

CE&A.

1. No access to $P(X, Y)$? \Rightarrow Hold back \rightarrow test set.
2. Time not differentiable? \Rightarrow Use a surrogate train
3. overfit? \rightarrow regularization / Reduce model order.
4. Expressivity \rightarrow piecewise linear (ReLU).
learnability (e.g. Ridge reg.).

Ridge Regression:

1. OLS problem $\arg\min \|y - Xw\|^2 + \lambda \|w\|^2$.
2. Shift Singular val. $(XX^T + \lambda I)^{-1} X^T y$
3. Maximum A Posteriori (MAP)
 $y = Xw + \sqrt{\lambda} N, N \sim N(0, I)$.
4. Fake data. $\hat{X} = \begin{bmatrix} X \\ \sqrt{\lambda} I_d \end{bmatrix} \quad \hat{y} = \begin{bmatrix} y \\ 0_d \end{bmatrix}$
5. Fake tests. $\tilde{X} = [X \quad \sqrt{\lambda} I_n]$.

↳ restrictions: least-squares linear reg.

1. $\eta < \frac{1}{\sigma^2}$.
2. Takes $\log(1/\epsilon) \rightarrow \eta \sigma^2 \epsilon$ iters
3. For fastest overall: $\eta = \frac{1}{\sigma_1^2 + \sigma_s^2}$

MWD: adjust the magnitude of the parameter update in each dim \rightarrow optimizer takes bigger & more confid. updates. less oscillation.

Initialization (keep input $\sim N(0,1)$)

1. weights
 - 1) Xavier init: $\sigma^2 = \frac{1}{d}$ (d: fan in, non-zero weights).
 - 2) He init: $\sigma^2 = \frac{2}{d}$ (He reLU "zero" hat).
2. bias
 - 1) treated as d+1 weight.
 - 2) $b=0$
 - 3) small random num ($b=0.01$)

LD.

1. vanilla gradient (LD)
2. SGD with momentum.
cope with shake. (other methods: low pass filter & Avg)
3. Adaptive method: Adam.

{ diff wrt to diff dim
or decay.

$$\begin{aligned} \tilde{a}_{k+1} &= (1-\beta)\tilde{a}_k + \beta \nabla_w L(w^k) \\ \tilde{V}_{k+1} &= (1-\beta^2)\tilde{V}_k + \beta^2 \nabla_w \left(\left(\frac{\partial L(w)}{\partial w_i} \right)^2 \right) \\ w_{k+1}[i] &= w_k[i] - \eta \frac{\tilde{a}_{k+1}[i]}{\sqrt{\tilde{V}_{k+1}[i] + \epsilon}} \end{aligned}$$

track the size of gradient
for stable.

CNN.

1. key ideas
 - 1) locality (neighbors matter)
 - 2) invariance (translation \rightarrow)
 - 3) support hierarchical structure & multi-resolution understanding

$$2. \text{ output } w \left\lfloor \frac{W+2P-K}{S} + 1 \right\rfloor.$$

weights: $K = K_c \times F$.

3. pooling: no weights compared with stride.
 - Max-pool: invariant feats v. retain strong feats highlight the presence of an object.
 - Avg-pool: capture smooth transitions / subtle change

Dropout:

1. improve generalization of network, overfit \downarrow
2. dropout rate: cause a discontinuity when 0+, cuz not full rank when $dr=0$. mult. solu.
3. Dropout & normalization.
 $L(w) = \|y - Pxw\|^2 + p(1-p)\|w\|^2$.
4. weights muted by $\frac{1}{1-p}$.

UNN.

1. CNN vs UNN (neighbors distinguishable?)
 - 1) weight-sharing: UNN neighbors
 - 2) pooling: UNN cluster and merge.
 - 3) Normalization: { layer (over channel)
batch (sample nodes as a batch)}
2. ① A: adjacency mat. D: degree mat.
A Normalized = AD^{-1}
 $A_{SymNorm} = D^{-\frac{1}{2}} A D^{\frac{1}{2}}$
② $L = D - A$ (Laplacian mat)
 $L_{ij} = \begin{cases} \deg(u_i), & i=j \\ -1, & i \neq j \text{ \& v.i. adj. to } u_j \\ 0, & \text{else.} \end{cases}$
 $L_{SymNorm} = I - A_{SymNorm}$
 $= D^{-\frac{1}{2}} L D^{-\frac{1}{2}}$
3. Be careful:
 - { which pts are for testing when calc. loss.
 - 100% test all with misclassified samples?

RNN. (FIR)

1. CNN: Finite Impulse Response Filters.
what if infinite? How to represent? \Rightarrow hidden stat.

$$RNN: IIR. \quad y[t] = a * y[t-1] + b * x[t].$$

2. E.g. Kalman Filter

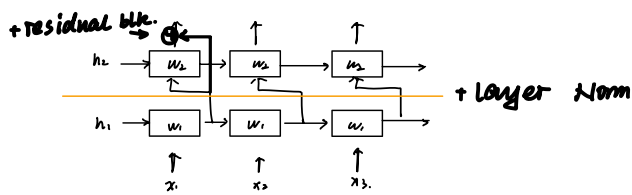
$$\begin{cases} \tilde{x}_t = \tilde{c}\tilde{x}_{t-1} + \tilde{v}_t \text{ (know } \tilde{v}_t, c) \\ \tilde{h}_{t+1} = A\tilde{h}_t + Bx_t \text{ (to learn } A, B) \end{cases}$$

case 1: data $\{c, \tilde{h}_t, \tilde{x}_t\} \Rightarrow$ comp \tilde{h}_t and \tilde{h}_0 .

case 2: data $\{\tilde{x}_0\} \Rightarrow$ comp \tilde{x}_t and \tilde{x}_t .

3. Challenges in RNN: gradient explode / vanish.

- 1) Saturating activation fns (e.g. sigmoid, tanh)
- 2) Gradient clipping.
- 3) Normalizations



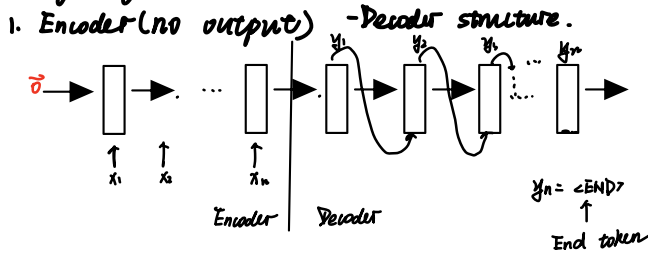
→ Interpreted as self-supervision as AE.

$X_{d=n} = U \Sigma V^T$
 ↳ Project into $\vec{u}_1 \dots \vec{u}_k$ (reduced subspace)

$$\min_{W, U} \|y - U W W^T x\|^2$$

*: $U: d \times k, W: k \times d$, ensure $\dim \leq k$.

Seq2Seq model



Math.

$$1. \quad \nabla_x f(x) = \begin{bmatrix} \frac{df}{dx_1} \\ \vdots \\ \frac{df}{dx_n} \end{bmatrix} \quad \nabla_x^T f(x) = \begin{bmatrix} \frac{df}{dx_1} & \dots & \frac{df}{dx_n} \end{bmatrix}$$

$$2. \quad \frac{d}{dx} x^T A x = (A + A^T) x$$

$$3. \quad \frac{d}{dx} w^T x = w$$

$$4. \quad \frac{d}{dA} (x^T A y) = x y^T$$

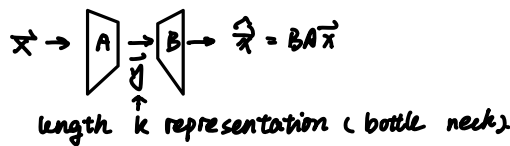
$$\frac{d}{dA} (x^T A^T y) = y x^T$$

Note we're writing derivative in column! (Transpose if in row)

① Note: during training, may feed in $h_T y$.

② If $\text{len}(\text{output}_{\text{true}}) \neq \text{len}(\text{output}_{\text{AT}})$: stop early.

2. Autoencoder (expect output → input).



↳ Vanilla AE.

→ Denoise AE. $\begin{cases} \text{inputs: } \vec{x}_i + \vec{n} \\ \text{expected outputs: } \vec{x}_i \end{cases}$

Do not add noise to test set.

→ Masked AE. $\begin{cases} \text{inputs: } \vec{x}_i + \text{mask} \\ \dots : \vec{x}_i \end{cases}$

① After pretrain, AE's decoder is useless

② Pre-trained sometimes worse perf.

③ Autoencoder representation more useful than raw inputs when few labels for downstream.

④ For language problem, loss: cross-entropy.

⑤ With imgs, more effective to mask patch instead of pixels. (Models'll get decent loss by copying neighbors pixels.)

⑥ Can think of denoising & masked AE as vanilla with data augmentation.

Attention

1. queryable softmax pooling: differentiable.

$$\begin{cases} \text{Compute "distance": } e_i = \frac{\langle \vec{h}_i, \vec{q} \rangle}{\sqrt{d}} \\ \text{Compute "weight": } \alpha_i = \frac{e_i}{\sum e_j} \\ \text{Compute "output": } \sum \alpha_i v_i \end{cases}$$

2. $\begin{cases} \text{self attention: key-val pairs from decoder} \\ \text{cross attention: .. encoder} \end{cases}$

Self-supervision (unsupervised learning)

1. PCA: de-mean data

↳ ① data mat → SVD. $\frac{d}{X} = U \Sigma V^T$

② Compute $v_i^T \vec{x}$ to get i th principle feature of a new pt \vec{x} .