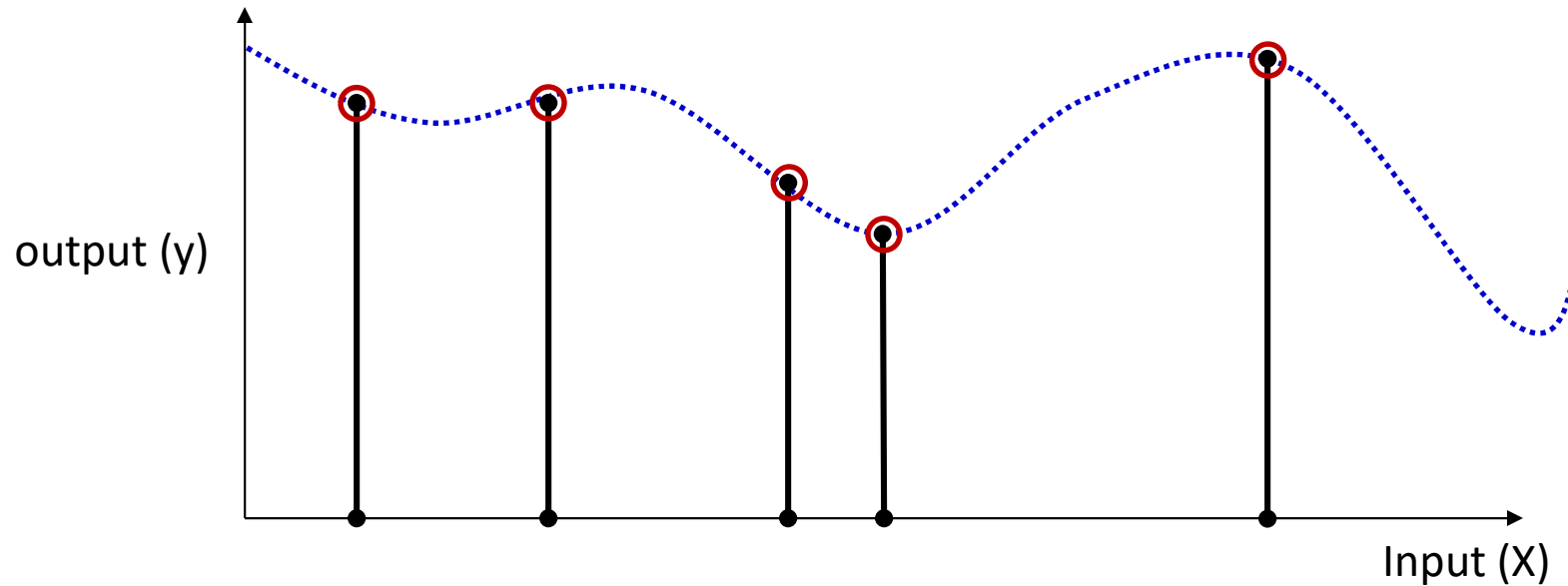
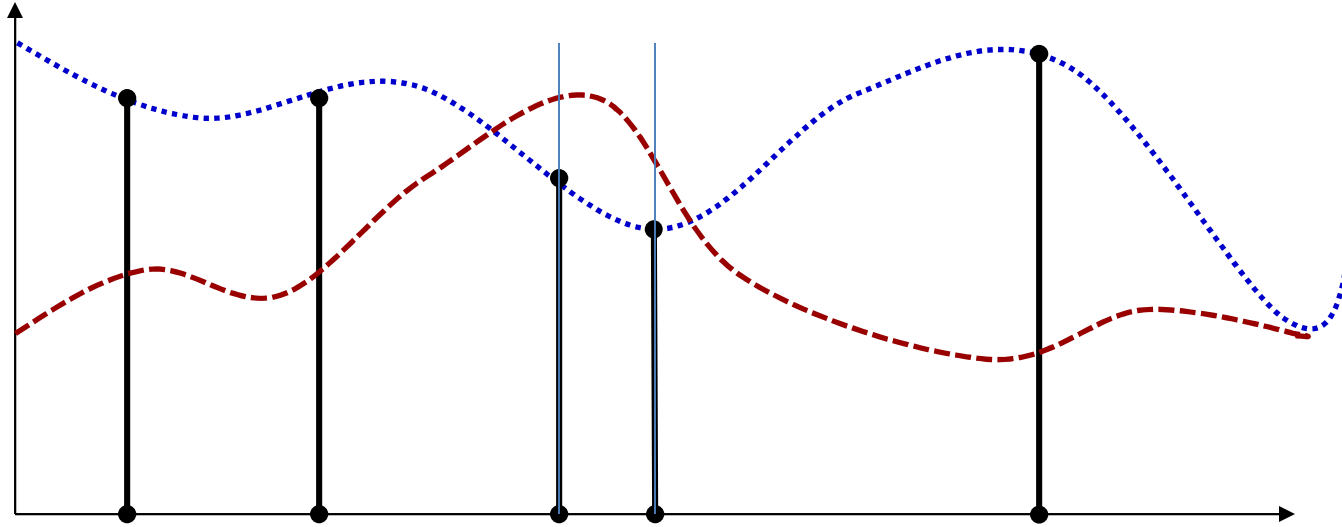


# The training formulation



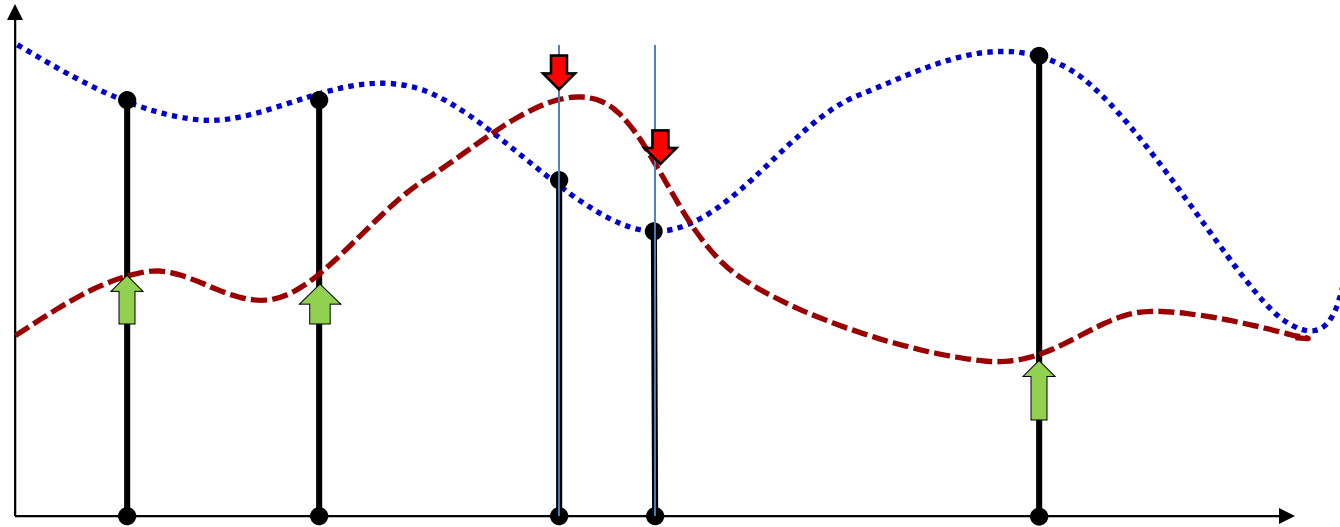
- Given input output pairs at a number of locations, estimate the entire function

# Gradient descent



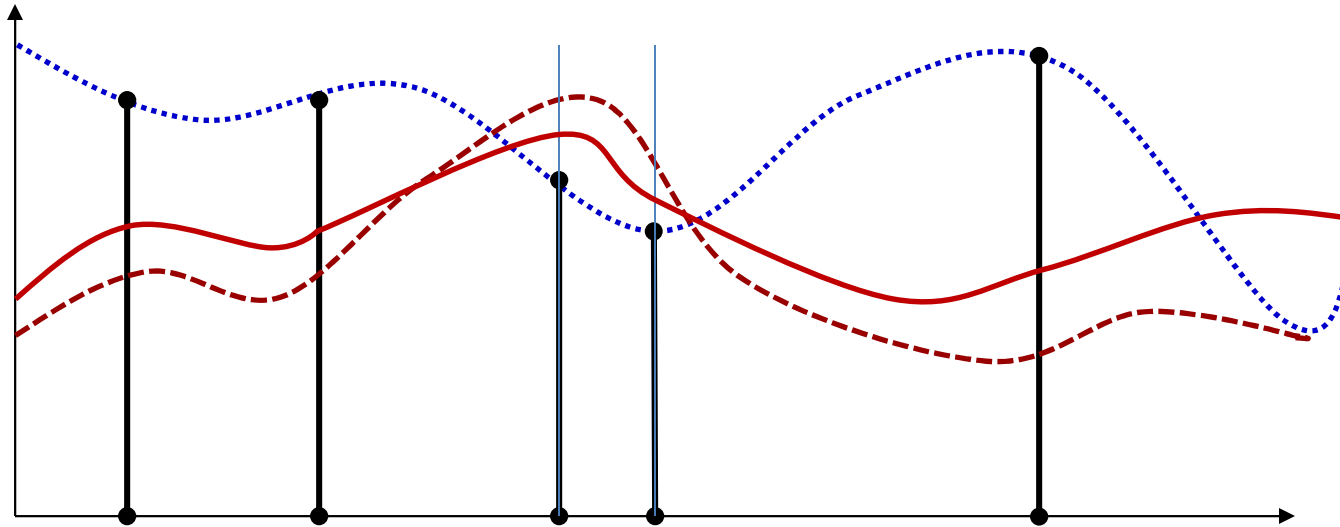
- Start with an initial function

# Gradient descent



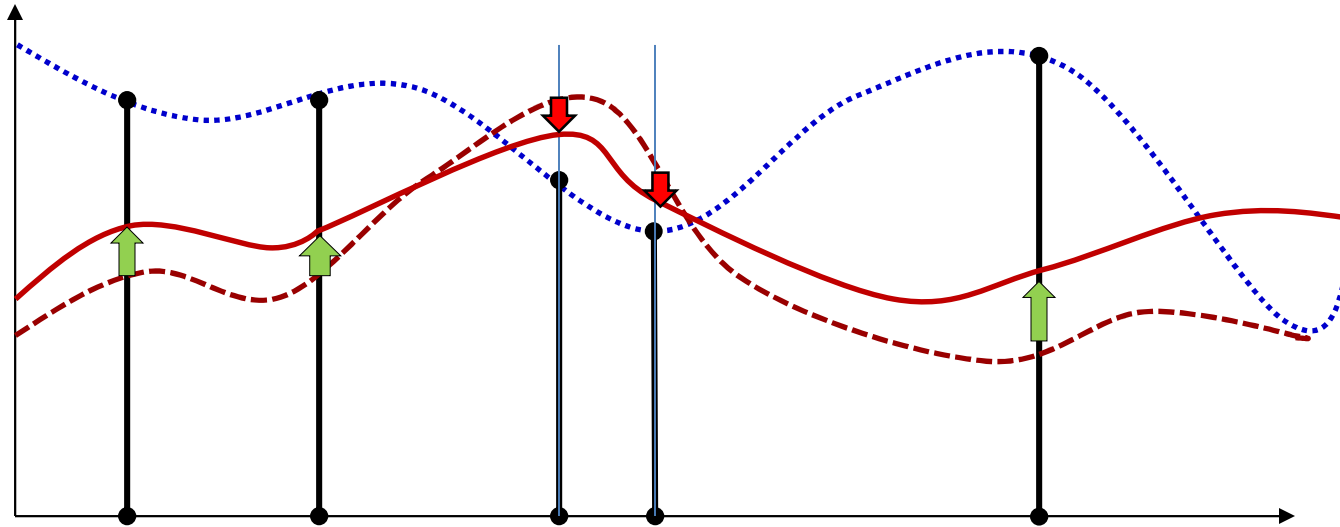
- Start with an initial function
- Adjust its value at *all* points to make the outputs closer to the required value
  - Gradient descent adjusts parameters to adjust the function value at *all* points
  - Repeat this iteratively until we get arbitrarily close to the target function at the training points

# Gradient descent



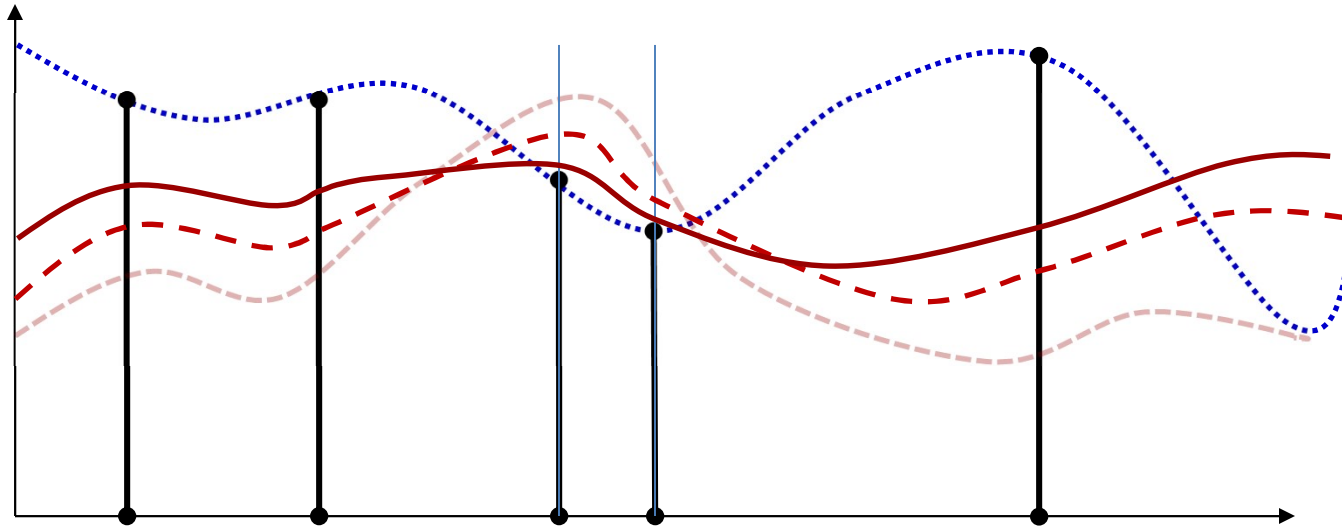
- Start with an initial function
- Adjust its value at *all* points to make the outputs closer to the required value
  - Gradient descent adjusts parameters to adjust the function value at *all* points
  - Repeat this iteratively until we get arbitrarily close to the target function at the training points

# Gradient descent



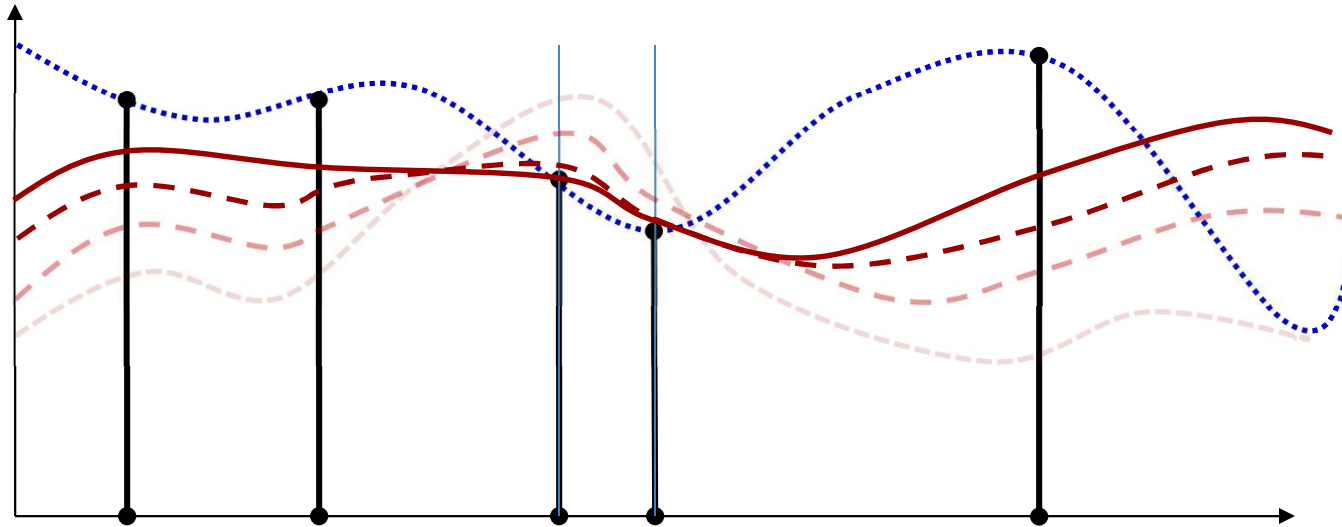
- Start with an initial function
- Adjust its value at *all* points to make the outputs closer to the required value
  - Gradient descent adjusts parameters to adjust the function value at *all* points
  - Repeat this iteratively until we get arbitrarily close to the target function at the training points

# Gradient descent



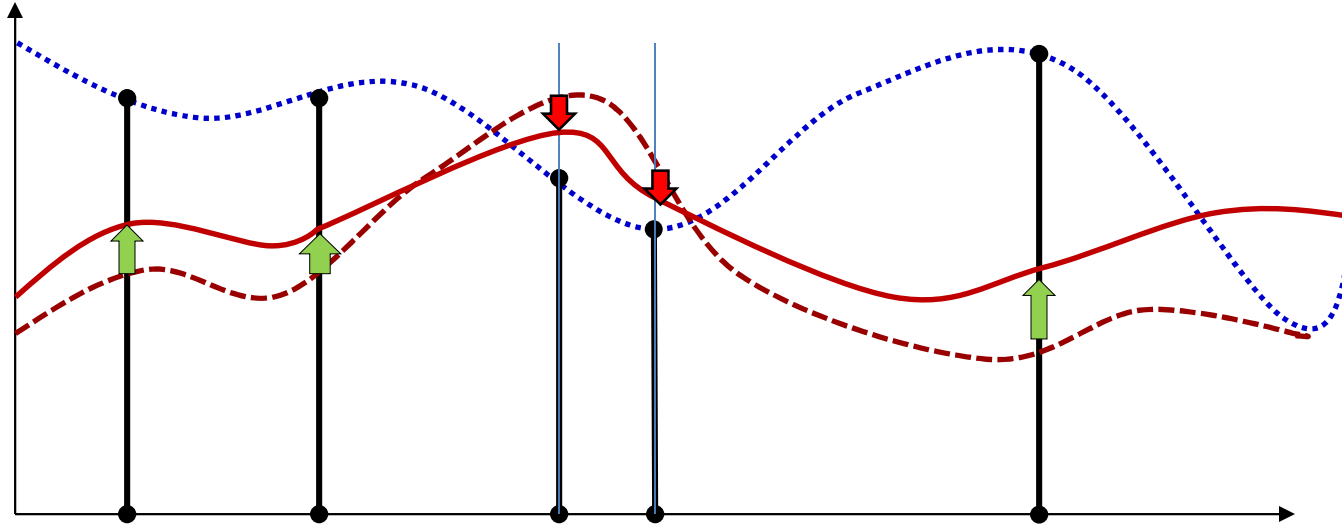
- Start with an initial function
- Adjust its value at *all* points to make the outputs closer to the required value
  - Gradient descent adjusts parameters to adjust the function value at *all* points
  - Repeat this iteratively until we get arbitrarily close to the target function at the training points

# Gradient descent



- Start with an initial function
- Adjust its value at *all* points to make the outputs closer to the required value
  - Gradient descent adjusts parameters to adjust the function value at *all* points
  - Repeat this iteratively until we get arbitrarily close to the target function at the training points

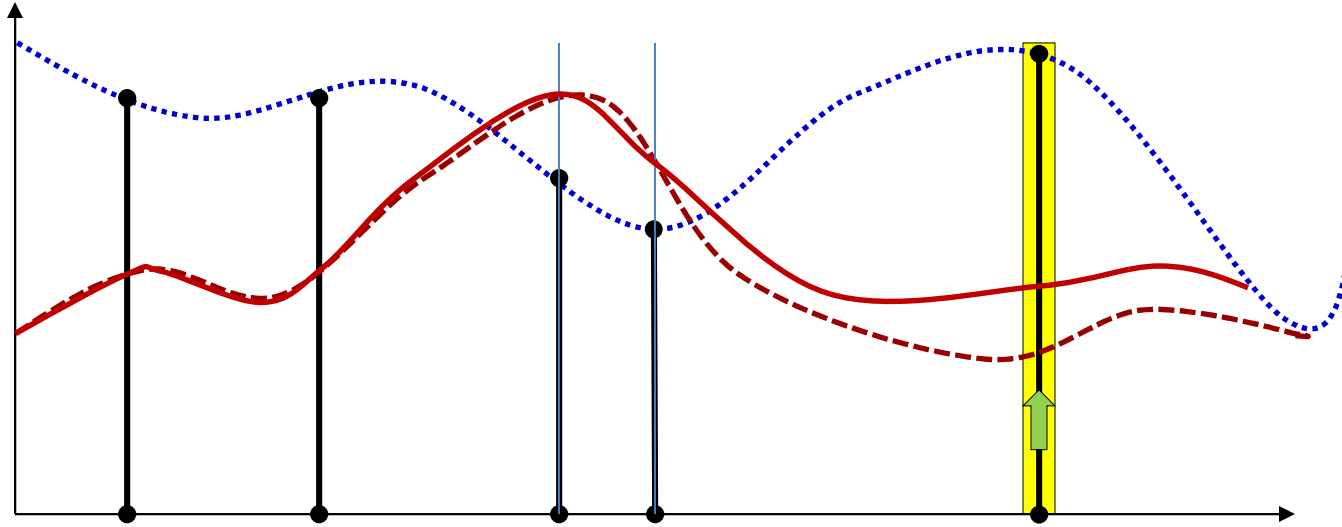
# Effect of number of samples



- Problem with conventional gradient descent: we try to simultaneously adjust the function at *all* training points
  - We must process *all* training points before making a single adjustment
  - **“Batch”** update

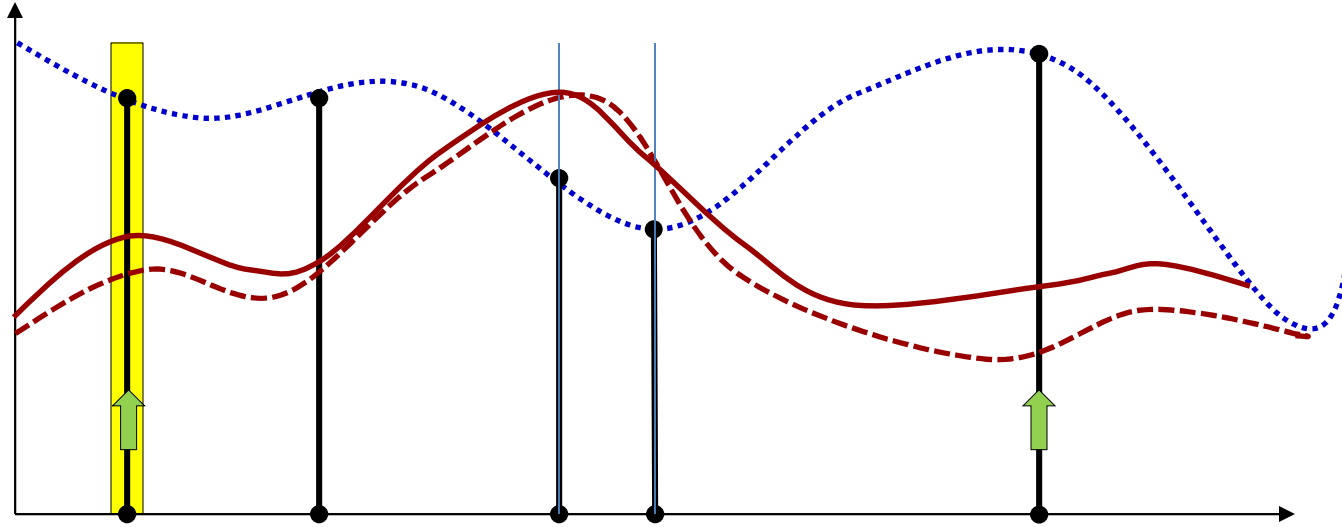


# Alternative: Incremental update



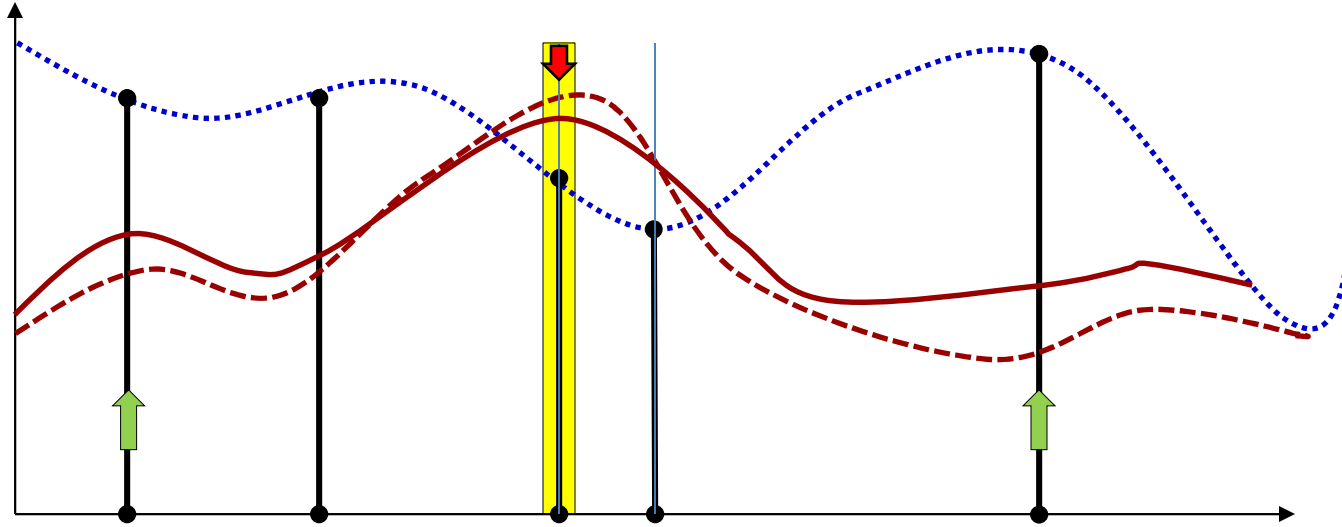
- Alternative: adjust the function at one training point at a time
  - Keep adjustments small

# Alternative: Incremental update



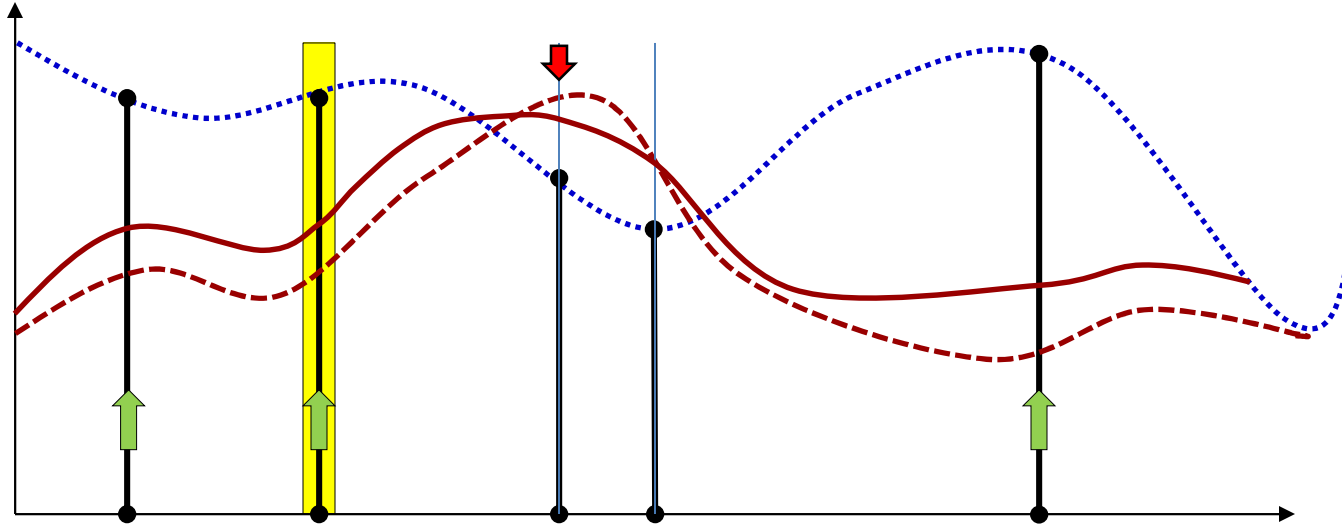
- Alternative: adjust the function at one training point at a time
  - Keep adjustments small

# Alternative: Incremental update



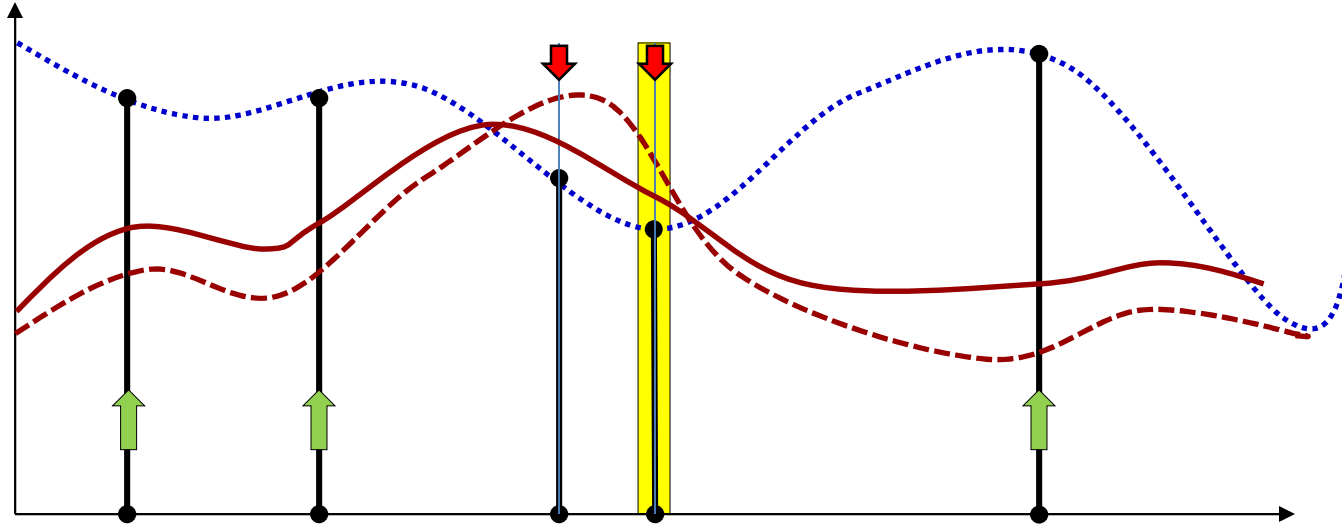
- Alternative: adjust the function at one training point at a time
  - Keep adjustments small

# Alternative: Incremental update



- Alternative: adjust the function at one training point at a time
  - Keep adjustments small

# Alternative: Incremental update



- Alternative: adjust the function at one training point at a time
  - Keep adjustments small
  - Eventually, when we have processed all the training points, we will have adjusted the entire function
    - With *greater* overall adjustment than we would if we made a single “Batch” update

# Incremental Update

- Given  $(X_1, d_1), (X_2, d_2), \dots, (X_T, d_T)$
- Initialize all weights  $W_1, W_2, \dots, W_K$
- Do:
  - For all  $t = 1:T$ 
    - For every layer  $k$ :
      - Compute  $\nabla_{W_k} \text{Div}(\mathbf{Y}_t, \mathbf{d}_t)$
      - Update
$$W_k = W_k - \eta \nabla_{W_k} \text{Div}(\mathbf{Y}_t, \mathbf{d}_t)^T$$
- Until *Loss* has converged

# Incremental Updates

- The iterations can make multiple passes over the data
- A single pass through the entire training data is called an “epoch”
  - An epoch over a training set with  $T$  samples results in  $T$  updates of parameters

# Incremental Update

- Given  $(X_1, d_1), (X_2, d_2), \dots, (X_T, d_T)$
- Initialize all weights  $W_1, W_2, \dots, W_K$

- Do:

Over multiple epochs

One epoch



- For all  $t = 1:T$

- For every layer  $k$ :

- Compute  $\nabla_{W_k} \text{Div}(Y_t, d_t)$

- Update

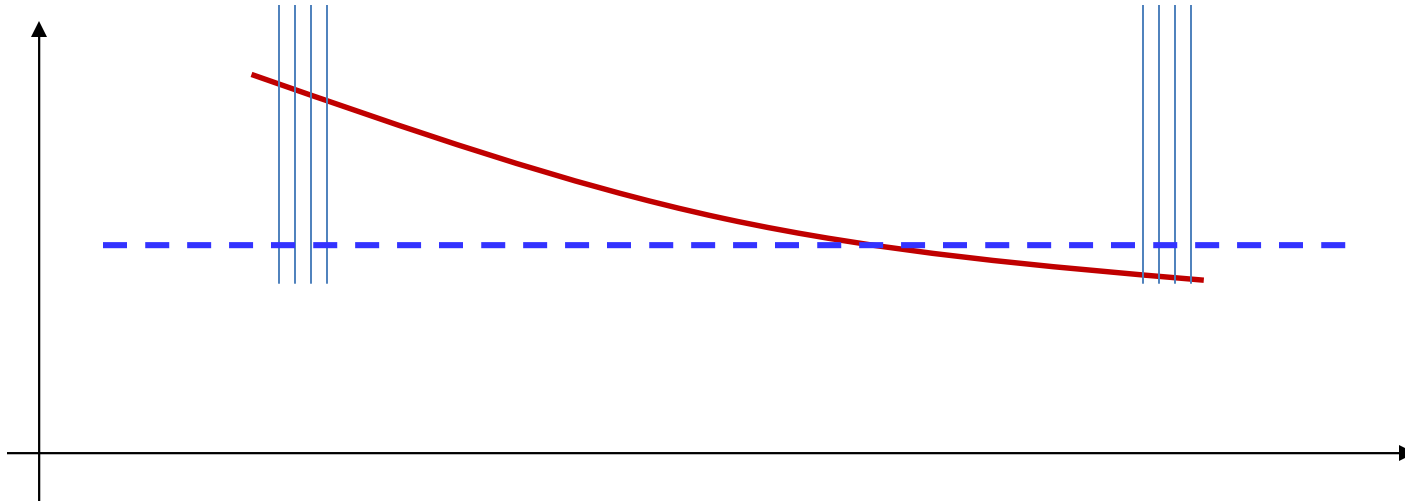
$$W_k = W_k - \eta \nabla_{W_k} \text{Div}(Y_t, d_t)^T$$

One update

- Until *Loss* has converged

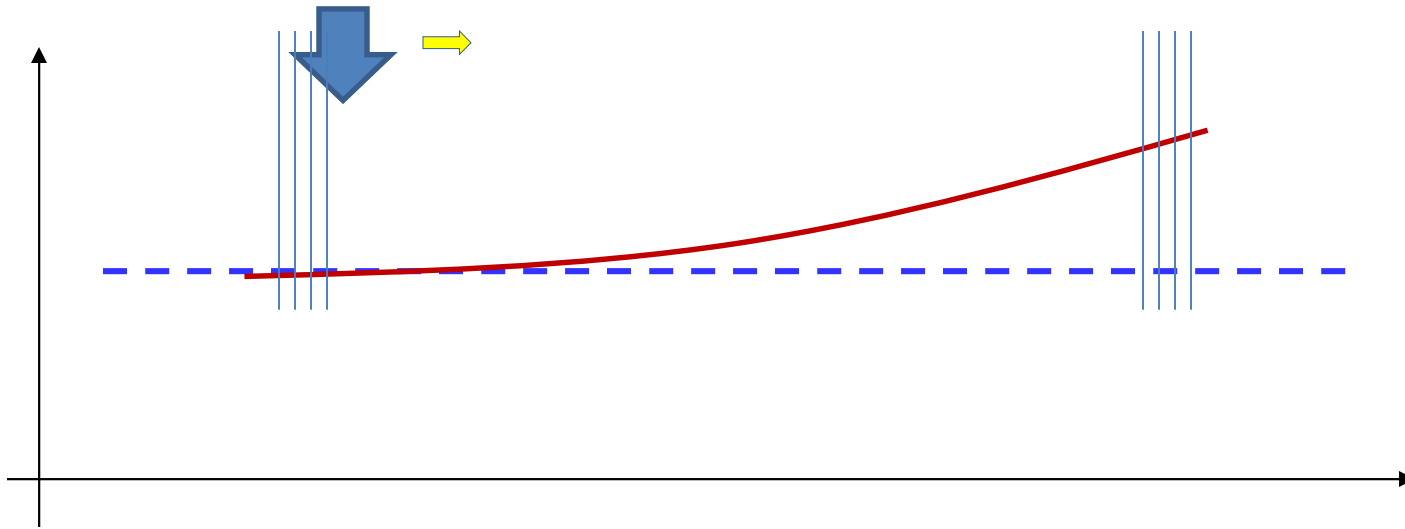


# Caveats: order of presentation



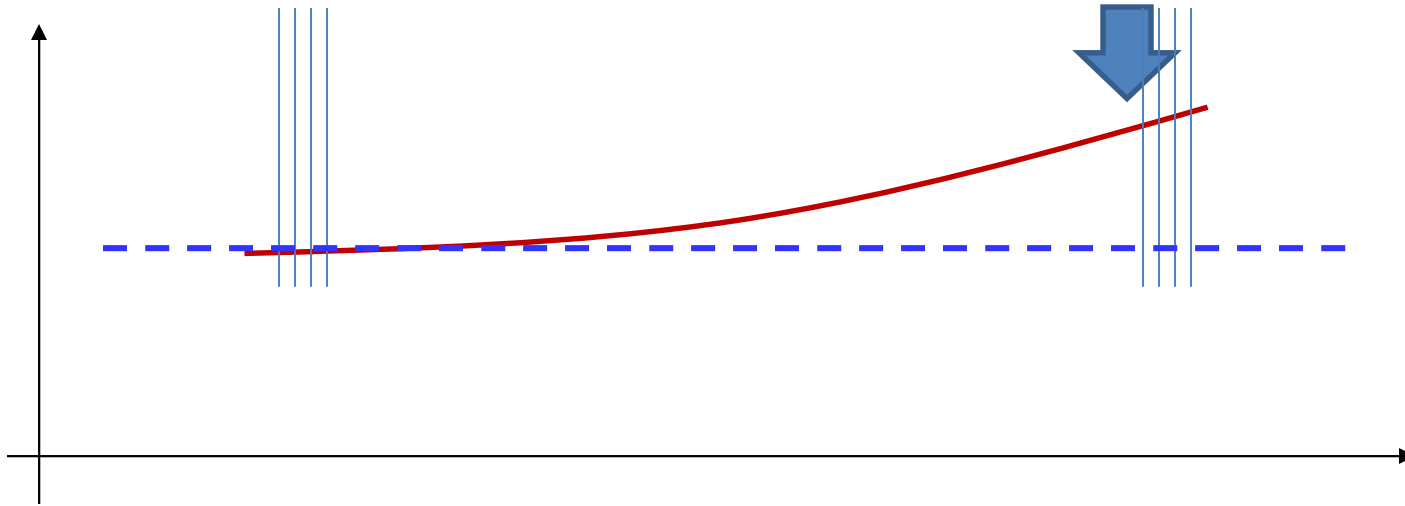
- If we loop through the samples in the same order, we may get *cyclic* behavior

# Caveats: order of presentation



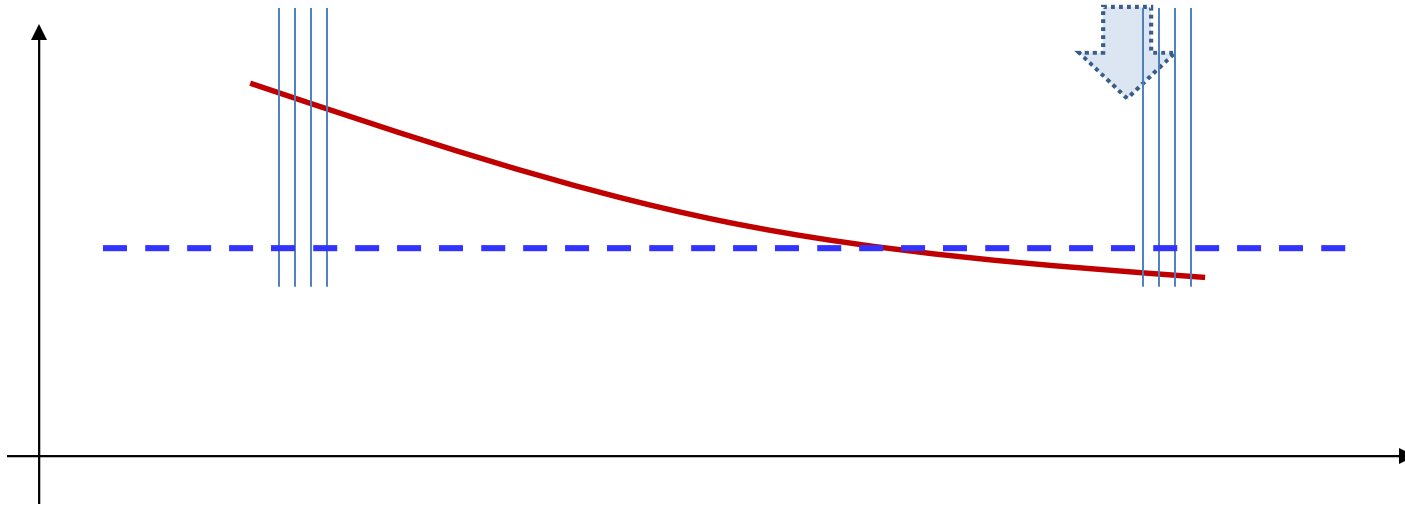
- If we loop through the samples in the same order, we may get *cyclic* behavior

# Caveats: order of presentation



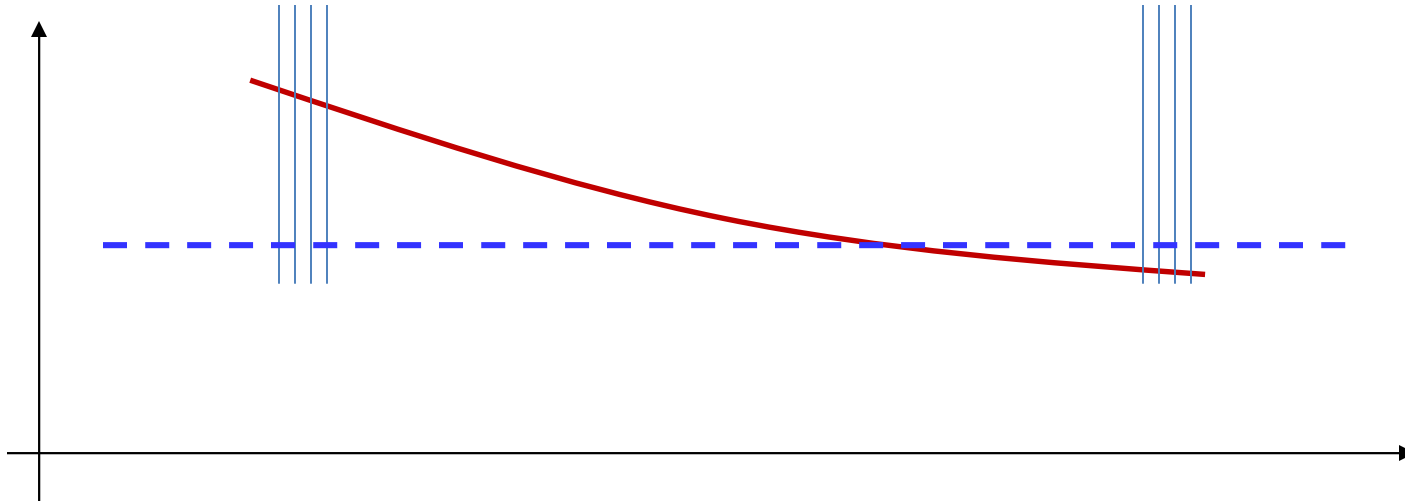
- If we loop through the samples in the same order, we may get *cyclic* behavior

# Caveats: order of presentation



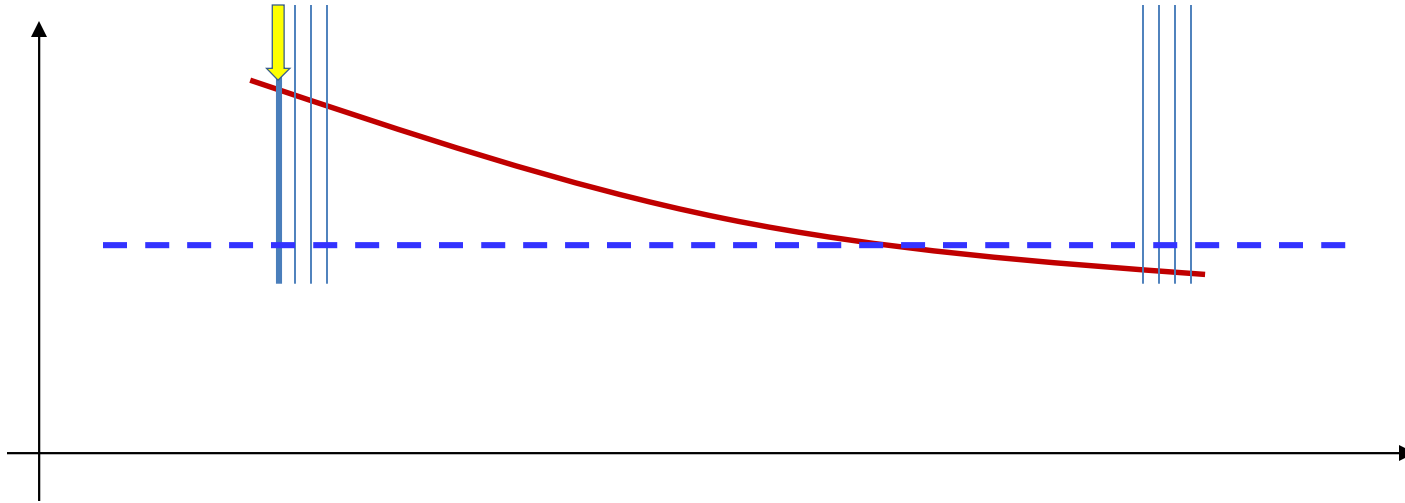
- If we loop through the samples in the same order, we may get *cyclic* behavior

# Caveats: order of presentation



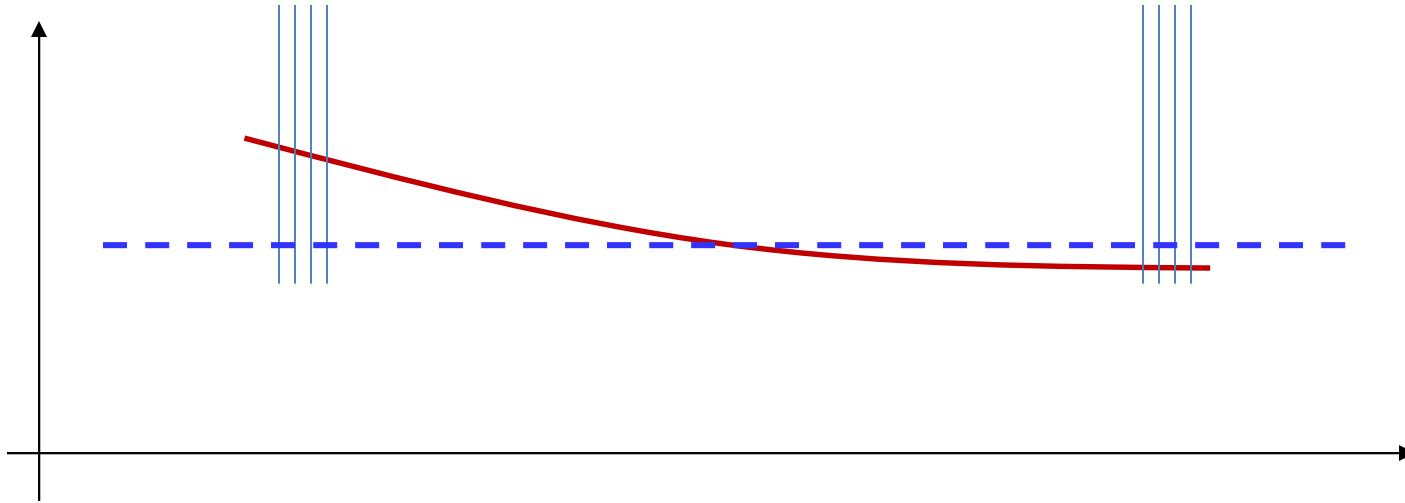
- If we loop through the samples in the same order, we may get *cyclic* behavior
- We must go through them *randomly* to get more convergent behavior

# Caveats: order of presentation



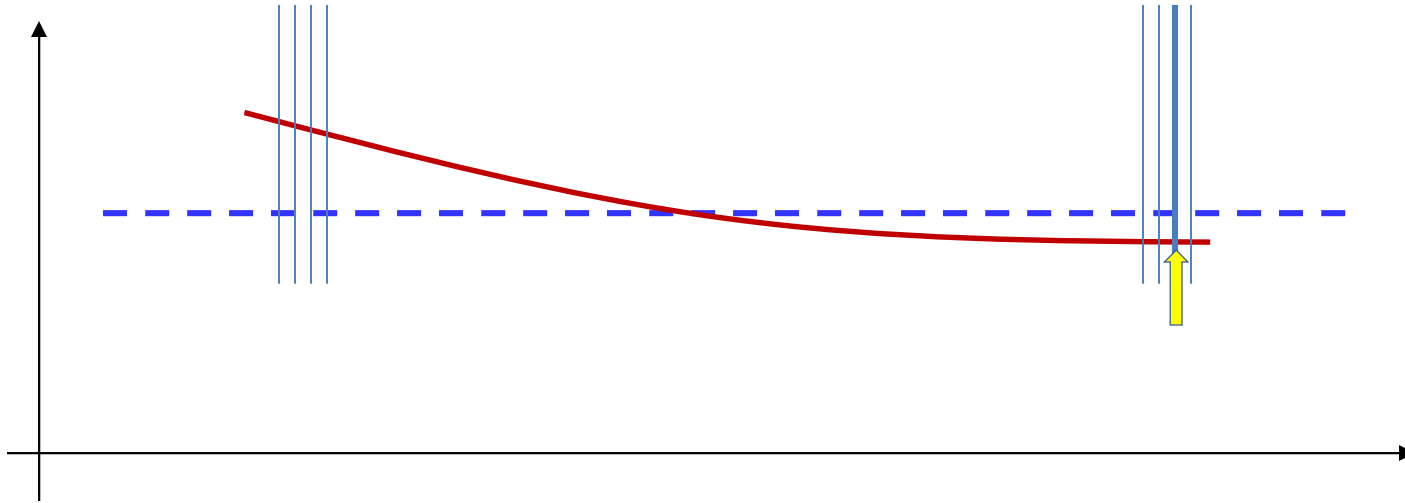
- If we loop through the samples in the same order, we may get *cyclic* behavior
- We must go through them *randomly* to get more convergent behavior

# Caveats: order of presentation



- If we loop through the samples in the same order, we may get *cyclic* behavior
- We must go through them *randomly* to get more convergent behavior

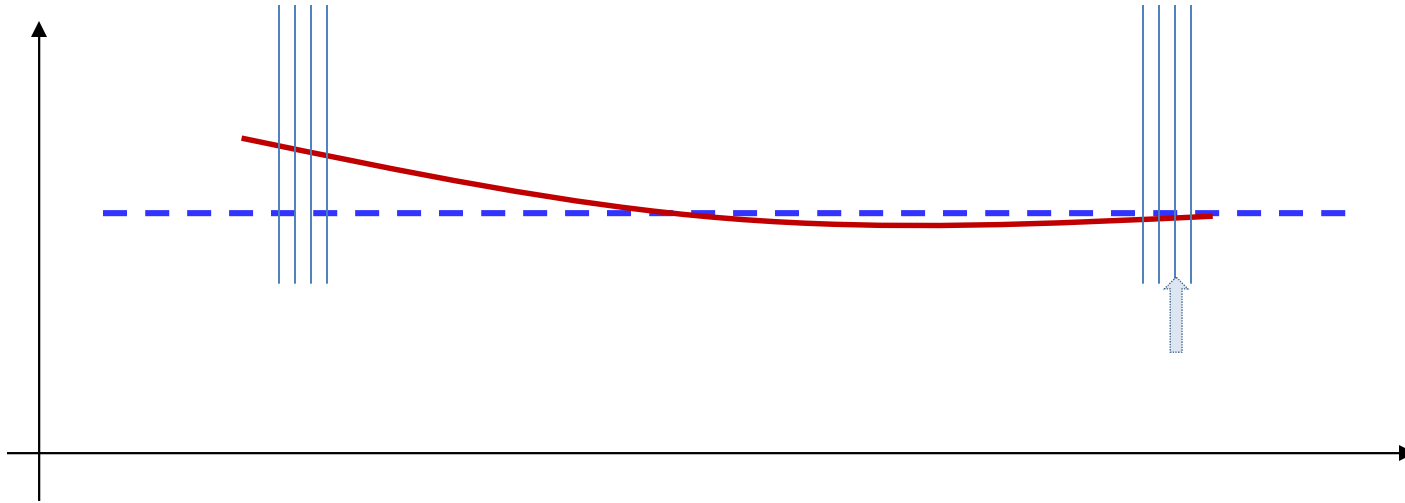
# Caveats: order of presentation



- If we loop through the samples in the same order, we may get *cyclic* behavior
- We must go through them *randomly* to get more convergent behavior



# Caveats: order of presentation



- If we loop through the samples in the same order, we may get *cyclic* behavior
- We must go through them *randomly* to get more convergent behavior

# Incremental Update: Stochastic Gradient Descent

- Given  $(X_1, d_1), (X_2, d_2), \dots, (X_T, d_T)$
- Initialize all weights  $W_1, W_2, \dots, W_K$
- Do:
  - Randomly permute  $(X_1, d_1), (X_2, d_2), \dots, (X_T, d_T)$
  - For all  $t = 1:T$ 
    - For every layer  $k$ :
      - Compute  $\nabla_{W_k} \text{Div}(\mathbf{Y}_t, \mathbf{d}_t)$
      - Update
$$W_k = W_k - \eta \nabla_{W_k} \text{Div}(\mathbf{Y}_t, \mathbf{d}_t)^T$$
- Until *Loss* has converged