LEC6  DSA5204

Re-cap: $\begin{cases} FCNN \\ RNN \\ CNN \end{cases}$ $\begin{cases} GD \\ SGD \text{ (computationally efficient)} \\ SGD \text{ with momentum} \end{cases}$

same epoch, more updates!

gradient → "external force"

$\downarrow$                              $\downarrow$

Model                          Optimization Algorithm

Future:   ① How to improve **performance** ?

② Outside super-vised learning scope.

---

① How to improve model performance ?

$\begin{cases} a) \text{ Model Architecture} \\ b) \text{ Training Method ( Optimization )} \\ c) \text{ Data ( Augmentation )} \end{cases}$

1. **Regularization** → low testing error

→ probably high training error

$\Downarrow$

**Inductive bias** such that model can behave better with limited data

Previously: **ERM Framework** ( Empirical Risk Minimization )

$$\hat{\theta}_{ERM} = \arg\min_{\theta} R_{emp}(\theta)$$

Regularization Framework

$$\hat{\theta}_{Reg} = \arg\min_{\theta} R_{emp}(\theta) + \alpha \underline{\Omega(\theta)}$$

$\downarrow$

regularization / penalty term

**a)** **$\ell_2$-reg** : $\Omega(\theta) = \frac{1}{2}\|\theta\|_2^2$

$\hookrightarrow$ Gaussian prior on $w$ (statistics perspective)

Toy example: (Ridge (linear) regression)

$$\begin{cases} \Omega(w) = \frac{1}{2}\|w\|_2^2 \\ R_{emp}(w) = \frac{1}{2}\|Xw - y\|_2^2 \end{cases} \implies$$

Property:
have closed-form sol$^n$

**Rmk:**
penalize more on large value than $\ell_1$-reg, But cannot guarantee sparsity!

$\implies \hat{w}_{ridge} = (X^TX + \alpha \mathbb{1})^{-1} X^T y$

Note: $\hat{w}_{LR} = (X^TX)^{-1} X^T y$

$X \in \mathbb{R}^{n \times m}$

$X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix}$

$\implies X^TX = \sum_{i=1}^{n} x_i x_i^T$

Covariance matrix

we have: $\|\hat{w}_{LR}\| \geq \|\hat{w}_{ridge}\|$

**Pf:** consider $X^TX$ (symmetric matrix)

$\downarrow$

eigen value decomp.

$$\begin{cases} \text{eigenvalue} & \lambda_1 \geq \cdots \geq \lambda_m \\ \text{eigenvector} & u_1, \ldots, u_m \end{cases}$$

Principal Component
(direction)

then: $\hat{w}_{LR} = (X^TX)^{-1} X^T y$

$\quad = (X^TX)^{-1} \cdot \sum_{i=1}^{m} \beta_i u_i \implies \|\hat{w}_{LR}\|_2^2 = \sum_{i=1}^{m} \frac{\beta_i^2}{(\lambda_i)^2}$

$\beta_i = \langle u_i, X^T y \rangle$

$\quad = \sum_{i=1}^{m} \frac{\beta_i}{\lambda_i} u_i$

$\hat{w}_{ridge} = (X^TX + \alpha \mathbb{1})^{-1} X^T y$

$\quad = (X^TX + \alpha \mathbb{1})^{-1} \sum_{i=1}^{m} \beta_i u_i \implies \|\hat{w}_{ridge}\|_2^2 = \sum_{i=1}^{m} \frac{\beta_i^2}{(\lambda_i + \alpha)^2}$

$\quad = \sum_{i=1}^{m} \frac{\beta_i}{\lambda_i + \alpha} u_i$

**Note:** if $\alpha$ is large, then $\frac{\beta_i}{\lambda_i + \alpha} \to 0$

$\implies \hat{w}_{ridge} \to 0$

Also, $\ell_2$-regularization can be viewed as <u>Weight Decay</u>
   for non-linear case ( <u>we use GD</u>)

$$\theta^{(k+1)} = \theta^{(k)} - \varepsilon \nabla_\theta \, R_{emp} (\theta^{(k)}) - \underbrace{\varepsilon \cdot \alpha \, \theta^{(k)}}_{\text{weight decay part}}$$

<u>In application</u>, <u>Weight Decay</u> $\iff$ <u>$\ell_2$-regularization</u>

Statistics
Perspective   $\boxed{\text{Laplacian prior on } w} \nwarrow$

b) $\ell_1$ - regularization ( non-smooth but convex)

<u>Toy example</u>: ( LASSO ( linear) regression )

$$\begin{cases} \Omega(w) = \|w\|_1 \\ R_{emp}(w) = \frac{1}{2} \|Xw - y\|_2^2 \end{cases} \implies \boxed{\text{no closed-form sol}^n}$$

<u>Example</u>: $\ell_1$ v.s. $\ell_2$ regularization

$\longrightarrow \boxed{\theta_i^* = \underset{\theta}{\arg\min} \; R_{emp}(\theta)}$

$$R_{emp}(\theta) = \frac{1}{2} \sum_{i=1}^m \lambda_i (\theta_i - \theta_i^*)^2 \qquad \longrightarrow \underline{\text{without reg}}$$

1. <u>$\ell_2$-norm</u>: $\boxed{\widetilde{R}_{emp}(\theta)} = \frac{1}{2} \sum_{i=1}^m \lambda_i (\theta_i - \theta_i^*)^2 + \frac{1}{2}\alpha \sum_{i=1}^m \theta_i^2$
   
   *convex + differentiable ( smooth)*

$$\frac{\partial \widetilde{R}_{emp}}{\partial \theta_i} = \lambda_i (\theta_i - \theta_i^*) + \alpha \theta_i$$

(convexity)

$$\hat{\theta}^{\ell_2} \in \underset{\theta}{\arg\min} \; \widetilde{R}_{emp}(\theta) \iff \frac{\partial \widetilde{R}_{emp}}{\partial \theta_i}(\hat{\theta}^{\ell_2}) = 0$$

$$\implies \hat{\theta}_i^{\ell_2} = \frac{\lambda_i}{\lambda_i + \alpha} \theta_i^*$$

2. <u>$\ell_1$-norm</u> $\boxed{\widetilde{R}_{emp}(\theta)} = \frac{1}{2} \sum_{i=1}^m \lambda_i (\theta_i - \theta_i^*)^2 + \alpha \sum_{i=1}^m |\theta_i|$

$\longrightarrow$ convex but non-smooth

$$\hat{\theta}^{\ell_i} \in \arg\min_{\theta} \; \widetilde{R}_{emp}(\theta)$$

(highly non-trivial)

$$\begin{pmatrix} \partial \|\cdot\|_1(\hat{\theta}_1^{\ell_i}) \\ \vdots \\ \partial \|\cdot\|_1(\hat{\theta}_m^{\ell_i}) \end{pmatrix}$$

$$\Leftrightarrow \quad 0 \in \partial \widetilde{R}_{emp}(\hat{\theta}^{\ell_i})$$

$$\Leftrightarrow \quad 0 \in \begin{pmatrix} \lambda_1(\hat{\theta}_1^{\ell_i} - \theta_1^*) \\ \vdots \\ \lambda_m(\hat{\theta}_m^{\ell_i} - \theta_m^*) \end{pmatrix} + \alpha \; \boxed{\partial \|\cdot\|_1(\hat{\theta}^{\ell_i})}$$

$$\Leftrightarrow \quad \theta_i \in (\hat{\theta}_i^{\ell_i} - \theta_i^*) + \frac{\alpha}{\lambda_i} \partial \|\cdot\|(\hat{\theta}_i^{\ell_i})$$

$$\Leftrightarrow \quad \hat{\theta}_i^{(\ell)} = P_{\frac{\alpha}{\lambda_i}\|\cdot\|_1}(\theta_i^*)$$

$$= \begin{cases} \theta_i^* - \frac{\alpha}{\lambda_i}, & \theta_i^* \geq \frac{\alpha}{\lambda_i} \\ 0, & o/w. \\ \theta_i^* + \frac{\alpha}{\lambda_i}, & \theta_i^* \leq -\frac{\alpha}{\lambda_i} \end{cases}$$

$$u = P_f(x) = \arg\min_y f(y) + \tfrac{1}{2}\|y - x\|_2^2$$

$$\Leftrightarrow \quad 0 \in u - x + \partial f(u)$$

---

② **Regularization on NN**

→ we seldom regularize on <u>bias term</u> $b$

→ we may choose <u>different strength of regularization</u>   for each layer

$$\boxed{\alpha_i \to i\text{-th layer}}$$

---

③ **Early Stopping for NN** ⟶ <mark>under certain assumption,</mark>
<mark>it is equivalent to</mark> <u>$\ell_2 - reg$</u> !

<span style="color:red">implicit regularization</span>

$\downarrow$

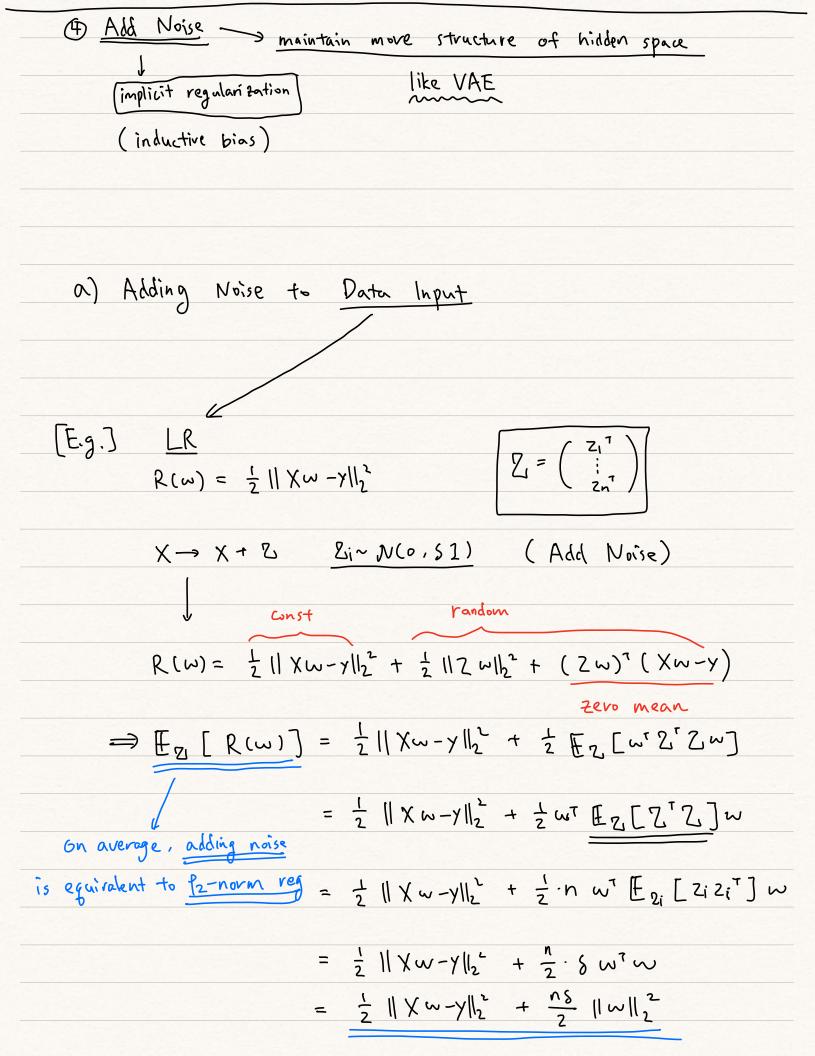$\boxed{\text{require validation set}}$ !!

<u>Normally</u>, we use validation set to monitor the time point
to stop !

**Variant1** : record the optimal epoch number

$\downarrow$

Stop criterion

[retrain]

**Variant2** : continue training with full dataset after early stop

$\downarrow$

record optimal loss function value (training)

stop criterion

[E.g.] Early Stop for Linear Reg

$$R_{emp}(\theta) = \frac{1}{2}\lambda(\theta - \theta^*)^2 \implies \nabla R_{emp}(\theta) = \lambda(\theta - \theta^*)$$

Consider **GD** :
$$\theta_{k+1} = \theta_k - \varepsilon\lambda(\theta_k - \theta^*)$$
$$= (1 - \varepsilon\lambda)\theta_k + \varepsilon\lambda\theta^*$$

$$\implies \theta_{k+1} = (1-\varepsilon\lambda)^{k+1}\theta_0 + \left[1 - (1-\varepsilon\lambda)^{k+1}\right]\theta^*$$

$\rightarrow$ Stop at iteration $l$ :  $\hat{\theta} = \theta_l = \boxed{(1-\varepsilon\lambda)^l}\,\theta_0 + \left[1 - (1-\varepsilon\lambda)^l\right]\theta^*$

(variant)

$\rightarrow$ $L_2$-regularization :  $\tilde{R}(\theta) = \frac{1}{2}\lambda(\theta - \theta^*)^2 + \frac{1}{2}\alpha(\theta - \theta_0)^2$

$$\nabla\tilde{R}(\theta) = \lambda(\theta - \theta^*) + \alpha(\theta - \theta_0) = 0$$

$$\implies \tilde{\theta} = \boxed{\frac{\alpha}{\alpha+\lambda}}\,\theta_0 + \left(1 - \frac{\alpha}{\alpha+\lambda}\right)\theta^*$$

Note $\hat{\theta} = \tilde{\theta} \iff \dfrac{\alpha}{\alpha+\lambda} = (1-\varepsilon\lambda)^l \iff \boxed{\alpha = \dfrac{\lambda(1-\varepsilon\lambda)^l}{1 - (1-\varepsilon\lambda)^l}}$

[early Stop] $\xleftarrow{\text{equivalent}}$ [$l_2$-regularization]

under certain condition $\{$ LR model

manual regularization strength

[regularization strength]

④ **Add Noise** $\longrightarrow$ maintain move structure of hidden space

$\downarrow$

implicit regularization           like VAE

( inductive bias )

a) Adding Noise to **Data Input**

[E.g.]   **LR**

$$R(w) = \frac{1}{2} \| Xw - y \|_2^2$$

$$Z = \begin{pmatrix} z_1^T \\ \vdots \\ z_n^T \end{pmatrix}$$

$$X \longrightarrow X + Z \qquad z_i \sim N(0, \delta 1) \qquad (\text{Add Noise})$$

$\downarrow$

$$R(w) = \underbrace{\frac{1}{2} \| Xw - y \|_2^2}_{\text{const}} + \underbrace{\frac{1}{2} \| Zw \|_2^2 + (Zw)^T(Xw - y)}_{\text{random}}$$

$\underbrace{\phantom{(Zw)^T(Xw-y)}}_{\text{zero mean}}$

$$\Rightarrow \underline{\mathbb{E}_Z[R(w)]} = \frac{1}{2} \| Xw - y \|_2^2 + \frac{1}{2} \mathbb{E}_Z[w^T Z^T Z w]$$

On average, adding noise is equivalent to $\ell_2$-norm reg

$$= \frac{1}{2} \| Xw - y \|_2^2 + \frac{1}{2} w^T \underline{\underline{\mathbb{E}_Z[Z^T Z]}} w$$

$$= \frac{1}{2} \| Xw - y \|_2^2 + \frac{1}{2} \cdot n \, w^T \mathbb{E}_{z_i}[z_i z_i^T] w$$

$$= \frac{1}{2} \| Xw - y \|_2^2 + \frac{n}{2} \cdot \delta \, w^T w$$

$$= \frac{1}{2} \| Xw - y \|_2^2 + \frac{n\delta}{2} \| w \|_2^2$$

② Another approach $\Rightarrow$ <u>Label Smoothing</u>

③ Adding Noise to <u>Weight</u>