

ECE 598NSG/498NSU

Deep Learning in Hardware

Fall 2020

Training DNNs – SGD, Backprop and its
Variants

Naresh Shanbhag

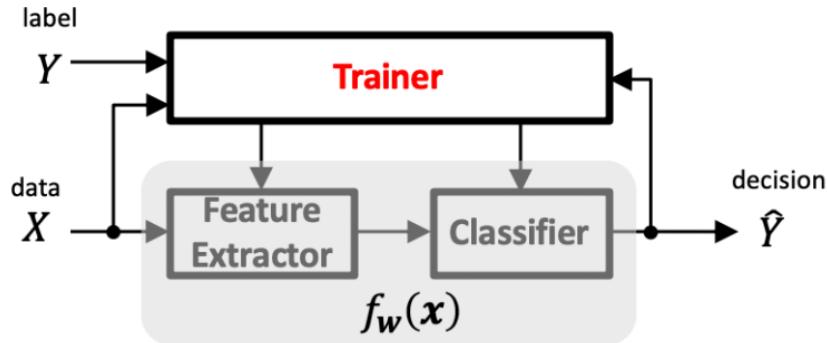
Department of Electrical and Computer Engineering
University of Illinois at Urbana-Champaign

<http://shanbhag.ece.uiuc.edu>

Today

- Overview of DNN training – the big picture
- The Stochastic Gradient Descent (SGD) training algorithm
- The back-prop algorithm

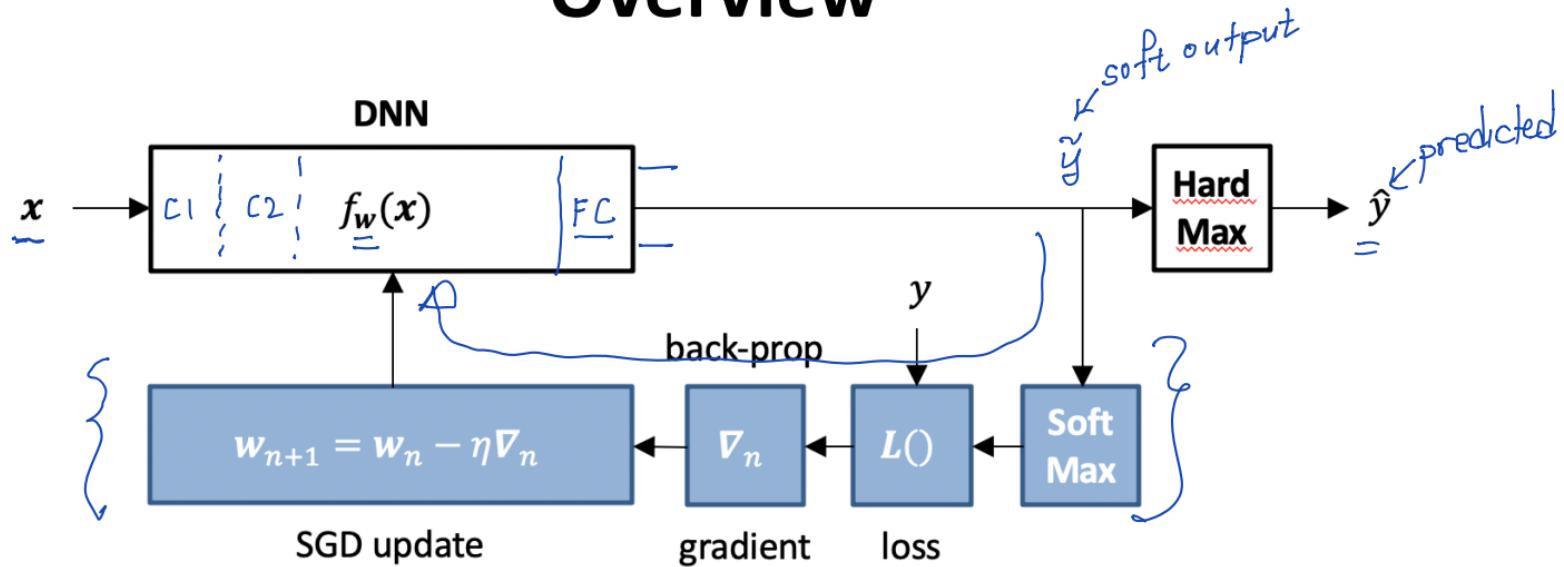
Supervised Training



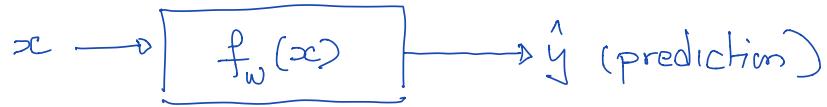
- **Sample:** $z = (x, y) = (\text{data}, \text{label})$;
- **Training sample (example):** samples used in training.
- **Loss function:** $L(\hat{y} = f_w(x), y)$; sample-specific value
- Family of functions parametrized by w : $\mathcal{F}: f_w(x) = \hat{y}$
- Loss function evaluated on the training set: $Q(w)$
- **learning (convergence) curves; learning rate;** stability

DNN Training Set-up

Overview



- Forward pass – generates predicted label \hat{y}
- Backward pass – SGD-based back-prop to update weights



$f_w(x)$: multi-stage & non-linear.

Q: How to find w : we get accurate predictions?

Q: What is a measure of accuracy? \rightarrow depends on the task.

Task \rightarrow classification

$$P_e = \Pr \left\{ \begin{array}{c} \hat{y} \neq y \\ \downarrow \\ \text{predicted label} \end{array} \right\}$$

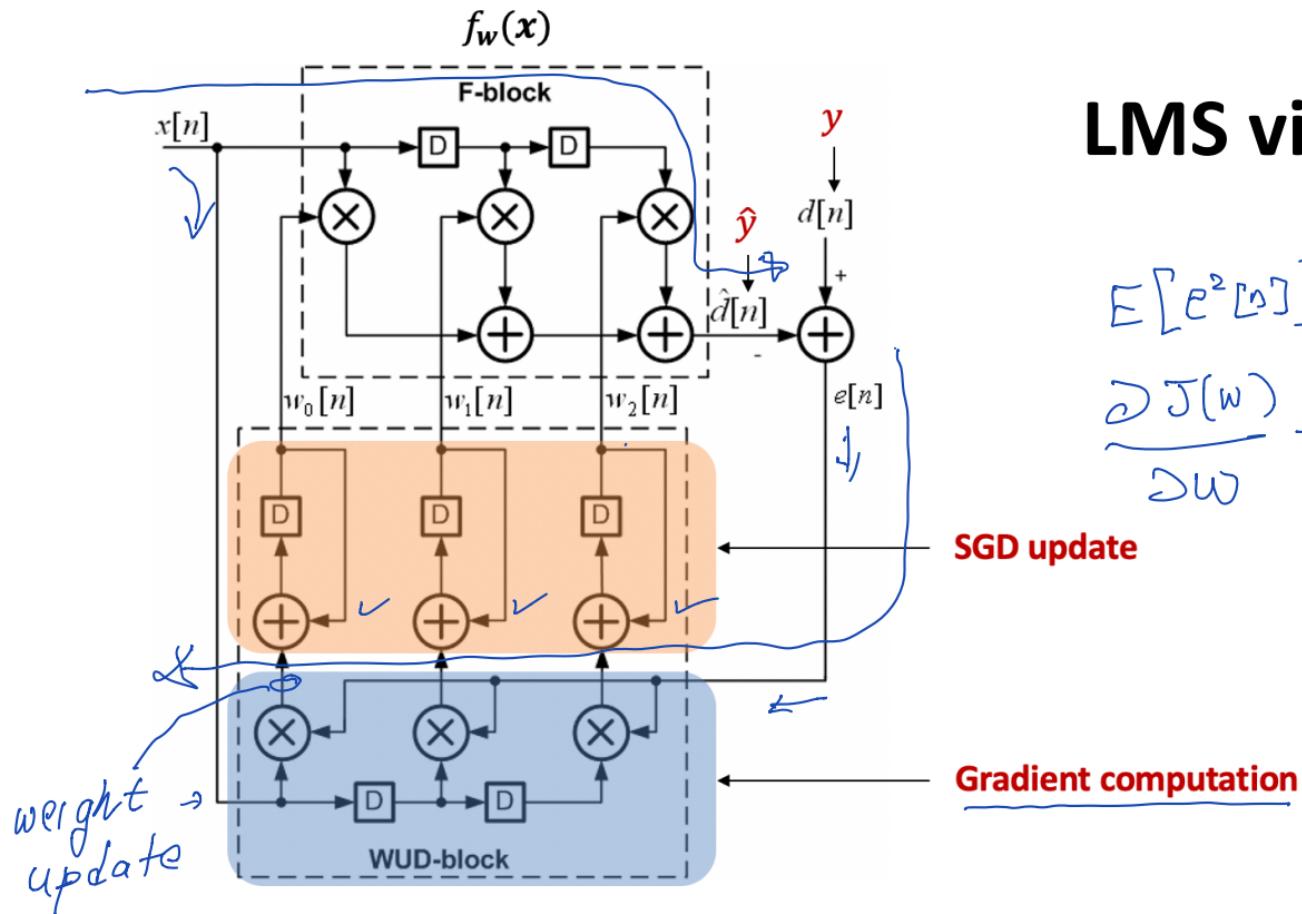
↑ true label

: misclassification or error rate

$1 - P_e$: accuracy

minimize P_e ? How do we adjust w : P_e is minimized? $\rightarrow \hat{y} \rightarrow y$

LMS via DNN Lens

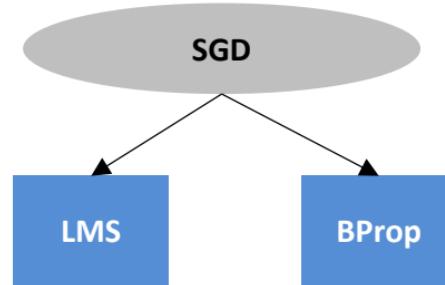
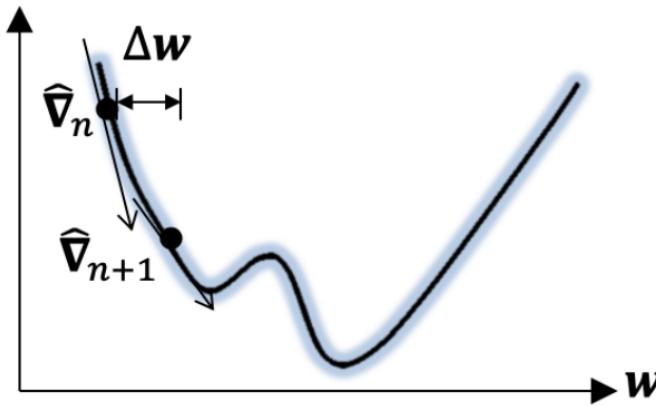


$$E[e^2[n]] = J(w)$$

$$\frac{\partial J(w)}{\partial w} \rightarrow \underbrace{\partial e[n]}_{\partial w} \times \underbrace{x[n]}_{\partial w}$$

$J(\mathbf{w})$

Stochastic Gradient Descent



- Calculate instantaneous gradient of loss function; update weight parameter by adding a small increment in the negative gradient

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \mu(-\hat{\nabla}_n)$$

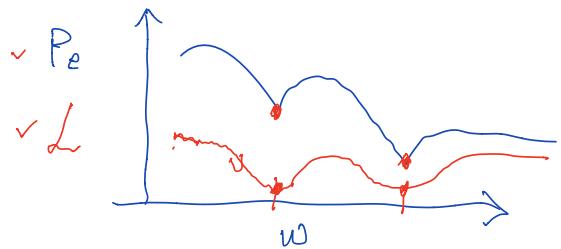
↗ -ve of gradient of
the loss fn. wrt. w.

- $\hat{\nabla}_n = \frac{\partial \hat{J}(\mathbf{w}_n)}{\partial \mathbf{w}_n}$ ⇒ is the gradient of the *instantaneous loss function* wrt \mathbf{w}_n
- μ : step-size or learning rate

minimize P_e : directly is hard :: $\underline{\underline{P_e(w)}}$

$$\frac{\partial P_e(w)}{\partial w} \rightarrow w_{n+1} = w_n + \mu \left(-\frac{\partial P_e(w)}{\partial w} \right)$$

Loss function: $\mathcal{L}(f_w(x), y) = \mathcal{L}(\hat{y}, y)$
 → \mathcal{L} is a surrogate/proxy for P_e



LMS: $MSE = \mathbb{E}[(\hat{y} - y)^2] \rightarrow y, \hat{y} \in \mathbb{R}$

For DNNs: \mathcal{L} is "cross entropy loss"

Entropy: of a RV X (discrete RV) → probability mass function $P_X(x)$

$$X \in \{-1, +1\} \rightarrow P_X(-1) = \frac{1}{2}; P_X(+1) = \frac{1}{2}$$

$$\text{Entropy of } X: H(X) = - \sum_i P_X(x_i) \log_2 P_X(x_i)$$

$$H(X) = \sum_i \underbrace{P_X(x_i)}_{\text{p min. # of bits on average}} \underbrace{\log_2 \left(\frac{1}{P_X(x_i)} \right)}_{\text{# of bits needed to rep. } x_i} = \mathbb{E}_{P_X} \left[\log_2 \left(\frac{1}{P_X(x)} \right) \right]$$

\downarrow p min. # of bits on average
 needed to represent x .

\cdot # of bits
 needed to rep. x_i

Cross-entropy: $X, Y \in \{-1, +1\}$, $X \sim \underline{P_X^Y}$, $Y \sim \underline{Q_X}$

$$\begin{aligned} H(X, Y) &= - \sum_i P_X(x_i) \log_2 (P_X(x_i)) \\ &= + \sum_i \underline{P_X(x_i)} \underbrace{\log_2 \left(\frac{1}{Q_X(x_i)} \right)}_{\text{bits assigned}} \end{aligned}$$

average # of bits needed to represent X when we assume $\underline{Q_X}$ as its distribution.

$$\underline{H(X, Y)} \geq \underline{H(X)}$$

In DNNs: $P_X \rightarrow$ true distribution of labels
 $\underline{Q_X} \rightarrow$ distribution generated network

update $\omega \rightarrow$ nudge Q_X towards P_X , \rightarrow "gap" or "distance" between Q_X & P_X is minimized.

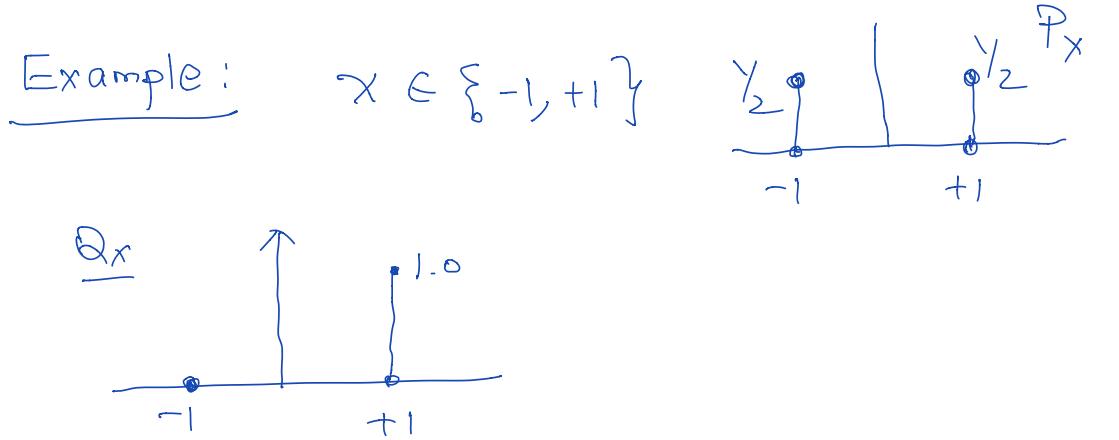
\downarrow
Kullback-Leibler divergence

$$\begin{aligned} D_{KL}(P_X || Q_X) &= \mathbb{E}_{P_X} \left(\log_2 \frac{P_X}{Q_X} \right) \\ &= \sum_i P_X(x_i) \log_2 \frac{P_X(x_i)}{Q_X(x_i)} \\ &= \sum_i P_X \log_2 P_X - \underbrace{\sum_i P_X \log_2 Q_X}_{= -H(X)} = -H(X) + H(X, Y) \end{aligned}$$

$$D_{KL}(P_x \parallel Q_x) = -H(x) + H(x, Y)$$

$\min D_{KL}(P_x \parallel Q_x) = \min H(x, Y) \rightarrow \text{minimize CE loss.}$

$$H(x, Y) \geq H(x) \rightarrow D_{KL}(P_x \parallel Q_x) \geq 0$$

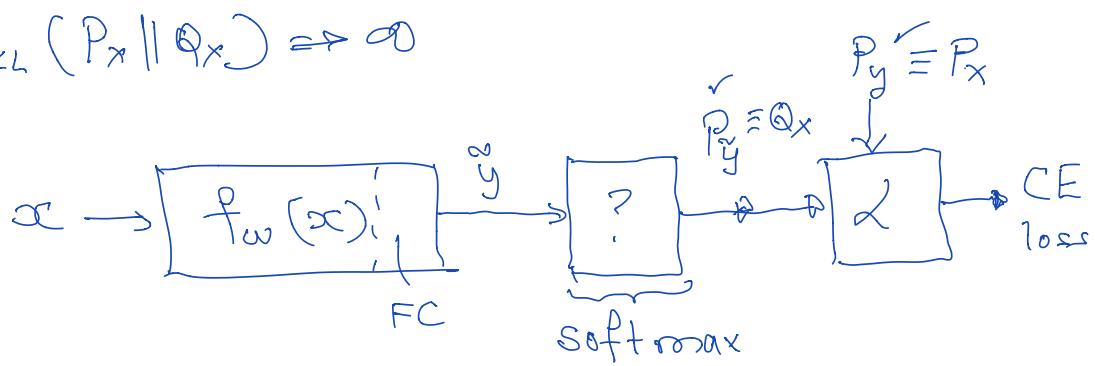


$$H(x) = -\left(\frac{1}{2} \times \log_2\left(\frac{1}{2}\right) + \frac{1}{2} \times \log_2\left(\frac{1}{2}\right)\right) = 1$$

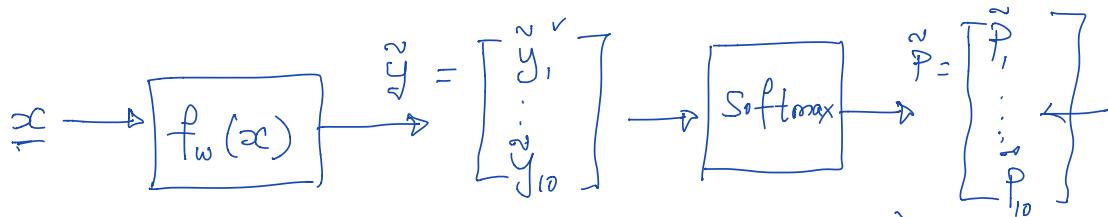
$$H(x, y) = -\left(\frac{1}{2} \times \log_2(0) + \frac{1}{2} \times \log_2(1)\right) \rightarrow \infty$$

\downarrow \downarrow
 $-\infty$ 0

$$D_{KL}(P_x \parallel Q_x) \Rightarrow \infty$$



DNN Example: 10-way classifier



$$\sum_i \tilde{P}_i = 1$$

$$\max(\tilde{y}) \rightarrow \hat{y}$$

$\tilde{P}_i = \frac{e^{\tilde{y}_i}}{\sum_i e^{\tilde{y}_i}}$ → softmax → generates the predicted distribution of class labels.

$\underline{\tilde{P}} = [0, \underbrace{1, 0, \dots, 0}]$ - one-hot encoded vector
class 2 is true label.

$$\begin{aligned} d(\hat{y}, y) &= -\log_2 \tilde{P}_i : i \leftarrow \text{true class} \\ &= -\sum_i \underbrace{P_i \log_2 \tilde{P}_i}_{\text{1-hot}} \end{aligned}$$

$$\frac{\partial (-\log_2 \tilde{P}_i)}{\partial \tilde{P}_i} = -\frac{1}{\tilde{P}_i \ln 2}$$

Ex: binary classifier

$$\hat{y}_1 = w_1 x_1 + w_2 x_2 ; \quad \hat{y}_2 = w_3 x_1 + w_4 x_2$$

True class = 1

$$d(\hat{y}, y) = -\log_2 \hat{P}_1$$
$$\hat{P}_1 = \frac{e^{\hat{y}_1}}{e^{\hat{y}_1} + e^{\hat{y}_2}}$$

$$\left(\frac{\partial d(\cdot)}{\partial w_1} \right) = ? \quad \text{— using chain rule}$$