum_y P(Y|X) = 1 \rightarrow \sum_y P(Y|X) - 1 = 0$

Eg



Add Lagrangian multiplier

bad guys $\downarrow (-\infty, \infty)$

$$\sum_y P(Y|X) \cdot \log P(Y|X) - \sum_{k=1}^K \lambda_k \left[\sum_y P(Y|X) \cdot f_k(x, y) - f_k(x, Y^*) \right] - \lambda_0 \left[\sum_y P(Y|X) - 1 \right]$$

\rightarrow can be +ve/-ve

$$\sum_y P(Y|X) \cdot \log P(Y|X) - \sum_{k=1}^K \lambda_k \left[\sum_y P(Y|X) \cdot f_k(x, y) - f_k(x, y^*) \right] - \lambda_0 \left[\sum_y P(Y|X) - 1 \right]$$

bad guys $(-\infty, \infty)$ can be +ve/-ve

$$\frac{\partial L}{\partial P(Y|X)} = \log P(Y|X) + 1 - \sum_{k=1}^K \lambda_k [f_k(x, y)] - \lambda_0 = 0 \quad - (1)$$

$$\log P(Y|X) = -1 + \sum_{k=1}^K \lambda_k [f_k(x, y)] + \lambda_0$$

$$P(Y|X) = \exp \left(-1 + \sum_{k=1}^K \lambda_k [f_k(x, y)] + \lambda_0 \right)$$

$$= \exp \left(\sum_{k=1}^K \lambda_k f_k(x, y) \right) \cdot \exp(\lambda_0 - 1)$$

$$= \frac{\exp \left(\sum_{k=1}^K \lambda_k f_k(x, y) \right)}{z}, \quad z = \frac{1}{\exp(\lambda_0 - 1)}$$

$$= \frac{\exp \left(\sum_{k=1}^K \lambda_k f_k(x, y) \right)}{\sum_y \exp \left(\sum_{k=1}^K \lambda_k f_k(x, y') \right)}, \quad \begin{aligned} P(Y_1|X) &= \frac{1}{2} \\ P(Y_2|X) &= \frac{3}{4} \\ P(Y_3|X) &= \frac{5}{8} \end{aligned}$$

$$\max_{\lambda} \sum_y P(Y|X) \cdot \log P(Y|X) - \sum_{k=1}^K \lambda_k \left[\sum_y P(Y|X) \cdot f_k(x, y) - f_k(x, y^*) \right] - \lambda_0 \left[\sum_y P(Y|X) - 1 \right]$$

bad guys $(-\infty, \infty)$ can be +ve/-ve

$$= \max_{\lambda} \sum_y P(Y|X) \left[\log P(Y|X) - \sum_{k=1}^K \lambda_k - \lambda_0 \right] + \sum_{k=1}^K \lambda_k \cdot f_k(x, y^*) + \lambda_0$$

-1 from eqn(1)

$$= \max_{\lambda} - \sum_y P(y|X) + \sum_{k=1}^K \lambda_k f_k(x, y^*) + \lambda_0$$

$$= \max_{\lambda} \log P(Y^*|X) \rightarrow \text{Learning Problem.}$$

Add a \log in front.

For (x, y^*)

$$L = \log P(Y^*|X) = \sum_{k=1}^K \lambda_k f_k(x, y^*) - \log \sum_y \exp \left(\sum_{k=1}^K \lambda_k f_k(x, y) \right)$$

$$\frac{\partial L}{\partial \lambda_j} = f_j(x, y^*) - \frac{1}{z} \sum_y \left[\exp \left(\sum_{k=1}^K \lambda_k f_k(x, y) \right) \cdot f_j(x, y) \right]$$

$$= f_j(x, y^*) - \sum_y \frac{\exp \left(\sum_{k=1}^K \lambda_k f_k(x, y) \right)}{\sum_y \exp \left(\sum_{k=1}^K \lambda_k f_k(x, y') \right)} \cdot f_j(x, y)$$

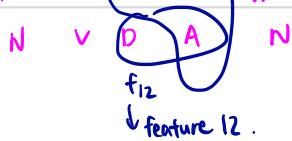
$$= f_j(x, y^*) - \sum_y P(y|X) \cdot f_j(x, y)$$

\downarrow
Actual time you observe

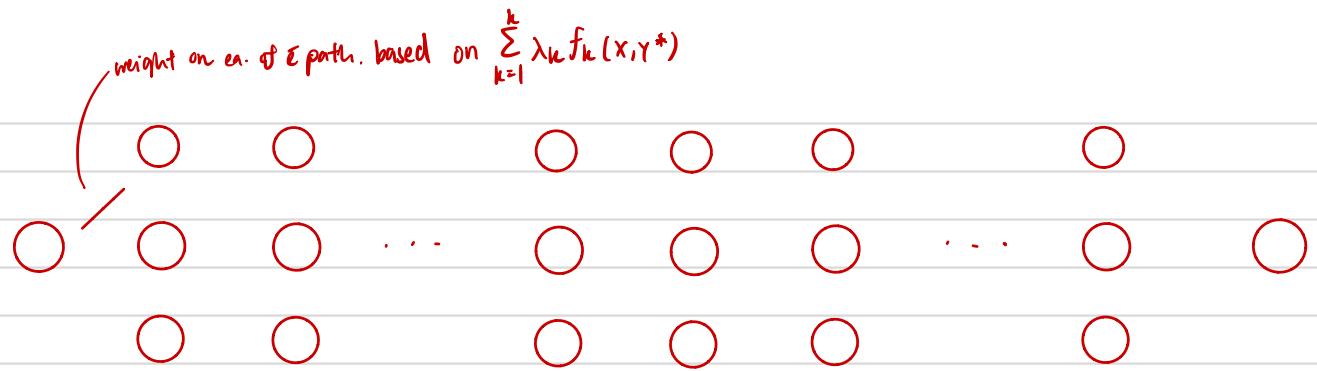
\downarrow
Expected time you observe y . \rightarrow calculated using forward-backward algo

\hookrightarrow due to lack of closed form form

Eg. Earth is \boxed{A} fine invents f_{23}



\downarrow feature 12.



$$\text{For } \{-\log \sum \exp [\sum_{k=1}^k \lambda_k f_k(x, y)]\}$$

$$\text{Compare to } \log \sum_{a \in S} \exp(a) = 21.000\dots$$

$$\max_{a \in S}(a) = 21 \quad \text{where } S = \{2, 5, 21\}$$

Then $L = \log P(Y^* | X) = \sum_{k=1}^k \lambda_k f_k(x, y^*) - \max(\sum_{k=1}^k \lambda_k f_k(x, y))$

+ Δ if 0 then training set is linearly separable (Perception Loss) and $y=+1 \& -1$

score of first path is ϵ + Δ if not 0 and obtain hinge loss

score of best path = Structured SVM.

Introduction of Latent Variable.

Faith - Noun - Proper Noun? or Improper noun. (ie latent variable) learn from data-set.

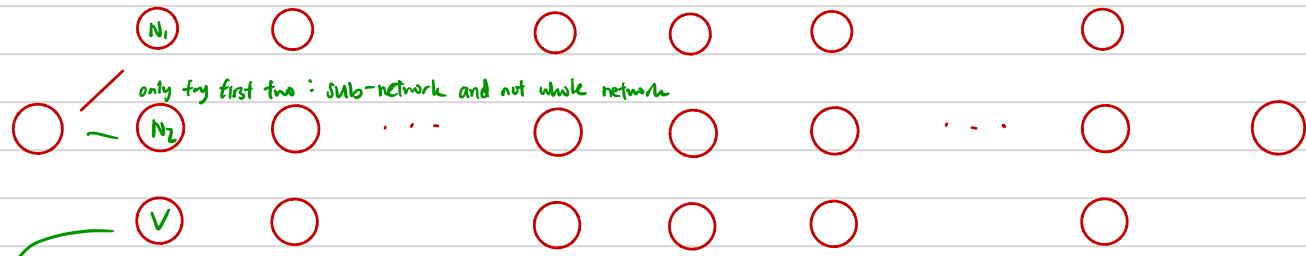
$$\sum_h P(Y^*, h | X) = \sum_h \frac{\exp(\sum_{k=1}^k \lambda_k f_k(x, h, Y^*))}{\sum_{h'} \exp(\sum_{k=1}^k \lambda_k f_k(x, h', Y'))}$$

$$L = \log \sum_h \exp(\sum_{k=1}^k \lambda_k f_k(x, h, Y^*)) - \log \sum_{h'} \exp(\sum_{k=1}^k \lambda_k f_k(x, h', Y'))$$

$$\frac{\partial L}{\partial \lambda_j} = \sum_h P(h | X, Y^*) \cdot f_j(X, Y^*, h) - \sum_{h' \neq h} P(h' | X) \cdot f_j(X, Y^*, h')$$

↓ Given X, Y^*
↓ Obtain Feature Count.

↓ Given X and n th else
↓ Obtain Feature Count

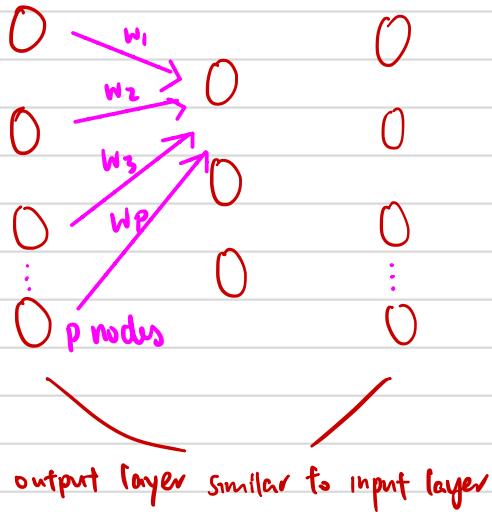


λ : tuned for optimal path in sub-network. and max eqn
Tendency to go verb is low.

(X, Y) can be diff from (X, Y^*)

W2D5

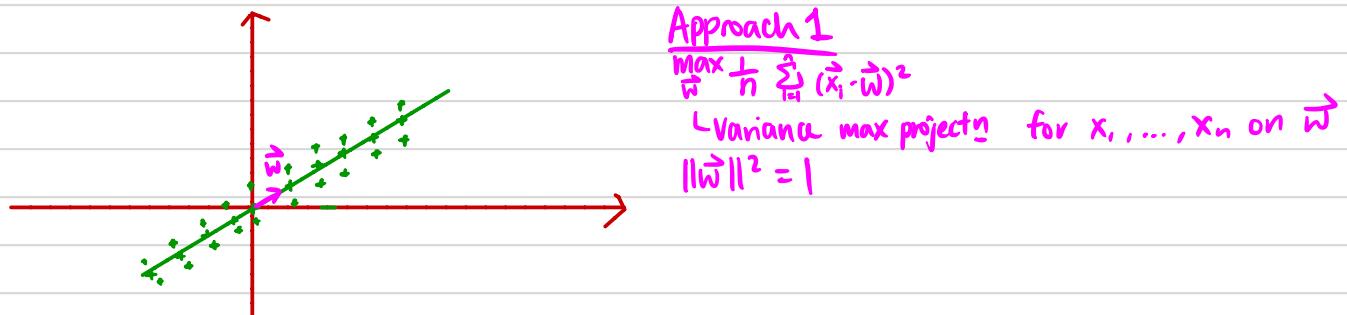
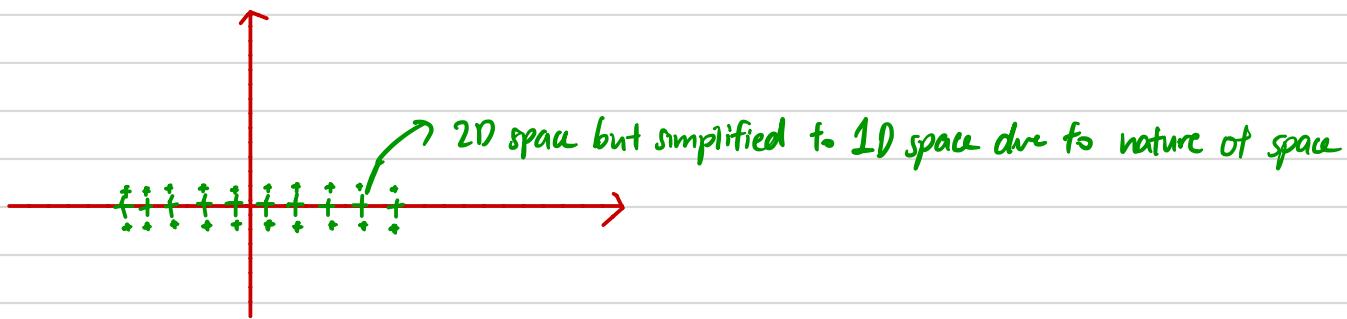
Autoencoder



Unsupervised Regression
(Left out from Normal Class)

		movies			
		2	5		
		1		4	
users		2	1		
		1	2	5	2
		5	2	5	1
		3		1	3

Rating for movie



Approach 2

$$\min_{\vec{w}} \frac{1}{n} \sum_{i=1}^n \| \vec{x}_i - (\vec{x}_i \cdot \vec{w}) \cdot \vec{w} \|^2$$

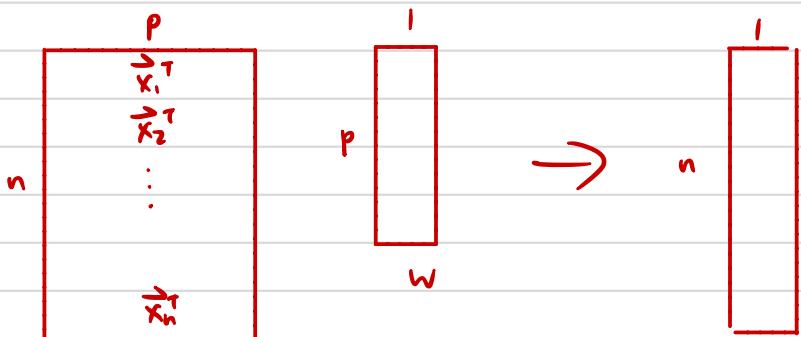
$$= \frac{1}{n} \sum_{i=1}^n \| \vec{x}_i \|^2 - 2 \sum_{i=1}^n (\vec{x}_i \cdot \vec{w}) \vec{w} + \sum_{i=1}^n \| (\vec{x}_i \cdot \vec{w}) \vec{w} \|^2$$

Approach 2

$$\begin{aligned}
 & \min_{\vec{w}} \frac{1}{n} \sum_{i=1}^n \| \vec{x}_i - (\vec{x}_i \cdot \vec{w}) \vec{p} \|_2^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \| \vec{x}_i \|^2 - 2 \vec{x}_i (\vec{x}_i \cdot \vec{w}) \vec{p} + \| (\vec{x}_i \cdot \vec{w}) \vec{p} \|_2^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \| \vec{x}_i \|^2 - 2 (\vec{x}_i \cdot \vec{w})^2 + (\vec{x}_i \cdot \vec{w})^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \| \vec{x}_i \|^2 - (\vec{x}_i \cdot \vec{w})^2 + (\vec{x}_i \cdot \vec{w})^2
 \end{aligned}$$

Approach 1

$$\begin{aligned}
 & \max_{\vec{w}} \frac{1}{n} \sum_{i=1}^n (\vec{x}_i \cdot \vec{w})^2 \\
 &= \frac{1}{n} (\vec{X}_W)^T \vec{X}_W \\
 &= \frac{1}{n} \vec{w}^T \vec{X}^T \vec{X}_W \\
 &= \vec{w}^T \left(\frac{\vec{X}^T \vec{X}}{n} \right) \vec{w} \quad \text{true semi definite} \\
 &= \vec{w}^T A \vec{w} = \vec{w}^T \lambda \vec{w} = \lambda \vec{w}^T \vec{w} = \lambda
 \end{aligned}$$



→ Lagrangian Multiplier

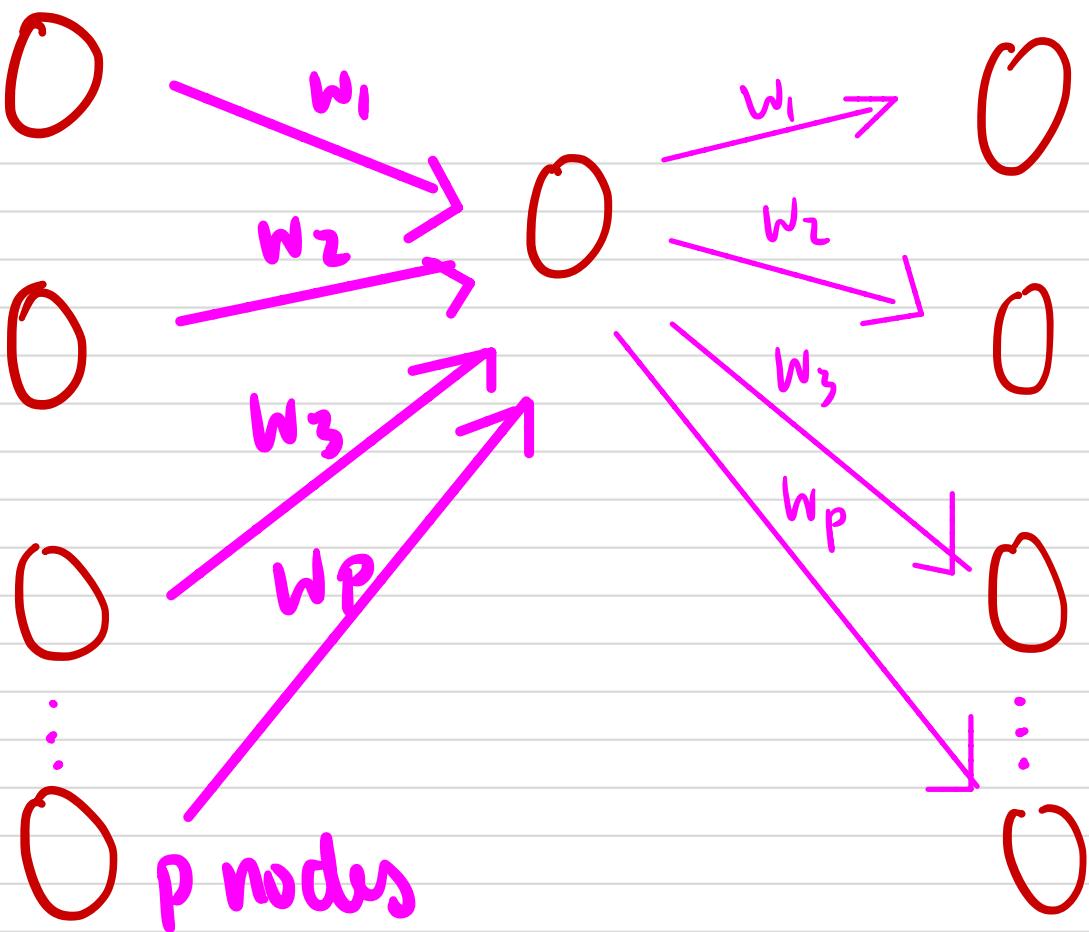
$$\begin{aligned}
 \lambda &= \vec{w}^T A \vec{w} - \lambda (\vec{w}^T \vec{w} - 1) \quad , \quad \vec{w}^T \vec{w} = 1 \\
 \frac{\partial \lambda}{\partial \vec{w}} &= 2 A \vec{w} - 2 \lambda \vec{w} = 0 \quad \rightarrow \quad A \vec{w} = \lambda \vec{w}
 \end{aligned}$$

↙ Eigenvalue (no imaginary value)
↙ Eigenvector

Soln : Use SVD (Single Value Decomposition) ↴ $\vec{v}^T \vec{v} = \vec{v} \vec{v}^T = 1$

$$X = U \Sigma V^T$$

$$\begin{aligned}
 A &= \frac{1}{n} V \Sigma U^T U \Sigma V^T \\
 &= V \frac{\Sigma^2}{n} V^T
 \end{aligned}$$

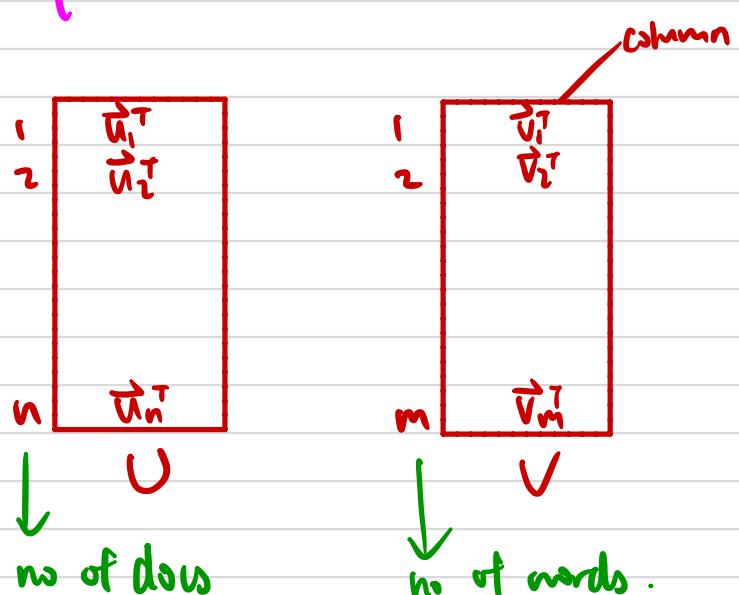


Principal Component Analysis.

	movies				m
users	1	2	5	4	
n	1	2	1	5	5
1	2	5	2	5	2
5	2	5		1	1
3		1		3	3

Collaborative Filtering.

$$\text{Aim: } \vec{u}_3, \vec{v}_2 \approx 2$$



$$\arg \min_{U, V} \sum_{(i, j) \in D} (\vec{u}_i \cdot \vec{v}_j - r_{ij})^2 + \lambda \sum_i (\vec{u}_i)^2 + \lambda \sum_j (\vec{v}_j)^2$$

↓
 Find similar scores & input into empty movie rating scores.
 ↗ Upgrades Linear → Ridge Regression .

Regularized Term

values fit don't
veer too off

Latent Semantic Analysis