

LEC7 DSAS204

Re-cap:

Regularization $\begin{cases} \downarrow \text{training loss} \\ \uparrow \text{generalization} \end{cases}$

- { ① parameter norm penalty
- ② early stop
- ③ inject noise

Note: up to now, we do not change the model architecture.

1. Ensemble

a) Bagging (Bootstrap Aggregating)

→ train m models through different subsets of data

→ then, aggregate the results

$\begin{cases} \text{regression : mean} \\ \text{classification : vote} \end{cases}$

idea: combine (uncorrelated) weak models to enhance performance
↓
small variance

→ Theoretical Analysis (Toy Example)

consider a simple scalar model:

$$f_i(x) = f^*(x) + \varepsilon_i(x)$$

assume

$$\begin{cases} \mathbb{E}[\varepsilon_i(x)] = 0 \\ \mathbb{E}[\varepsilon_i^2(x)] = \sigma(x)^2 \\ \mathbb{E}[\varepsilon_i(x) \varepsilon_j(x)] = 0 \end{cases}$$

aggregate model: $\bar{f}(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$

Expectation Error

consider:

$$\begin{cases} E(x) := \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(f_i(x) - f^*(x))^2] \\ \bar{E}(x) = \mathbb{E}[(\bar{f}(x) - f^*(x))^2] \end{cases}$$

$$\begin{aligned} E(x) &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\varepsilon_i^2(x)] \\ &= \sigma(x)^2 \end{aligned}$$

$$\begin{aligned} \bar{E}(x) &= \mathbb{E}\left[\left(\frac{1}{n} \sum_{i=1}^n \varepsilon_i(x)\right)^2\right] \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}[\varepsilon_i^2(x)] + \frac{1}{n^2} \sum_{i \neq j} \mathbb{E}[\varepsilon_i(x) \varepsilon_j(x)] \\ &= \frac{1}{n} \sigma(x)^2 \end{aligned}$$

Rmk: The unrealistic assumption is: uncorrelation of model

↓

Trade-off between Bias and Correlation

$$\begin{cases} \rightarrow \text{Large Sample Size} \Rightarrow \begin{cases} \text{Bias} \downarrow \\ \text{correlation} \uparrow \end{cases} \\ \rightarrow \text{Small Sample Size} \Rightarrow \begin{cases} \text{Bias} \uparrow \\ \text{correlation} \downarrow \end{cases} \end{cases}$$

2. Dropout → can be viewed as "Ensemble"

efficient

(approximate)

overcome the disadvantage of "Ensemble"

Store many sub-models

★ Suggest to put Dropout at the end

Dropout → apply random mask to weights

Average to achieve aggregation

Mathematically, it can be formulated as:

$$f(x; W, b, \mu) = G(W \cdot (\mu \circ x) + b)$$

node mask

for each layer

→ μ can be viewed as d-dim Bernoulli R.V.

→ $\mu_i \sim \text{Bernoulli}(p)$

$\begin{cases} \text{w.p. } p, & \text{we keep the node} \\ \text{w.p. } (1-p), & \text{we ignore the node (mask)} \end{cases}$

→ this only works for "Training Step"

→ While Testing (inference), replace W with $p \circ W$

Rmk:

1. The Entire Model Ensemble is : (Stochastic Version)

$$\{ f(x; \theta, \mu) : \mu_i \in \{0, 1\} \quad i \in [m] \}$$

2. There are 2^m models (distinct)

3. each of the model is not independent since they share the same θ

4. Ensemble $\rightarrow \frac{1}{n} \sum_{i=1}^n f_i(x)$

Stochastic Ensemble $\rightarrow \underline{\underline{\mathbb{E}_{\mu} [f(x; \theta, \mu)]}}$

\rightarrow in training, use Dropout to Approximate $\mathbb{E}_{\mu} [\cdot]$ & introduce independence to decrease correlation

\rightarrow in testing (inference), use "Weight Scaling Inference" to approximate $\mathbb{E}_{\mu} [\cdot]$

\swarrow

$W \rightarrow p \circ W$

\rightarrow Theoretical Analysis (Toy Example)

\downarrow

We want to show:

for LR, Dropout \iff L2-norm regularization

\downarrow

(Simple Case)

[E.g.] \rightarrow Linear Regression with Dropout

$$f(x; w, \mu) = \boxed{\frac{1}{p}} w^T (\mu \circ x) \quad \text{compensate for } \mathbb{E}[\mu_j] = p$$

here $\mu_j \sim \text{Bernoulli}(p)$ $j \in [d]$

\Rightarrow Empirical Risk (w.r.t $\mu \sim \text{Bernoulli}(p)$)
 \rightarrow we want to minimize to achieve \hat{w}

$$\tilde{R}(w) = \mathbb{E}_{\mu} \left[\frac{1}{2} \| X \cdot \left(\frac{1}{p} D_{\mu} \right) \cdot w - y \|_2^2 \right]$$

Here, $D_{\mu} = \text{diag}(\mu_1, \dots, \mu_d)$

$$= \mathbb{E}_{\mu} \left[\frac{1}{2} \| \underbrace{Xw - y}_{\text{standard LR}} + \underbrace{X \left[\frac{1}{p} D_{\mu} - I \right] w}_{\text{new term}} \|_2^2 \right]$$

$$= \frac{1}{2} \| Xw - y \|_2^2 + \frac{1}{2} \mathbb{E}_{\mu} \left[\| X \left[\frac{1}{p} D_{\mu} - I \right] w \|_2^2 \right]$$

$$+ \mathbb{E}_{\mu} \left[(Xw - y)^T X \left[\frac{1}{p} D_{\mu} - I \right] w \right]$$

$$\boxed{\mathbb{E} \left[\frac{1}{p} D_{\mu} \right] = I}$$

$= 0$

$$= \frac{1}{2} \| Xw - y \|_2^2 + \frac{1}{2} \mathbb{E}_{\mu} \left[\| X \left[\frac{1}{p} D_{\mu} - I \right] w \|_2^2 \right]$$

$R(p)$

$$\underline{R(p) := \frac{1}{p} D_{\mu} - I_d}$$

$$= \frac{1}{2} \| Xw - y \|_2^2 + \frac{1}{2} w^T \mathbb{E}_{\mu} \left[R(p)^T \underline{X^T X} R(p) \right] w$$

$$\downarrow$$
$$\underline{X^T X \in \mathbb{R}^{d \times d}}$$

(covariance matrix)

$$\text{consider } \left(\mathbb{E}_{\mu} \left[R(p)^T X^T X R(p) \right] \right)_{ij} \quad \left(R(p)_{ij} = \delta_{ij} \left(\frac{\mu_i}{p} - 1 \right) \right)$$

$$= \mathbb{E}_{\mu} \left[\sum_{k, \ell} \left(\frac{\mu_i}{p} - 1 \right) \delta_{ik} (X^T X)_{k\ell} \left(\frac{\mu_{\ell}}{p} - 1 \right) \delta_{\ell j} \right]$$

$$= (X^T X)_{ij} \mathbb{E}_\mu \left[\left(\frac{\mu_i}{p} - 1 \right) \left(\frac{\mu_j}{p} - 1 \right) \right]$$

$$= (X^T X)_{ij} \cdot \delta_{ij} \left(\frac{1}{p} - 1 \right) \underbrace{\frac{\text{var}(\mu_i)}{p^2}}_{\text{var}(\mu_i)/p^2}$$

$$= \frac{1}{2} \|Xw - y\|_2^2 + \frac{1}{2} \left(\frac{1}{p} - 1 \right) \sum_{i=1}^d (X^T X)_{ii} w_i^2$$

\Rightarrow weighted L_2 -norm regularization
(the weight is data-driven)

3. Batch Normalization (BN)

$$\hat{y}(x; w) = x w_1 w_2 \dots w_L$$

① Motivation: [1-D Deep LR] \rightarrow In slide!

② Idea: 1. For Deep NN, Learning Rate depends on the scale of weights & activations at different layers

2. A constant learning rate cannot fulfil our requirements as training goes on!

③ BN Layer:

$H = \{h^{(1)}, \dots, h^{(B)}\}$ is a batch of hidden state values

$$\rightarrow \text{BN}(h^{(i)}; \gamma, \beta) = \gamma \odot \left(\frac{h^{(i)} - \mu}{\sigma} \right) + \beta$$

$$= \sigma \odot (\tilde{h}^{(i)}) + \beta$$

$$\tilde{h}^{(i)} \rightarrow \begin{cases} \text{mean } 0 \\ \text{std } 1 \end{cases}$$

$$\Rightarrow H = \{h^{(1)}, \dots, h^{(B)}\} \rightarrow \begin{cases} \mu & \text{mean} \\ \sigma^2 & \text{variance} \end{cases}$$

$\Rightarrow \sigma$ & β are learnable parameters

④ Training & Inference

→ Training : simple

→ Inference (Testing) : use moving average (MA)

↓

record recent B samples and update

⑤ Effect (observation & intuition)

→ break the inter-dependency of different layers

→ the scaling is controlled individually with σ & β

→ It is far beyond naive scaling operation

→ σ & β are learnable

→ gradient can flow!

4. Data Augmentation

→ input-output dist: $(x, y) \sim \mu$

→ $(x^{(i)}, y^{(i)}) \sim \mu \quad i=1, 2, \dots, N$ (draw N samples)

→ Empirical Loss

$$R_{\text{emp}}(\theta) = \frac{1}{N} \sum_{i=1}^N L(f(x^{(i)}, y^{(i)}))$$

→ Expected Loss

$$R_{\text{pop}}(\theta) = \mathbb{E}_{(x,y) \sim \mu} [L(f(x), y)]$$

(Statistical)

⇒ Learning Theory tells us:

$$\text{generalization gap} \leq \frac{\text{Complexity of Model}}{\text{\# of samples}}$$

⇒ better generalization



{ reduce complexity
move data

⇒ More Data ← Data Augmentation



Introduce Prior Knowledge to model

(Soft Encode Prior Knowledge)

5. Learning Rate Decay



Control the Noise for SGD

→ Choice Criterion for ϵ_k

(k^{-1} decay schedule)

$$\begin{cases} \sum \epsilon_k = +\infty \\ \sum \epsilon_k^2 < +\infty \end{cases}$$

$$\longrightarrow \epsilon_k = \frac{a}{1+bk}$$

(one choice)

→ In practice, we do not favor this choice

→ one reason: our ultimate goal is not $\min_{\theta} R_{\text{emp}}(\theta)$

but $\min_{\theta} R_{\text{pop}}(\theta)$

→ one choice is: exponential decay schedule

$$\xi_k = \xi_0 \gamma^{\lfloor \frac{k}{k_0} \rfloor}$$

initial learning rate: ξ_0

decay ratio: $\gamma \in [\frac{1}{10}, \frac{1}{2}]$

iteration to wait before
decay: k_0