

## Practical considerations:

① ②

When we are looking at local gradient in BPA, it depends on derivative of activation function. Hence the Act. fun must be continuous and differentiable. [Condition]

Therefore we use Sigmoidal non linear function. It has 2 forms :

① Logistic sigmoidal act. function.

$$\text{for a neuron } j : \phi_j(v_j(n)) = \frac{1}{1 + \exp(-a v_j(n))} \quad -\infty < v_j < \infty \quad a > 0 \text{ (condition)}$$

from restriction on boundary  $v_j$

wkT a: controls the slope of Act<sup>n</sup> function.  
 $v_j$ : induced local field.

$$\phi_j'(v_j(n)) = \frac{a \exp(-a v_j(n))}{[1 + \exp(-a v_j(n))]^2}$$

$$\text{wkT} \quad y_j(n) = \phi_j(v_j(n)).$$

$$\phi_j'(v_j(n)) = a y_j(n) [1 - y_j(n)] \quad \text{--- Q.}$$

If  $\phi$  neuron  $j$  is in output layer we simply write.

$$y_j(n) = \phi_j(v_j(n)). = o_j(n)$$

& local gradient:

$$\delta_j(n) = o_j(n) \phi_j'(v_j(n)).$$

$$= a [o_j(n) - o_j(n)] o_j(n) [1 - o_j(n)]$$

(numerical)  
If  $j$  is in hidden layer.

$$d_j(n) = \phi_j'(v_j(n)) \sum_k d_k(n) w_{kj}(n)$$

$$= a y_j(n) [1 - y_j(n)] \sum_k d_k(n) w_{kj}(n).$$

The range of o/p  $y_j(n)$  will be 0 to 1

when  $y_j(n) = 1/2$ , The  $\sum d_k(n) w_{kj}(n)$  will be equal to

$$\phi_j'(v_j(n)) = a \frac{1}{2} [1 - \frac{1}{2}] = a/4.$$

$$y_j(n) = 0 \quad \phi_j'(v_j(n)) = 0.$$

$$y_j(n) = 1 \quad \phi_j'(v_j(n)) = 0.$$

$\Rightarrow \phi_j'$  has maximum value at  $y_j(n) = 1/2$ .  
at this point the w value adjustment  
will be maximum in the  $0 \leq y \leq 1$ .

## ② Hyperbolic tangent function:

$$\phi_j(v_j(n)) = a \tanh(b v_j(n)) \quad (a, b) > 0.$$

$$\begin{aligned} \phi_j'(v_j(n)) &= ab \operatorname{sech}^2(b v_j(n)) \\ &= ab (1 - \tanh^2(b v_j(n))) \\ &= ab \frac{b}{a} [a - y_j(n)][a + y_j(n)] \end{aligned}$$

## Generalization & cross validation

If neuron  $j$  is at output layer  
local gradient

$$\delta_j(n) = g'(n) \phi'_j(v_j(n))$$

$$= \frac{b}{a} [y_{dj}(n) - g(n)] [a - o_j(n)] [a + g(n)]$$

If neuron  $j$  is in a hidden layer

$$\delta_j(n) = \phi'_j(v_j(n)) \sum_k \delta_k(n) w_{kj}(n)$$

$$= \frac{b}{a} [a - y_j(n)] [a + y_j(n)] \sum_k \delta_k(n) w_{kj}(n)$$

i. Local gradient may calculate the w/o requirement of activation function knowledge.

② Rate of learning:

If we see carefully the BPA it is nothing but modified / approximated version of steepest descent method.

The trajectory in perfect SPM is smooth but in BPA is not smooth as it depends on  $\eta$ . If  $\eta$  is small, trajectory will be smooth, but the cost of a slow rate of learning. On other hand if  $\eta$  is too large in learning speed will be fast but of large change in  $w$ , which may lead to unstable (oscillatory).

A simple way to overcome this method

ii by consider the following

$$\Delta w_{ji}(n) = \alpha \Delta w_{ji}(n-1) + \eta f_j(n)y_i(n) \quad \text{--- ①}$$

NKT :  $w(k+1) = w(k) + \Delta w(k)$

where  $\Delta w(k) = -\eta \frac{\partial E(k)}{\partial w}$  in general case

local gradient

$$f_j(n) = -e_j(n) \nabla \phi^T(v_f)$$

$$\frac{\partial E(k)}{\partial w} = -\eta f_j(n) y_j(n)$$

In Eq<sup>n</sup> ① 2nd term has  $\eta$  measures how much the current situation affects the next step.

1st term has  $\alpha$  : momentum which means how much past steps affect the next step.

Eq ① : generalized delta rule

Now find the solution for eq ① since it has first difference.

$$\Delta w(n) - \alpha \Delta w(n-1) = \eta f(n)y(n)$$

$$n=0$$

$$\Delta w(0) = \alpha \Delta w(0) + \eta f(0)y(0) = \eta f(0)y(0)$$

$$n=1$$

$$\Delta w(1) = \alpha \Delta w(1) + \eta f(1)y(1)$$

$$= \cancel{\alpha \Delta w(0)} + \alpha \eta f(0)y(0)$$

$$+ \eta f(1)y(1)$$

$$\begin{aligned} \Delta w(2) &= \alpha \Delta w(1) + \eta \delta(2)y(2) \\ &= \alpha^2 \Delta w(1) + \alpha \eta \delta(0)y(0) + \eta \delta(1)y(1) \\ &\quad + \eta \delta(2)y(2) \end{aligned}$$

If  $w(0) = 0$

and so on.



$$\Delta w(n) = \eta \sum_{t=0}^n \alpha^{n-t} \delta(t)y(t)$$

$$\Delta w(n) = -\eta \sum_{t=0}^n \alpha^{n-t} \frac{\partial \mathcal{E}(t)}{\partial w(t)}$$

The following observations can be made w.r.t above eq<sup>n</sup>.

- ① The above eq<sup>n</sup> converges if  $0 < |\alpha| < 1$ .  
If  $\alpha = 0$ , BPA operates w/o momentum.
- ② If  $\frac{\partial \mathcal{E}(t)}{\partial w(t)}$  has same sign in the consecutive iteration then the sum grows in magnitude as a result  $\Delta w$  will be adjusted by a large amount.
- ③ If  $\frac{\partial \mathcal{E}(t)}{\partial w(t)}$  has opposite sign in the consecutive iteration, sum will decrease. So  $w$  is adjusted by a small amount.

(3)

## Sequential and Batch mode of training:

The complete presentation of the entire training set during the learning process is called epoch.

The learning process is maintained on an epoch by epoch basis until the synaptic weights and bias level of the network stabilize and avg. sq. error over the entire training set converge to some minimum value.

BP learning may proceed in one of 2 ways.

① Sequential mode: It is also referred as on-line/pattern / stochastic mode.

In this mode of operation weight updating is performed after the presentation of each training sample.

② Batch mode:

here weight updating is performed after the presentation of all the training examples that constitute an epoch.

here we consider average squared error.

$$E_{av} = \frac{1}{2N} \sum_{n=1}^N \sum_{j \in c} e_j^2(n).$$

The  $e_j(n)$  pertains to ofp neuron  $j$  for training example  $n$ .

Inner summation: w.r.t  $j$  over all neurons in the ofp layer of  $N^{th}$  outer —————;  $n$  to  $n$  is performed over the entire training set in the epoch at hand.

$$\Delta w_{ji} = -\eta \frac{\partial^2 e_i}{\partial w_{ji}}$$

$$= -\eta \sum_{n=1}^N g_i(n) \frac{\partial e_i(n)}{\partial w_{ji}}$$

Advantages of sequential mode:

- ① It less storage for each synapses connection.
- ② Easy to implement.
- ③ Provides effective solution to large & difficult pblm.
- ④ ability to take the adv. of redundancy by presenting example one at a time.

Disadvantage:

- ① Stochastic nature of presenting samples to NW is a pblm ∵ it is trapped in a local minimum.

### ③ Stopping Criteria:

In general the BPA cannot be shown to converge and there are no well defined criteria for stopping its operation.

2 ways to stop BP learning

The BPA is considered to have converged

- (i) When the Euclidean norm of the gradient vector reaches sufficiently small gradient threshold.

Desirable:

- (ii') When the absolute rate of change in the averaged squared error per epoch is sufficiently small.

## Heuristic For Making the BPA perform better.

1. Sequential v/s batch mode update:

For large & redundant data, sequential processing is faster.

For smaller data & NW, B-processing is faster.

2. Maximizing information content:

Aim: Every training example presented to the BPA should be chosen on the basis that its information content is the largest possible for the task at hand

To achieve this aim there are 2 ways.

- ①. The use of an example that results in the largest training error.

∴ for updating w-values local gradient will be used which is a function of error.

When an example produces the larger error which leads to more larger local gradient, so faster will be the learning process.

- ②. The use of example that is radically different from all those previously used.

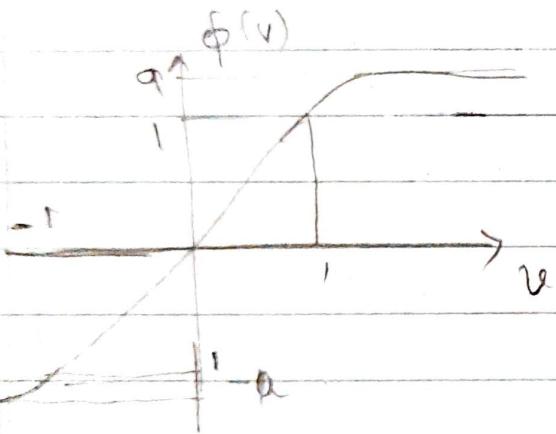
### 3. Activation function: choice of

As we discussed earlier activation function plays an imp. role in training NN using RPA.  
If we use tangent hyperbolic activation fun learning will be faster than logistic. Because THyperbolic is antisymmetric while as logistic is non symmetric about 0.

$$\text{if } \phi(-v) = -\phi(v).$$

$$\phi(v) = a \tanh(bv).$$

We can choose particular value for  $a = 1.7159$  and  $b = 2/3$ . Because the following useful opties.



$$\textcircled{1} \quad \phi(1) = 1$$

$$\phi(-1) = -1$$

\textcircled{2} At the origin the slope (i.e effective gain) of the activation fun is close to 1  
 $\phi'(0) = ab$ .

$$= 1.1424$$

\textcircled{3}  $\phi''(v) = 0$  attains its maximum value at  $v=0$ .

$|v| < 1$  : Region is linear

$|v| > 1$  : Region is defined as Saturation

#### ④ Target value:

It is important to choose the target value within the activation function range i.e.  $-1 \leq y_{dj} \leq 1$ . For example,

Let us assume the desired output at neuron  $i$  is  $y_{dj}$ , if its it should be offset by some amount  $\epsilon$ .

$$\text{ie } y_{dj} = a - \epsilon$$

$$\text{or } y_{dj} = -a + \epsilon$$

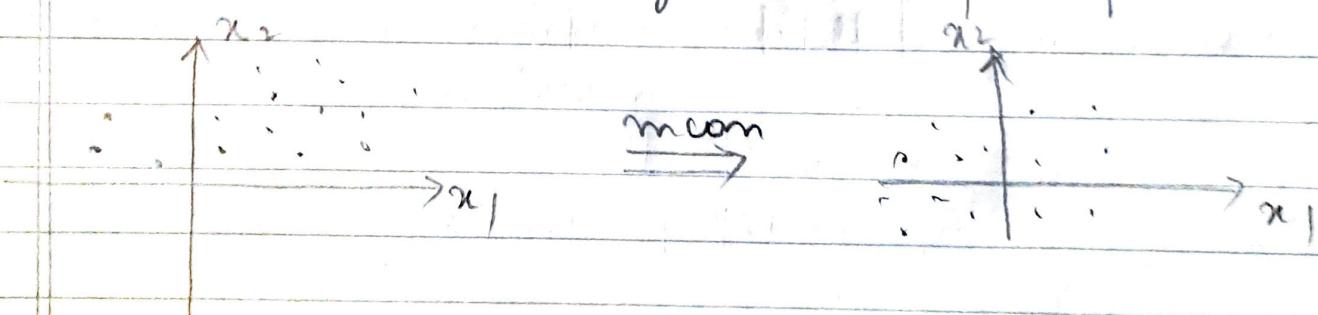
so that it will in the linear region of activation function. Otherwise the BPA tends to drive the free parameters ( $w$ ) of the NN to infinity and thereby slow down the learning process drastically.

#### ⑤ Normalization of inputs:

The training samples must be preprocessed before presenting to the NN.

The preprocessing steps are:

- ① mean subtracting each input from mean



the data was  
not distributed

properly as a result  
after training process.

Now the sample coming from the NN must be

Data are balanced

non zero mean  $\rightarrow$  large  $\epsilon$  values  $\rightarrow$  condition no is large  $\rightarrow$  cost fun. surface will be steep

- ② de correlate the Training samples using PCA  
otherwise NN will be clumsy & slow
- ③ Covariance: it ensures the different weight in the new learned approximately the same speed.
- ④ Learning from hints:  $t$  is the target function which is learning basically means mapping unknown input-output pairs.
- During this process it exploits the information in the t-example. This process can be made faster by giving hints along with input examples. hints can be symmetric, invariance properties etc.
- ⑤ Learning rate:  $\eta_1 > \eta_2 > \eta_3 \dots \eta_n$

Initialization:

A good choice for the initial values of the synaptic weights and thresholds of the NW can be more useful in a successful NW design. Then what is a good choice?

If I chose large values for  $w$ , it is highly likely that the neurons in NW will be driven into saturation. If this happens the local gradient in BPA assume small values, which in turn will cause the learning process to slow down.

If  $w$ 's are assigned with small initial values, BPA may operate on a very flat area around the origin of the cost surface; when origin is a saddle point

Let the bias applied to each neuron in the network be set to zero  $\times v_i = \sum_{j=1}^m w_{ij}y_j$ ,  $\mu_y = E[y_i] = 0 \forall i$  and  $\sigma_y^2 = E[(y_i - \mu_i)^2] = E[y_i^2] = 1$

$$E[y_i y_k] = \begin{cases} 1 & k=i \\ 0 & k \neq i \end{cases}$$

and  $w$ 's are drawn from a uniformly distributed set of no. with '0' mean

$$\mu_w = E[w_{ij}] = 0 \quad \sigma_w^2 = E[(w_{ij} - \mu_w)^2] = E[\text{mean } \times \text{ variance of } v_i]$$

$$\mu_v = E[v_i] = E \sum w_i y_i = 0.$$

$$\sigma_v^2 = E[v_i^2] = \sum E[w_{ij}^2] \quad i=k.$$

$$\sigma_w^2 = m \sigma_w^2 = m \sigma_v^2 \quad \text{where } m \text{ is the no. of synapti}$$

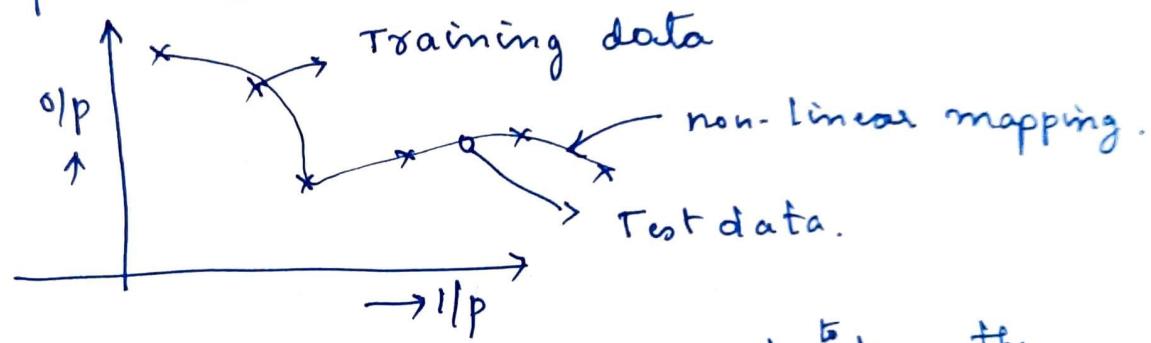
## Generalization & Cross validation.

### i) Generalization:

A NW is said to generalize, when the I/p - O/p mapping computed by NW is correct for test data not used in the training.

To understand the concept of generalization, we consider the NW as slm. which performs non linear input-output mapping.

The learning process is viewed as a curve-fitting problem. The problem of mapping simply boils down to that of interpolation



The training data is used to help the NW learn the non-linear input-output mapping represented by the curve.

When a test data is presented, the NW computes the O/p as a result of interpolation performed by the NW.

When a NN learns too many I/p - O/p samples, the NW may end up memorizing the training set. w/o trying to find the underlying function to be modelled.  
- overfitting / overtraining.

When we trained the N/W loses the ability to generalize.  
This normally happens when more no. of H. Neurons than necessary are used resulting in undesired contributions in the input space due to noise stored in it.

Generalization is influenced by:

- Size of training set and how close it is to the m.v.t. of interest.
- Arch. of NN.
- Physical complexity of the pblm.

Good generalization

$$N = O\left(\frac{W}{\epsilon}\right). \quad O(\cdot) : \text{order of quantity}$$

W: Total no. of free parameters.

$\epsilon$ : error permitted on the test data.

Cross Validation:

- Choosing a suitable topology for a NN, given application is a difficult problem.
- Usually, after a tedious trial & error process, overfitted topology is chosen which is prone to various drawbacks like:
  - ① high demand on computational resources
  - ② increase in training time
  - ③ non convergence of parameters
  - ④ Decrease in the generalisation capability of a NN.
  - ⑤ price for fw implementation ↑
  - ⑥ less efficient.
- If NN topology is too small, then it might not be able to fit the data. To overcome this pblm, we need to find minimal topologies with best performance.
- There are 2 ways to achieve
  - ① Growing approach.
  - ② Pruning Methods.
- Generally latter method starts the training of NN which is expected to be big enough to ensure successful training.

## Pruning Methods

① Deletion

② regularization

### ① Deletion (Hessian based NW pruning)

The basic idea of 2<sup>nd</sup> approach is to use inf<sup>n</sup> of 2<sup>nd</sup> deriv. of error surface in order to make a trade off b/w Err & error performance & network complexity.

To do this we construct a local model of error surface by predicting the effect of perturbation in  $\mathbf{w}$  around

$$E_{av}(\mathbf{w} + \Delta \mathbf{w}) = E_{av}(\mathbf{w}) + \frac{\partial E_{av}}{\partial \mathbf{w}} \Delta \mathbf{w} + \frac{1}{2!} \frac{\partial^2 E_{av}}{\partial \mathbf{w}^2} \Delta \mathbf{w}^2 + H.O.T$$

$\Delta \mathbf{w}$  is perturbation applied to training point  $\mathbf{x}$

Next, we need to identify a set of parameters whose deletion from NW will cause the least increase in the Err. We identify it as follows. by assumption the  $\Delta \mathbf{w}$  after training is such that  $\mathbf{g}^T \Delta \mathbf{w} = 0$  after training. ② error surface is around

local minima is highly quadratic & neglecting H.O.T

$$\Delta E_{av} = E(\mathbf{w} + \Delta \mathbf{w}) - E(\mathbf{w}) = \frac{1}{2} \Delta \mathbf{w}^T H \Delta \mathbf{w}$$

We adopt OBS procedure

Goal is let one of the synaptic weights to zero to minimize the incremental increase in  $E_{av}$  to do so we follow  $\mathbf{u}_i^T \Delta \mathbf{w} + w_i = 0$ .

Minimize the quadratic form w.r.t to incremental change in the weight vector  $\Delta \mathbf{w}$  subject to the constraint that  $\mathbf{u}_i^T \Delta \mathbf{w} + w_i$  is zero & then minimize it w.r.t index i

$$S = \frac{1}{2} \Delta \mathbf{w}^T H \Delta \mathbf{w} + \lambda (\mathbf{u}_i^T \Delta \mathbf{w} + w_i) \quad - A_0 H^{-1} \mathbf{u}_i$$

$$\frac{\partial S}{\partial \Delta \mathbf{w}} = H \Delta \mathbf{w}^T H^T \Delta \mathbf{w} + \lambda \mathbf{u}_i^T - \lambda w_i = 0 \Rightarrow \Delta \mathbf{w}^T H^T H \Delta \mathbf{w} = -\lambda \mathbf{u}_i^T \mathbf{u}_i \quad S = \frac{w_i - \lambda \mathbf{u}_i^T \mathbf{u}_i}{\mathbf{u}_i^T H^T H \mathbf{u}_i}$$

## Conjugate Gradient Method

- it is a second-order optimization method
- let us assume minimization of quadratic function

$$f(x) = \frac{1}{2} x^T A x - b^T x + c.$$

Where  $x$  is a  $W \times 1$  parameter vector

$A$  :  $W \times W$  symmetric positive definite matrix

$b$  :  $W \times 1$  vector

$c$  : constant

minimization of quadratic fun.  $f(x)$  is achieved by

$$f'(x) = 0.$$

$$\Rightarrow A x^* - b = 0$$

$$x^* = A^{-1} b.$$

Thus minimizing  $f \circ f(x)$  and solving linear system of eqns  $Ax^* = b$  are equivalent problems.

Definition: Note for a given matrix  $A$ , we say that a set of non zero vectors  $[s(0) \ s(1) \ \dots \ s(W-1)]$  is  $A$ -conjugate.

If the following condition is satisfied.

$$s^T(n) A s(j) = 0 \quad \forall n, j \text{ s.t } n \neq j.$$

For a given set of A-conjugate vectors  $s(0), s(1), \dots, s(n-1)$ , the corresponding conjugate direction method for unconstrained minimization of the quadratic cost function  $f(x)$  is defined by

$$x^{(n+1)} = x^{(n)} + \eta^{(n)} s^{(n)} \quad n=0, 1, \dots, n-1,$$

where  $x^{(0)}$  is an arbitrary starting vector

$\eta^{(n)}$  : scalar defined by

$$f(x^{(n)} + \eta^{(n)} s^{(n)}) = \min_{\eta} f(x^{(n)} + \eta s^{(n)}).$$

The procedure of choosing  $\eta$  so as to minimize the form  $f(x^{(n)} + \eta s^{(n)})$  for some fixed  $n$  is referred to as line search.

$$f(x^{(n)} + \eta^{(n)} s^{(n)}) = \frac{1}{2} [x^{(n)} + \eta^{(n)} s^{(n)}]^T A [x^{(n)} + \eta^{(n)} s^{(n)}] - b^T [x^{(n)} + \eta^{(n)} s^{(n)}] + c$$

$$f(x^{(n)} + \eta^{(n)} s^{(n)}) = \frac{1}{2} [x^{(n)} A x^{(n)} + \eta^{(n)} s^{(n)} A s^{(n)}] + \frac{1}{2} \eta^{(n)} s^{(n)} T A s^{(n)} + \eta^{(n)} s^{(n)} T A s^{(n)} - b^T x^{(n)} - \eta^{(n)} b^T s^{(n)} + c$$

$$\frac{\partial f(\cdot)}{\partial \eta} = 0 + \frac{1}{2} s^{(n)} A x^{(n)} + \frac{1}{2} x^{(n)} T A s^{(n)} + \eta^{(n)} s^{(n)} T A s^{(n)} - b^T s^{(n)}$$

$$0 = s^{(n)} A x^{(n)} - b^T s^{(n)} + \eta^{(n)} s^{(n)} T A s^{(n)}$$

$$\eta^{(n)} s^{(n)} T A s^{(n)} = -s^{(n)} \{ A x^{(n)} - b \} \quad A = A^T = A^{-1}$$

$$\eta^{(n)} = -\frac{s^{(n)} A \{ x^{(n)} - A^{-1} b \}}{s^{(n)} A s^{(n)}} \quad x^{(n)} - x^* = \varepsilon$$

$$\eta^{(n)} = -\frac{s^{(n)} A \varepsilon^{(n)}}{s^{(n)} A s^{(n)}} \quad \text{where } x^* = A^{-1} b$$

starting from an arbitrary point  $x(0)$  the conjugate direction method is guaranteed to find optimal sol<sup>h</sup>  $x^*$  of  $f(x)=0$  in at most  $H$  iterations.

For the CG method to work we require the availability of set A-conjugate vectors  $s(0), s(1), \dots, s(H-1)$  if it is determined in a sequential manner at ~~iterating~~ step.

Define residual  $r(m) = b - Ax(m)$

$$\text{L.C of } r(m) \times s(m-1) \Rightarrow s(m) = r(m) + \beta(m) s(m-1) \quad m=1, \dots, H-1$$

Where  $\beta(k)$  is chosen s.t. A conjugacy  $\Rightarrow$  process.

$$\Rightarrow s^T(n-1) A s(n) = s^T(n-1) r(n) + \beta(n) s^T(n-1) A s(n-1).$$

$$0 = s^T(n-1) A s(n) + \beta(n) s^T(n-1) A s(n-1)$$

$$\Rightarrow \beta(n) = -\frac{s^T(n-1) A r(n)}{s^T(n-1) s(n-1)}.$$

~~conjugate~~ To avoid the dependency on A to generate s<sub>n</sub> vectors. we consider

Polak-Ribiere formula

$$\beta(n) = \frac{s^T(n)[r(n) - r(n-1)]}{s^T(n-1) r(n-1)}$$

Fletcher-Reeves formula

$$\beta(n) = \frac{r^T(n) r(n)}{s^T(n-1) r(n-1)}.$$