1. linear regression

$$f(x) = w_1 x_1 + \cdots + w_k x_k + b \longrightarrow f(x) = w^T x + b \quad \nearrow \hat{y}_i$$

least square method : $(w^*, b) = argmin_{(w,b)} \sum_{i=1}^{m} (f(x_i) - y_i)^2$

i.e. $MSE = \frac{1}{m} \| \hat{y} - y \|_2^2 = \frac{1}{m} \sum_{i=1}^{m} (\hat{y}_i - y_i)^2 \quad \Big\} E(w,b)$

$$E(w,b) = \sum_{i=1}^{m} (y_i - w x_i - b)^2$$

to make E the least : there are :

$$\begin{cases} \frac{\partial E}{\partial w} = \sum_{i=1}^{m} (2w x_i^2 - 2(y_i - b) x_i) = 0 \\ \frac{\partial E}{\partial b} = \sum_{i=1}^{m} (2b - 2(y_i - w x_i)) = 0 \end{cases}$$

another formate : $E(w,b) = E(\hat{w}) = argmin_{\hat{w}} (y - X\hat{w})^T \cdot (y - X\hat{w})$

$$\frac{dE}{d\hat{w}} = 0 \Rightarrow \hat{w} = (X^T X)^{-1} X^T y$$

2. Bias - variance trade-off

$\begin{cases} Bias \ 偏差 \quad bias(\hat{\theta}_m) = E(\hat{\theta}_m) - \theta \qquad 預測值期望 \\ Variance \ 方差: \ Var(\hat{\theta}) \longrightarrow Standard \ error \ [Var(\hat{\theta})]^{1/2} \end{cases}$
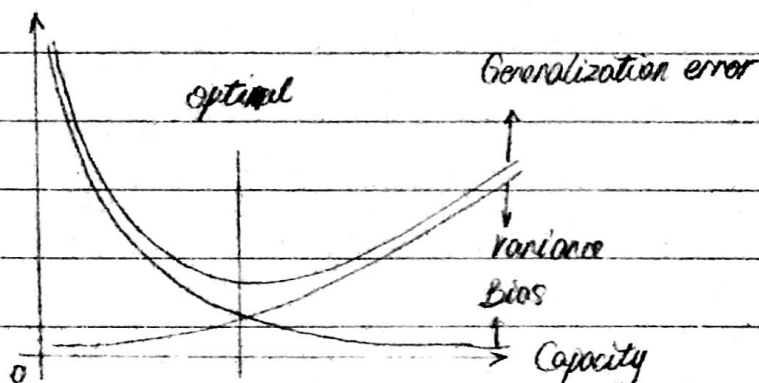
$\rightarrow$ MSE : mean square error 均方误差

$\begin{cases} S^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i - M)^2 & \text{Variance} \\ SSE = \sum_{i=1}^{m} w_i (y_i - \hat{y}_i)^2 & \text{The sum of squares due to error} \\ MSE = \frac{1}{n} SSE & \sim \\ RMSE = \sqrt{MSE} & \text{Root MSE} \end{cases}$

★ $MSE = E(\hat{\theta} - \theta)^2 = Bias^2 + Var \quad ?$

3. Backpropagation

An example: $a : a = 2 \xrightarrow{1} c : c = a + b \xrightarrow{2} e : e = c * d$

$b : b = 1 \xrightarrow{1} d : d = b + 1 \xrightarrow{3}$

Then use BP algorithm:  $c : (2)$  $d : (3)$  ← layer 2

$a : (2)$  $b : (2+3)$  ← layer 1

N.B.  $b : 2 * 1 + 3 * 1 = 5$

4. Early stopping

Generalization Performance  泛化性能

How to deal with overfitting
→ Cutting down parameters (parameter tying & sharing
↘ Cutting down dimensions. (weight decay & early stop

Algorithm:    train & output validation error

(improved) ↙          ↘ (not improved) → maybe for some times

store a copy of weights          Return the copy of weights
and go on.

5. Bagging

when to use it : if your models have <u>low biases</u> and <u>high variances</u>
(which are more easily to become overfit )

Algorithm : 1. select training dataset from original data set
N.B. the selecting is based on Bootstraping method (?)
so there may be a case that a part of dataset is never used
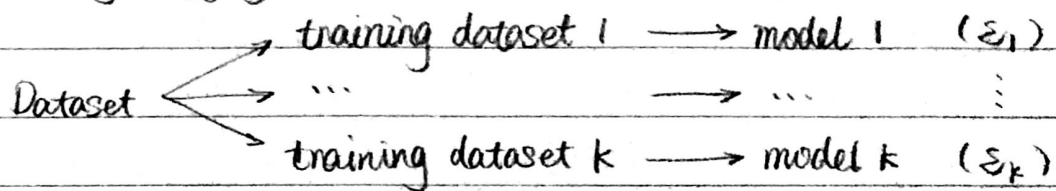
2. Each train dataset will be trained to a model

3. Summary such models and get the final mode
N.B. for classify model, we may use voting method
for regression methods, we use the mean for the result

The reason why Bagging performance better in results:

$$\text{Dataset} \begin{cases} \text{training dataset } 1 \longrightarrow \text{model } 1 \quad (\varepsilon_1) \\ \cdots \longrightarrow \cdots \\ \text{training dataset } k \longrightarrow \text{model } k \quad (\varepsilon_k) \end{cases}$$

We have : $E(\varepsilon_i^2) = V$ ; $E(\varepsilon_i \varepsilon_j) = C$

If we combine all these models :

$$E([\tfrac{1}{k} \Sigma_i \varepsilon_i]^2) = \tfrac{1}{k^2} E\left( \Sigma_i (\varepsilon_i^2 + \Sigma_{i \neq j} \varepsilon_i \varepsilon_j) \right)$$

$$= \tfrac{1}{k} V + \tfrac{k+1}{k} C. \qquad (*)$$

N.B. Perfectly correlated : $V = C$

Perfectly uncorrelated : $C = 0$.


## 6. Dropout

During forward propagation period, let neural nodes have certain probability to stop working (Dropout)


Algorithm: $r_j^{(l)} \sim \text{Bernoulli } (p)$

$$\tilde{y}^{(l)} = r^{(l)} * y^{(l)}$$

$$z^{(l+1)} = w^{(l+1)} \cdot \tilde{y}^{(l)} + b^{(l+1)}$$

$$y_i^{(l+1)} = f(z_i^{(l+1)})$$

MAKE UP POINTS

1. linear regression result dirivation      *ignored bias*

    MSE (mean square error) $= \frac{1}{m} \| \hat{y} - y \|_2^2$

    $\nabla MSE = 0 \Leftrightarrow \nabla \left[ \frac{1}{m} \sum (\hat{y}_i - y)^2 \right] = \nabla \left[ \frac{1}{m} \sum (y_i - w x_i)^2 \right] = 0$

         $\Leftrightarrow \frac{d}{dw} \sum (y_i - w x_i)^2 = \sum (x^2 w - xy - bx) * 2 = 0$

    OR $\Leftrightarrow \nabla \frac{1}{m} (y - wx)^T (y - wx) = \nabla \frac{1}{m} (y^T y - y^T wx - (wx)^T y + (wx)^T (wx)) = 0$

         $\Leftrightarrow \underline{-y^T x - x^T y} + 2w x^T x = -2 x^T y + 2w x^T x = 0$

         $\Leftrightarrow w = (X^T X)^{-1} X^T y$

                                           *prediction.*

2. bias $(\hat{\theta}_m) = E(\hat{\theta}_m) - \theta$     N.B. *The meaning of $\hat{\theta}_m$. (parameters, rather than*

    $Var(\hat{\theta}) = E \left[ (\hat{\theta} - E(\hat{\theta}))^2 \right]$

    $MSE = E \left[ (\hat{\theta} - \theta)^2 \right] = Bias(\hat{\theta})^2 + Var(\hat{\theta})$

       N.B. MSE means the difference between train & test

           while $Var(\hat{\theta})$ means the degree of results' aggregation

3. N.B pay attention to <u>Dirivative</u> and <u>partial Dirivative</u>

4. N.B. <u>validation data set required</u>

5 N.B. $E(\varepsilon_i^2) = V$ (Variance)     $E(\varepsilon_i \varepsilon_j) = C$ (Covariance)     协方差