

1. Robustness and Regularization

1.1 Adversarial Examples.

1.1.2 Prediction under attack.

$$w \in \mathbb{R}^n, x \in \mathbb{R}^n$$

$$\nabla_x f(x; w) = \frac{\partial w^T x}{\partial x} = w$$

$$\Rightarrow x' \leftarrow x - \epsilon \cdot w = x - \epsilon w$$

$$\Rightarrow f(x'; w) = w^T x' = \boxed{w^T (x - \epsilon w)}$$

1.2 Gradient Descent and Weight Decay

1.2.2 Closed form Ridge Regression Solution (ignore $\frac{1}{2n}$ for convenience).

$$\min_w \|xw - t\|_2^2 + \lambda \|w\|_2^2$$

$$\text{Expand} \Rightarrow (xw - t)^T (xw - t) + \lambda w^T w$$

$$= (w^T x^T - t^T) (xw - t) + \lambda w^T w$$

$$= (w^T x^T x w - w^T x^T t - t^T x w + t^T t) + \lambda w^T w$$

$$\text{take derivative w.r.t. } w \Rightarrow 2x^T x w - 2x^T t + 2\lambda w = 0.$$

$$(x^T x + \lambda I) w = x^T t$$

$$\Rightarrow \boxed{w_{\text{ridge}}^* = (x^T x + \lambda I)^{-1} x^T t}$$

1.2.3 Adversarial Attack under weight decay ID: $x \in \mathbb{R}^n, x \in \mathbb{R}$

$$\text{Given: } w_{\text{ridge}}^* = (x^T x + \lambda)^{-1} x^T t \quad f(x; w) = w^T (x - \epsilon w)$$

$$\Rightarrow f(x'; w_{\text{ridge}}^*) = [(x^T x + \lambda)^{-1} x^T t]^T [x - \epsilon (x^T x + \lambda)^{-1} x^T t] = 0$$

$$x - \epsilon (x^T x + \lambda)^{-1} x^T t = 0 \Rightarrow x = \epsilon (x^T x + \lambda)^{-1} x^T t$$

$$\boxed{\epsilon = \frac{x^T x + \lambda}{x^T t}}$$

\Rightarrow As $\lambda > 0$, the weight decay will make ϵ larger than the normal $\epsilon = \frac{x^T x}{x^T t}$
 which means it will perturb the data more and more robust to fool the model.

{2. Trade off resources in Neural Net training}

2.1 Effect of batch size

2.1.1

Batch size vs. learning rate

$$g_B(w) = \frac{1}{B} \sum_{i=1}^B g_i(w), \text{ where } g_i(w) = \nabla \mathcal{L}(w) + \epsilon_i, \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2).$$

From the minibatch gradient formula, it clearly shows that the noisy version of gradient is closer to the true gradient $\nabla \mathcal{L}(w)$ as the batch size B increases. which indicates more 'accurate' gradient can be generated by increasing the batch size. Therefore, the optimal learning rate will also increase with batch sizes since it can take large steps without losing accuracy and also accelerate the training time.

2.1.2

Training steps vs. Batch size.

(a) C is the most efficient batch size among other two points since A requires much more training time to finish the entire process since the batch size is too small and the results may contain a lot noises; B requires much more resources since larger batch size requires more computational power which will induce more cost and slower convergence.

(b) Point A: noise dominated \Rightarrow seek parallel compute.

Point B: curvature dominated. \Rightarrow use higher order optimizers

2.2 Model size, data size and compute

(a) C would be the best option. By looking at the right figure in Figure 2

and Figure 3, it clearly shows that the optimal model size is increasing with more compute and increasing model size will significantly increase the model performance. "A" may induce overfitting since some model with more training steps may make the test loss getting higher; "B" does not make sense since the fine-tuned learning rate in smaller batch sizes does not fit for larger batch size due to its inefficiency.

3. Dropout and Gaussian Noise

3.1 Linear regression with input dropout

$$\tilde{y}_m^{(i)} = \frac{1}{p} \sum_j m_j^{(i)} w_j x_j^{(i)}, \quad \text{where } m_j^{(i)} \stackrel{\text{i.i.d.}}{\sim} \text{Ber}(p).$$

$$\begin{aligned}
E_m[J] &= \frac{1}{2N} \sum_{i=1}^N E_m[(\tilde{y}_m^{(i)} - t^{(i)})^2] \\
&= \frac{1}{2N} \sum_{i=1}^N E_m[\tilde{y}_m^{(i)2} - 2\tilde{y}_m^{(i)}t^{(i)} + t^{(i)2}] \\
&= \frac{1}{2N} \sum_{i=1}^N E_m(\tilde{y}_m^{(i)2}) - \frac{1}{2N} \sum_{i=1}^N E_m[2\tilde{y}_m^{(i)}t^{(i)}] + \frac{1}{2N} \sum_{i=1}^N t^{(i)2} \\
&= \frac{1}{2N} \sum_{i=1}^N \text{Var}[\tilde{y}_m^{(i)}] + \frac{1}{2N} \sum_{i=1}^N E_m[\tilde{y}_m^{(i)}]^2 - \frac{1}{2N} \sum_{i=1}^N E_m[2\tilde{y}_m^{(i)}t^{(i)}] + \frac{1}{2N} \sum_{i=1}^N t^{(i)2} \\
&= \frac{1}{2N} \sum_{i=1}^N (\underbrace{E_m[\tilde{y}_m^{(i)}] - t^{(i)}}_{\text{①}})^2 + \frac{1}{2N} \sum_{i=1}^N \underbrace{\text{Var}[\tilde{y}_m^{(i)}]}_{\downarrow p(1-p)}
\end{aligned}$$

$$\text{Var}[\tilde{y}_m^{(i)}] = \sum_j \text{Var}\left[\frac{1}{p} m_j^{(i)} w_j x_j^{(i)}\right]$$

\downarrow
 $p(1-p)$

then, ① can be written as:

$$= \frac{1}{2N} \sum_{i=1}^N (\underbrace{E_m[\tilde{y}_m^{(i)}] - t^{(i)}}_{\text{②}})^2 + \frac{(1-p)}{2p} \sum_j \text{Var}[x_j] w_j^2 \quad \text{②}$$

3.2 Multiplicative Gaussian noise

$$\tilde{y}_\pi^{(i)} = \sum_j (1 + \pi_j^{(i)}) w_j x_j^{(i)}, \quad \text{where } \pi_j^{(i)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2).$$

Similar to ①: $\frac{1}{2N} \sum_{i=1}^N (\mathbb{E}_\pi[\tilde{y}_\pi] - t_i)^2 + \frac{1}{2N} \sum_{i=1}^N \text{Var}[\tilde{y}_\pi]$ ③

first term: ②: $\mathbb{E}[\tilde{y}_\pi] = \mathbb{E}\left[\frac{1}{P} \sum_j m_j w_j x_j\right] = \sum_j w_j x_j$ since $\mathbb{E}[m_j] = P$

③: $\mathbb{E}[\tilde{y}_\pi] = \mathbb{E}\left[\sum_j (1 + \pi_j) w_j x_j\right] = \sum_j w_j x_j$ since $\mathbb{E}_\pi[1 + \pi_j] = \mathbb{E}_\pi[1 + 0] = 1$

second term: $\frac{1}{2N} \sum_{i=1}^N \sum_j \text{Var}[(1 + \pi_j) w_j x_j]$
 \downarrow
 $\text{Var}[1 + \pi_j] = \sigma^2 \Rightarrow = \frac{\sigma^2}{2} \sum_j \text{Var}[x_j] w_j^2$

therefore, ③ can be written as.

$$\frac{1}{2N} \sum_{i=1}^N (\mathbb{E}_\pi[\tilde{y}_\pi] - t_i)^2 + \frac{\sigma^2}{2} \sum_j \text{Var}[x_j] w_j^2 \quad ④$$

Compare ③ and ④: when $\sigma = \sqrt{\frac{1-P}{P}}$

\Rightarrow multiplicative gaussian noise is same as applying input dropout and they share the same structure of bias-variance decomposition.