## 6.2 Weight norm penalty

Change on the cost function:

old cost function : $L(\underline{\theta})$   based on KL divergence

New regularized cost f.ⁿ :

$$L_r(\underline{\theta}) = L(\underline{\theta}) + \sum_{l=1}^{L} \lambda_l \, P(\underline{\underline{w}}_l)$$

$P(\underline{\underline{w}}_l) \geqslant 0$ : <u>Penalty terms</u> ; Penalize $\underline{\theta}$ with large $P(\underline{\underline{w}}_l)$

$\lambda_l \geqslant 0$   : <u>Regularization parameters</u>

→ A compromize between $\min L(\underline{\theta})$ and $\min P(\underline{\underline{w}}_l)$
→ no unique sol.ⁿ as a multiobjective, conflicting problems.

$\lambda_l = 0 \; \forall l$ :  no regularization
$\lambda_l$ is high → then $L(\underline{\theta})$ gets impacted.

<u>Common choice of $P(\underline{\underline{w}}_l)$</u>

(a) l2 regularization : use L2 norm of $vec(\underline{\underline{w}})$

$$P(\underline{\underline{w}}_l) = \| vec(\underline{\underline{w}}) \|_2^2 = \sum_j \sum_i w_{l,ij}^2 \; \hat{=} \; \text{weight energy}$$

→ prefer $\underline{\theta}$ with small weight energy

<u>Analysis</u>:
$$L_r(\underline{\theta}) = L(\underline{\theta}) + \lambda \|\theta\|^2 \quad \text{for simplicity.}$$

$$\nabla L_r(\underline{\theta}) = \nabla L(\underline{\theta}) + 2\lambda \underline{\theta}$$

$$\underline{\theta}^{t+1} = \underline{\theta}^t - \gamma^t . \; \nabla L_r(\underline{\theta})$$

L.2 regularization leads to weight decay during training

try to lower the → original weight ...

$$= \underline{\theta}^t - \gamma^t \nabla L(\underline{\theta}) - 2\lambda\gamma^t \underline{\theta}^t = (1 - 2\lambda\gamma^t)\underline{\theta}^t - \gamma^t \nabla L(\underline{\theta})$$

(b) $L1$ - regularization : use $l1$ norm of $\text{vec}(\underline{\underline{w}}_\lambda)$

$$P(\underline{\underline{w}}_\lambda) = ||\text{vec}(\underline{\underline{w}}_\lambda)||_1 = \sum \sum |w_{\lambda, ij}|$$

→ prefer sparse $\underline{\underline{w}}_\lambda$ with many zero elements.

Bias $\underline{b}_\lambda$ :
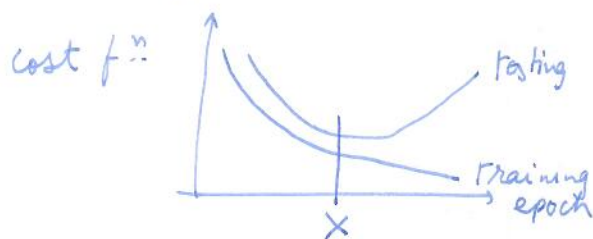
• no amplification of the input vector $\underline{x}_{\lambda-1}$

• no need for regularization

slide 6-5

6.3 <u>Early stopping</u>

Change on optimizer



stop the training. Here there is a divergence in the training error rate and test error rate

Slide 6-6

6.4 <u>Data augmentation</u>

change on dataset

Overfitting - more complex model which can memorize the training dataset. Theoretically infinite no of training data will never overfit any neural network

But we have limited dataset.

<u>Data augmentation</u> - Generate artificial but realistic training ~~angles~~ data / samples.

6.5    Ensemble learning

change on dataset / maodel / cost function / optimizer

slide 6-8


6.6    Dropout

change in model

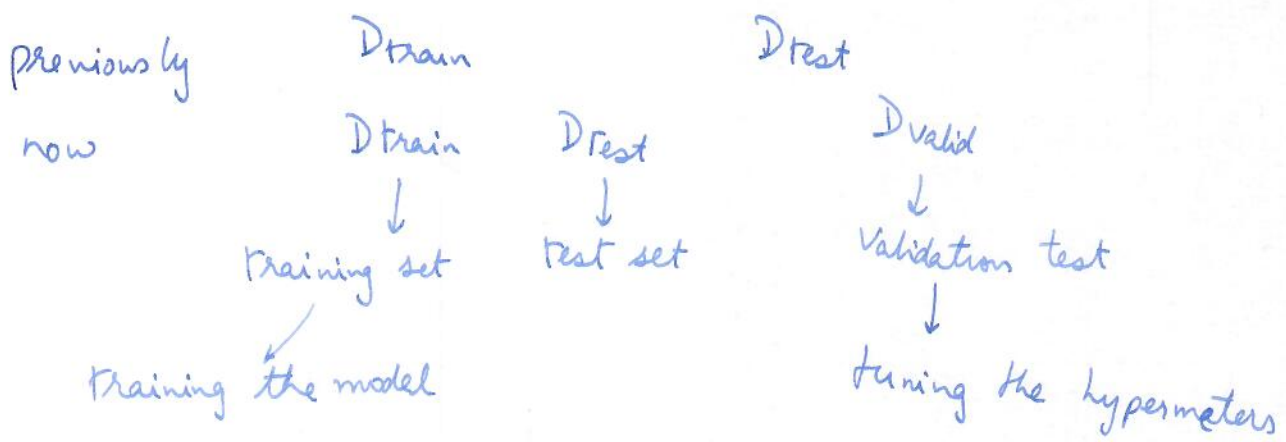An implicit ensemble learning method

Slide 6-9, 6-10

Co-adapted — ?


6.7    Hyperparameters optimization

6-17, 18, 19
      Slide

(a) Can't do ~~so~~ many hyper parameter optimization as they are
    mostly integer optimization. or discrete value optimisation

(b) Can be solved by test set or validation set

      dataset        D

previously        Dtrain              Dtest

now        Dtrain      Dtest              Dvalid
              ↓          ↓                  ↓
         training set  test set        Validation test
              ↙                              ↓
    training the model            tuning the hypermeters

    Early stopping is a type of hyperparameter optimization on
    no of epochs.

# Training and hyperparameter optimization

for $\underset{\sim}{\theta}^{\eta} = \ldots$

- learn $\underline{\theta}$ ~~$\theta\theta$~~ of $f(\underline{x}; \underline{\theta}, \underset{\sim}{\theta}^{\eta})$ from $D_{train}$

- Calculate validation error $\left(\underset{\sim}{\theta}\right)$ of $f(\underline{x}; \underline{\theta}, \underset{\sim}{\theta}^{\eta})$

   on $D_{val}$

end

$\min\limits_{\eta}$ validation error $(\underset{\sim}{\theta} \, \eta)$

Calculate Test error of $f(\underline{x}; \underline{\theta}; \underline{\eta})$ on $D_{test}$

If $err_{test} > err_{train}$ or $err_{val}$    then use more well-defined

optimization techniques.

Slide 6-20

Q Any mathematical approach for Bayesian optimisation?

Ilias