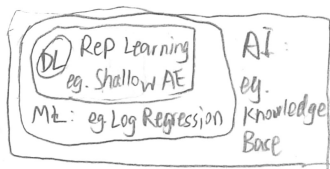


§ Historial

1958 Rosenblatt's Perceptron Algo (X XOR Problem)
 Backpropagation Algo (1986~1979) (X very deep)
 mid-90's SVM (outperforms NN)
 2006 Restricted Boltzman Machine to pretrain DNN



§ ML Basis

Scalar: $x \in \mathbb{R}$
 Vector: $\vec{x} \in \mathbb{R}^p$
 Matrices: $A \in \mathbb{R}^{m \times n}$
 Tensor: $A \in \mathbb{R}^{m \times n \times r}$

Vector Norm: $\|\vec{x}\|_p = (\sum |x_i|^p)^{1/p}$ ($p \geq 1$)

Frobenius Norm: $\|A\|_F = \sqrt{\sum_{i,j} A_{i,j}^2}$
 Dot product of 2 vect: $\vec{x}^T \vec{w}$
 Mat Mult: $\vec{z} = \vec{w} \vec{x} = (\vec{x}^T \vec{w}^T)^T$



Vector space (set of vectors)
 (linear) subspace: V that is a subset of some larger V

$$\text{Var}(f(x)) = E[(f(x) - E[f(x)])^2]$$

$$\text{Var}(x) = E[(x - E[x])^2] = E[x^2] - (E[x])^2$$

$$\text{Cov}(f(x), g(y)) = E[(f(x) - E[f(x)])(g(y) - E[g(y)])]$$

$$\text{Cov}(x, y) = E[(x - E[x])(y - E[y])] = \sigma_{xy}$$

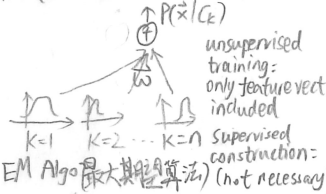
Structure of PR sys.



Bayes $P(Y|X) = P(X|Y)P(Y) \rightarrow$ prior: prob without any given info
 $P(X) \rightarrow$ Marginal Prob
 posterior

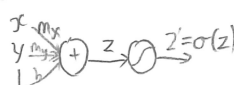
for K-class classification
 $k = \arg \max_{i=1}^K P(\vec{x}|w_i)$ (if $P(\vec{x}|w_k) > P(\vec{x}|w_i)$)
 or $k = \arg \max_{i=1}^K P(\vec{x}|w_i)P(w_i)$

GMM (Gaussian Mixture Model)



Model capacity
 Underfit \rightarrow dim reduction
 Overfit \rightarrow val project high-dim to low-dim
 Optimal Capacity
 SVM: maximize margin
 $f(x) = \sum_{i \in S} \alpha_i y_i \vec{x}_i^T \vec{x} + b$
 goal for binary classification

§ DNN Basic



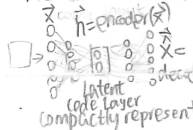
developed in 2010
 alleviate gradient vanishing
 allow hidden layer to output true 0 (speed up training, sparse rep)

DNN (MLP)
 Classif $\phi(\sigma)$
 Reg $\phi(\text{lin})$

Activation Func
 σ (sigmoid) $= \frac{1}{1 + e^{-z}}$
 $\text{ReLU} = \max(z, 0)$
 $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

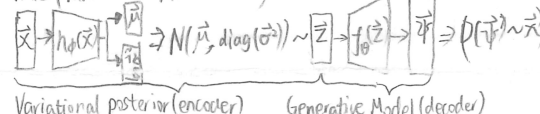
§ Deep Arch

AE (Autoencoders)

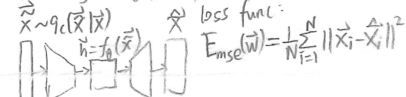


using $E_{mse}(\vec{w}) = \frac{1}{N} \sum_{i=1}^N \|\vec{x}_i - \hat{\vec{x}}_i\|^2$

VAE (Variational Autoencoder)

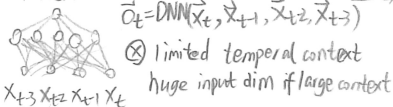


DAE (Denoising Autoencoder)



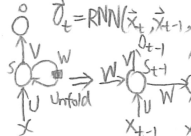
Sequential Data

Forced FNN



more flexible smarter

RNN



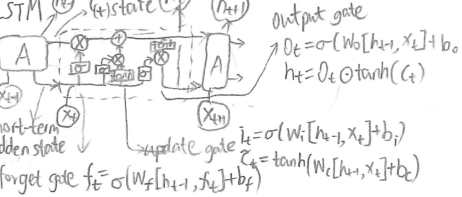
$$\vec{z}_t = f(\vec{u}_t \vec{x}_t + \vec{w}_t \vec{x}_{t-1})$$

$$\vec{a}_t = \text{softmax}(\vec{q}_t) \quad (\vec{q}_t = \vec{V} \vec{z}_t)$$

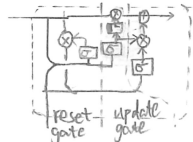
(Gen image description, mach tran classify data sequence)

Can only look back a finite number of steps

long-term dependency



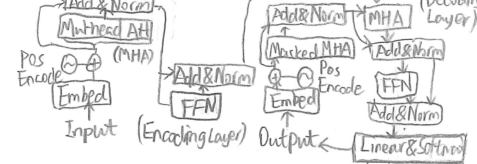
GRU (Gate Recurrent Unit)



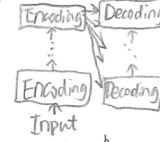
Encoder-Decoder Architecture



Transformer



(Decoding Layer)



§ DL Math Fundamentals

$$\vec{w} = \vec{w} - \eta \frac{\partial L(\vec{w})}{\partial \vec{w}}$$

Minibatch SGD

$$\vec{w}_b = \vec{w} - \eta \frac{\partial L(\vec{x}_b)}{\partial \vec{w}}$$

Momentum (reduce ziggle, acc training)
 $\Delta \vec{w} = \eta \frac{\partial L(\vec{w})}{\partial \vec{w}} + \gamma \Delta \vec{w}_{t-1}$ (especially noisy grad)

Loss Functions

$$L_{mse} = \frac{1}{N} \sum_{i=1}^N \|\vec{y}_i - \vec{t}_i\|^2$$

$$L_{cce} = -\frac{1}{N} \sum_{k=1}^K \sum_{i=1}^N t_{i,k} \log y_{i,k}$$

$$L_{ce} = (\text{special case of CCE})$$

DNN Problem

1. Gradient Vanishing
 Sol. ReLU Residual, Proper \vec{w} init (eg. Xavier)
 Batch norm, Layer norm, Scaled Exp LU
 Xavier: $\vec{w}_{ij}^{(l)} \sim N(0, \frac{2}{n_{l-1} + n_l})$
2. Gradient Exploding
 Sol. Clip the gradient
3. (Almost) No data per class (few shot learning)
 Sol. Transfer Learning, Meta Learning, Self-supervised

Regularization (Solution Space)

L2 Regularization

$$L(\vec{w}) = L_{cce}(\vec{w}) + \lambda \vec{w}^T \vec{w}$$

L1 Regularization

$$L(\vec{w}) = L_{cce}(\vec{w}) + \lambda \|\vec{w}\|_1$$

