

Day - 25, Oct 27, 2024.

- Non-linearity in Neural Networks, Purposes of Nonlinearities
- Activation functions → ReLU, ELU, leaky ReLU, Sigmoid, tanh
- Optimization Part : SGD, Adam, AdaGrad, AdaDelta
- Backpropagation Via Flow Graphs and Chain Rule Rule
- Flow Graphs, Recursive Application
- Regularized linear Regression

No Matter How Deep Neural Networks without Activation

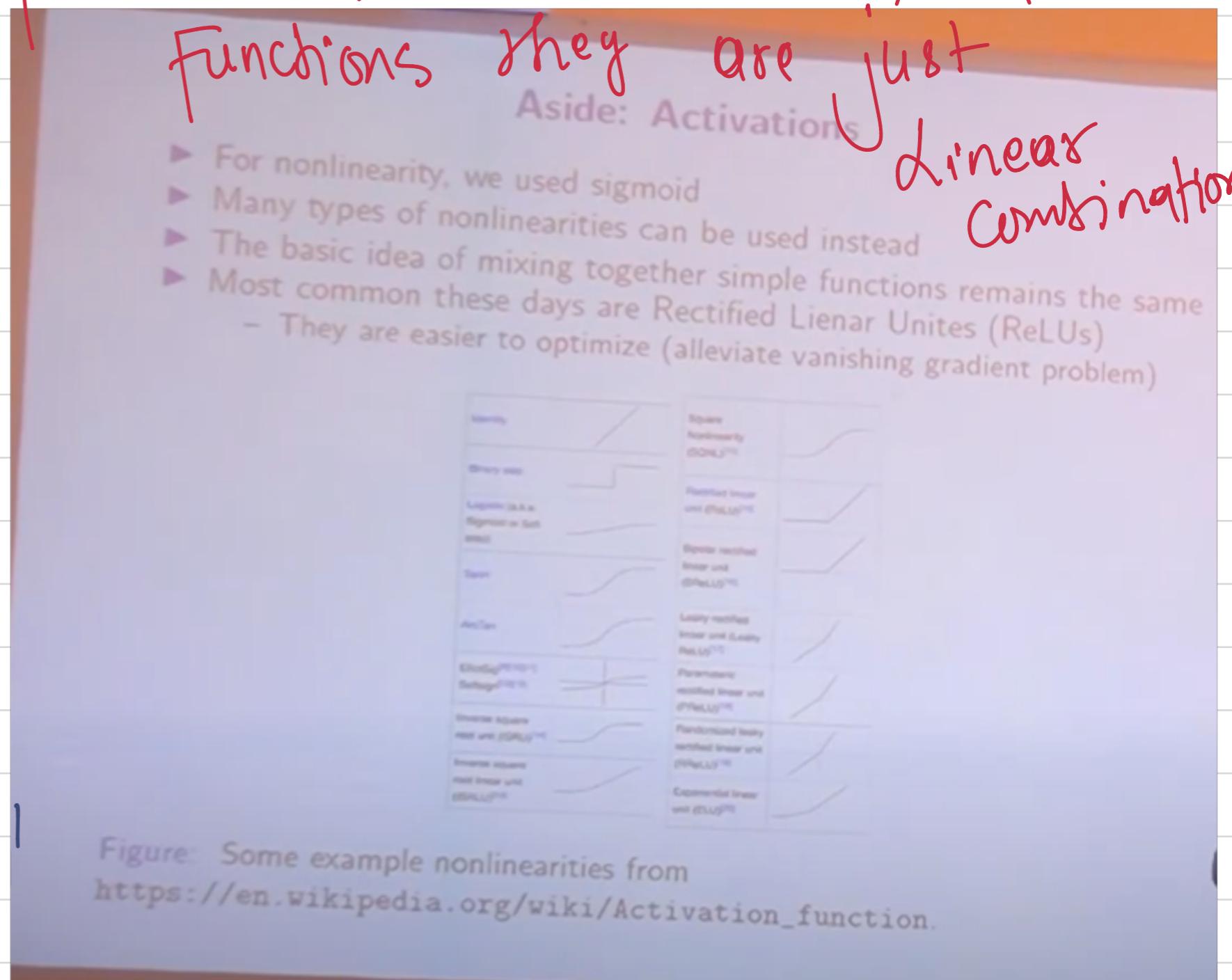
① Sigmoid: Squeezes the input 2 weights range $[0, 1]$

② tanh: $[-1, 1]$

③ ReLu: $[0, \infty)$ \rightarrow Simplicity, effectiveness, mitigate

gradient Vanishing problem

④ Leaky ReLu: Allow Smaller negative values $[-ve, +ve]$.



⑤ ELU (Exponential Linear Unit):

ReLU + Leaky ReLU → Smooth transitions for negative inputs.

⑥ SoftMax function:
used in the output layer for Multi-class problem

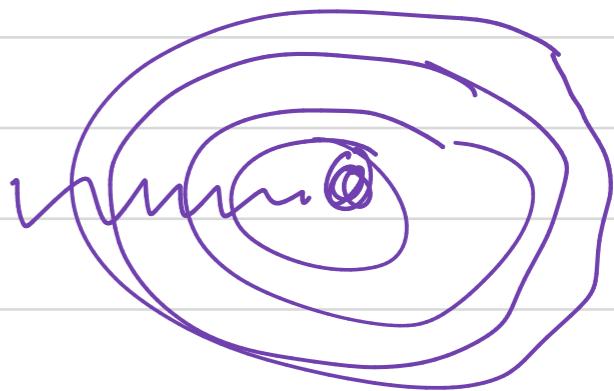
Why Non-linearity?

→ Because it introduces Complexity into network, allowing it to capture intricate features, make more accurate predictions, enable the network to learn hierarchical representations of data.

Optimizers: Reduce the loss of the loss function, find the minima values, optimal parameters to minimize loss functions

SGD → traditional iterative, takes step in the direction of steepest descent to reach the minimum of the loss function

Adam → adapts learning rate for each parameter, faster convergence & better performance, widely used



AdaGrad → adjusts the learning rate for each parameter based on the history of gradients, allowing for more aggressive updates for less frequent parameters.

Adadelta → A variant of AdaGrad better version overcome Shortcomings such as diminishing learning rate.

Speed Measurement in Common Scenario.
Adam > RMSprop > Adadelta > AdaGrad > SGD
 \downarrow
adjusts learning rate
parameter

Back propagation via flow charts

- Back propagation calculates the gradient of the loss function w.r.t. model's parameters and then updating the parameters to minimize the loss. It will be nightmare if we compute each one of them using only derivatives. So, we can use composition function get solved by chain Rule.
-

Chain Rule allows us to differentiate Composite functions.

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial w}$$

Chain Rule Revisited

- Chain rule for vector-valued transformations $x \xrightarrow{g} h \xrightarrow{f} y$,
- $$D(f \circ g)(x) = Df(h) Dg(x)$$

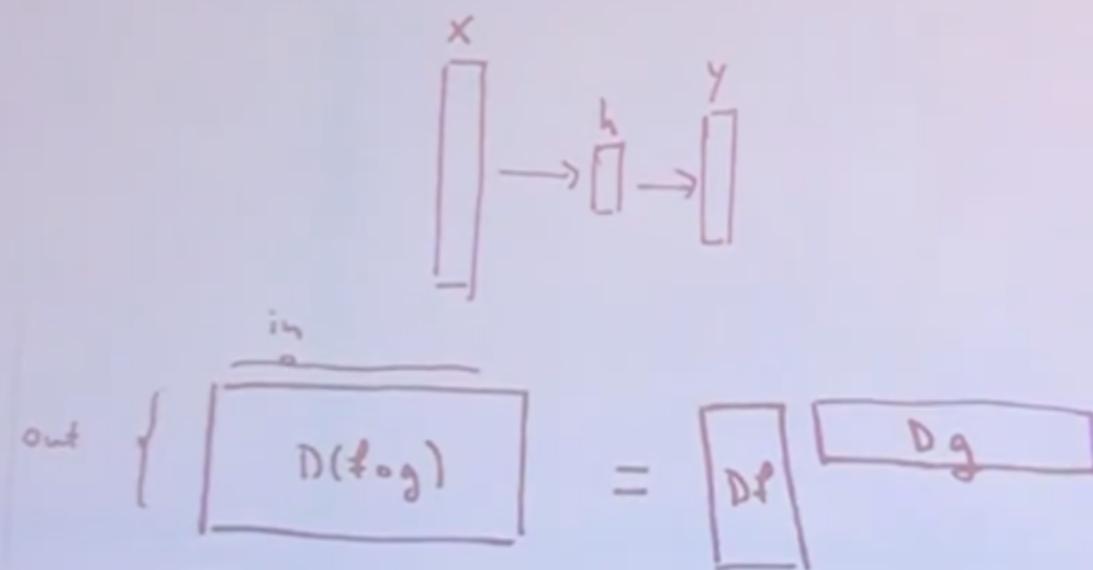


Figure: The first row shows the sequence of vector-valued transformations. The ij^{th} cell of the derivative matrices tells you how the i^{th} output changes when you perturb the j^{th} input slightly.

Optimization: Stochastic Gradient Descent

- Optimize by moving in a direction expected to decrease loss
- SGD: Estimated direction of gradient, after looking at small batch of data
- These days, other methods used too – Adam, AdaGrad, AdaDelta, ...

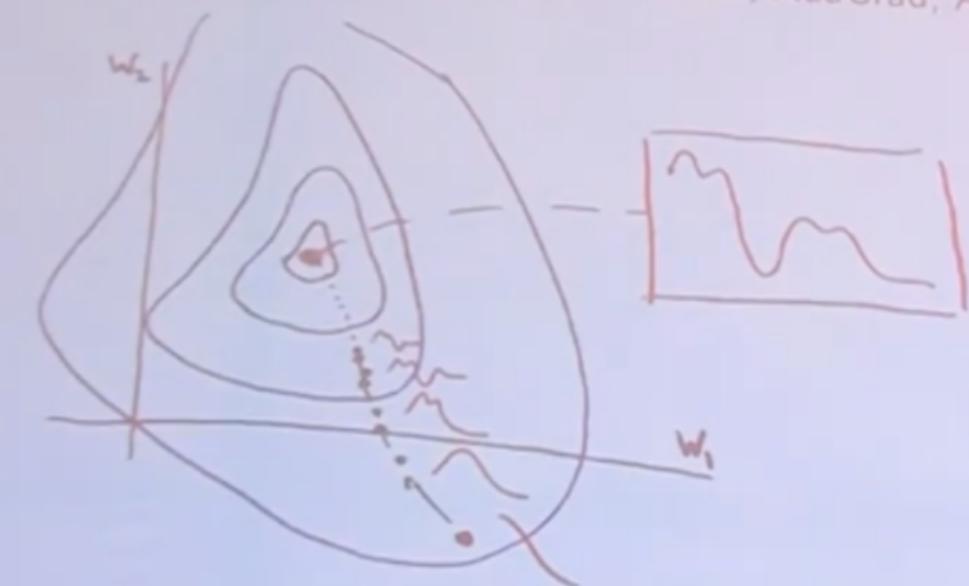


Figure: Each optimization step modifies model parameters to make the fitted function closer to the true (local) minimizer of the loss.

See how chain Rule helps to solve gradient of loss functions (Composite functions) into simpler form - SGD iteratively steps into steepest direction to minimize the loss function.

Backpropagation via Flow Graphs

- ▶ For each iteration of SGD, we need to estimate the gradient of the loss with respect to the parameters
- ▶ When you perturb W_k a little, it changes all the downstream h_k, h_{k+1}, \dots, h_K up to the prediction, which uses that final representation
- ▶ There is a lot of composition going on here... courage alone won't help you take these derivatives
- ▶ Chain rule saves the day again!

Flow Graphs

- ▶ A flow graph shows how different vectors are computed from one another
 - Nodes are vectors
 - Edges are specific operations on input, which when combined with other incoming edges gives the output
- ▶ Crucial: Operations along each edge should be easily differentiable with respect to input

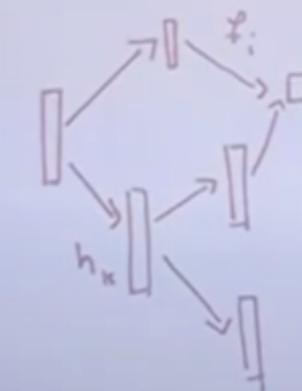


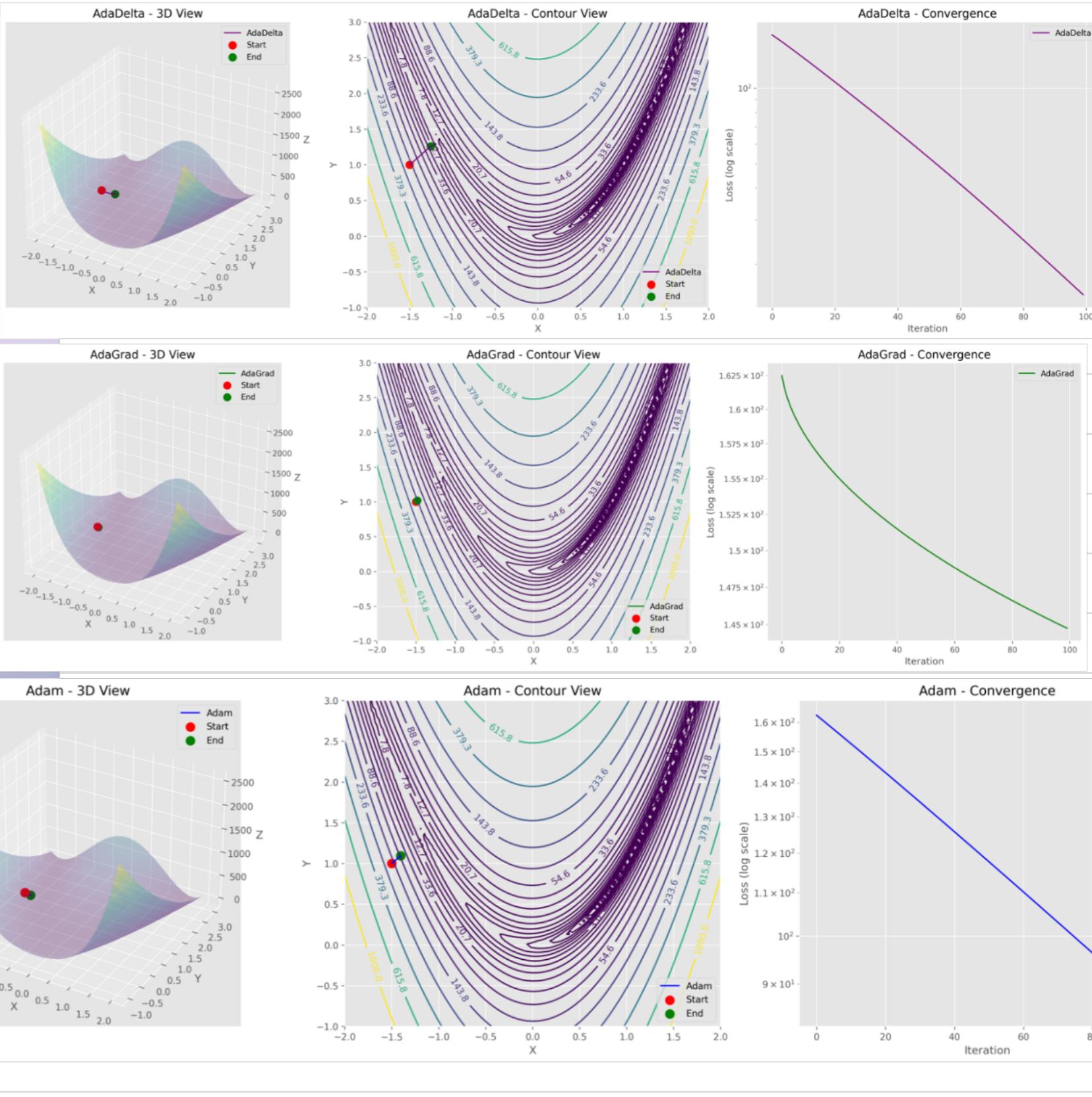
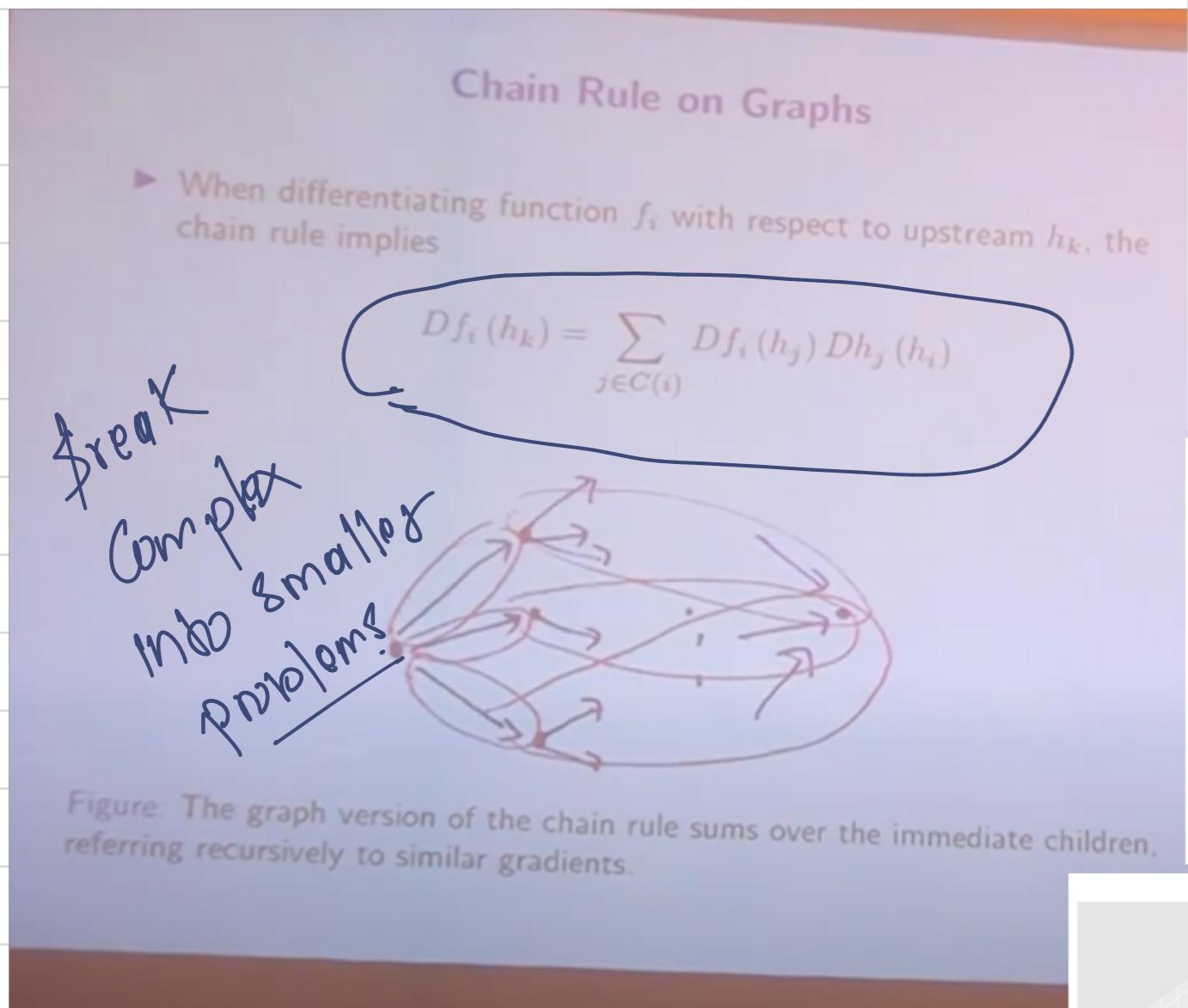
Figure: The main idea of a flow graph.

$$DF_i(h_k) = \sum (DF_i(h_j) + Dh_j(h_k)) \text{ for } j \in C(i)$$

$C(i)$ → represents the set of immediate children node i .

$h_k \rightarrow$ sum of products , F_i - short that the derivative of f_i .

So, the
solution is to
break child
nodes
first.



Chain Rule.

$$\frac{dL}{dw} = \frac{dL}{da} \cdot \frac{da}{dw}$$

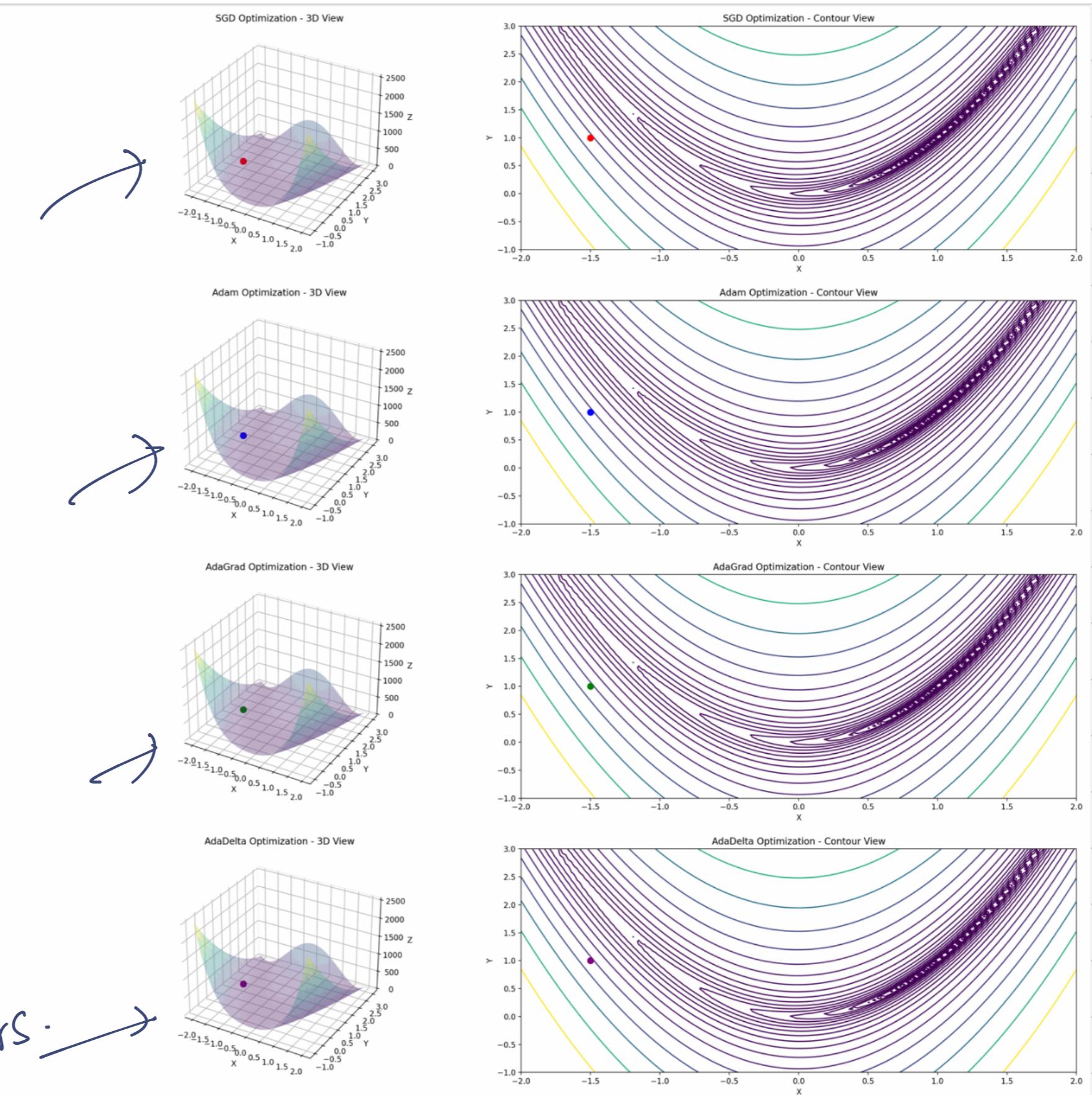
→ SGD → AdaGrad

Adam → AdaDelta

→ Back propagation via
→ Flow Graph

→ ReLU, Sigmoid, Softmax,
tanh.

figure showing optimizers:



From Textbook - (Knowledge Representation) PL: Inference Using Resolution

Resolution \rightarrow Powerful inference used in PL to prove or determine the validity of logical arguments.

Resolution Rule

Given two clauses:

C1: A ∨ B

C2: $\neg B \vee C$

C3: A ∨ C (infer a new clause)

So if we have two clauses, one containing a literal and the other containing negation we can resolve the complementary literals and combining the remaining literals.

KB: (Knowledge Base).

- $P \rightarrow Q$
- $Q + R$
- $\neg Q R$

① Convert into CNF.

- $\neg P \vee Q$
- $\neg Q \vee R$
- $\neg R$

② Apply Resolution:

- Resolve $(\neg P \vee Q)$ and $(\neg Q \vee R)$: $\neg P \vee R$
- Resolve $(\neg P \vee R)$ and $\neg R$: $\neg P$

Example:

$$A \rightarrow B$$

$$B \rightarrow C$$

A is true.

① Convert into CNF:

$$A \rightarrow B \text{ becomes } \neg A \vee B$$

$$B \rightarrow C \text{ becomes } \neg B \vee C$$

A remains A

Assume $\neg C$

- $\neg A \vee B$, $\neg B \vee C$, A and $\neg C$

Apply Resolution.

Resolve A with $\neg A \vee B$ to get B

ii B with $\neg B \vee C$ to get C.

ii C with $\neg C$ to get Contradiction

i) Convert KB and Negation of theorem in CNF

• CNF is conjunction of clauses (\wedge)

• Each clause is a disjunction of literals (\vee)

KB: $P \rightarrow Q, Q \rightarrow R, \neg R \rightarrow P$ prove $\neg P$.

① Convert Implications to Disjunctions:

• $P \rightarrow Q$ is equivalent to $\neg P \vee Q$

$P \rightarrow$ Cat is animal.

$\neg(P \vee Q) \rightarrow$ True.

$Q \rightarrow$ Dog is animal

Neither Cat nor Dog is bird.

So, • $Q \rightarrow R$ is $\neg Q \vee R$

So, Negate theorem $\neg P$ is P

$\neg(\neg P) \rightarrow P$.

C1: $\neg P \vee Q$

C2: $\neg Q \vee R$.

C3: P

Resolution Steps:

$$P \rightarrow Q \quad Q \rightarrow R, \quad TR.$$

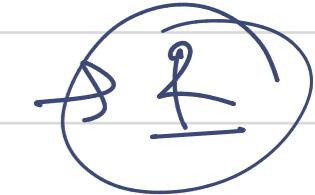
(1) Resolve $P \rightarrow \neg P \vee Q$

$$\cancel{P(\neg P \vee Q)}$$



(2) Resolve $Q \rightarrow \neg Q \vee R$

$$\cancel{Q(\neg Q \vee R)}$$



(3) Resolve R with $\neg R$. \rightarrow empty clause()

$$R(\neg R)$$

(4) $Q \rightarrow \neg P \vee Q$



(5) $R \rightarrow \neg Q \vee R$.

