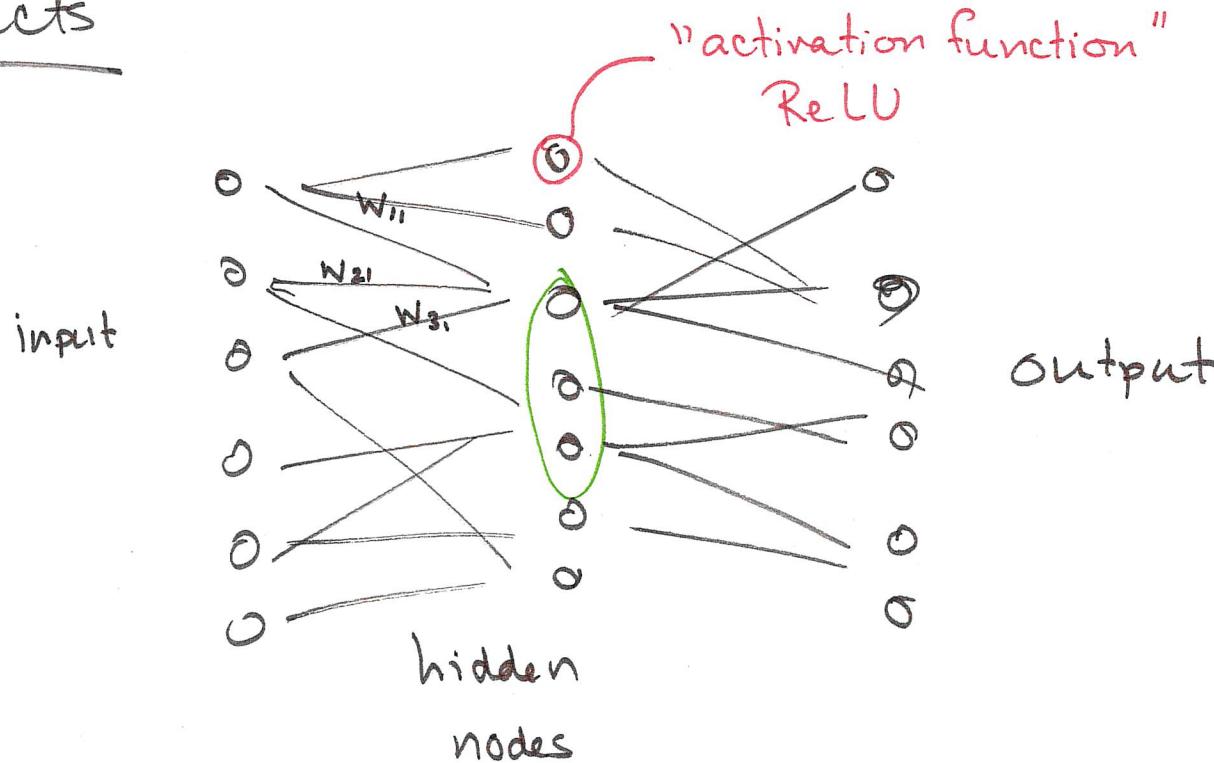


Projects



other ideas

- nonlinear extensions of PCA
 - kernel PCA
 - locally linear embeddings
 - nonnegative matrix factorization
 - autoencoder NN's

- sparse weights
or
small # of hidden nodes?
- data augmentation
- linear activation?

(autoencoder NN
w/ linear activation
= PCA)

- LASSO
 - group LASSO
 - elastic net (LASSO + Ridge)
 - multi-task learning.
- low-rank tensors

Lec 22 - SGD + Intro to Neural Networks

$$\underset{w}{\text{minimize}} \quad l(w) + \lambda r(w)$$

data fit regularizer

$f(w)$

$$\hat{w} = \arg \underset{w}{\min} \underline{f(w)}$$

} assume $f(w) = \sum_{i=1}^n f_i(w)$

e.g. $f(w) = \|y - Xw\|_2^2$
 $= \sum_i (y_i - x_i^T w)^2$

$$\Rightarrow f_i(w) = (y_i - x_i^T w)^2$$

SGD:

@ iteration k , choose $i_k \in \{1, \dots, n\}$

$$\hat{w}^{(k+1)} = \hat{w}^{(k)} - \frac{\tau}{2} \nabla f_{i_k}(\hat{w}^{(k)})$$

most commonly: $i_k \sim \text{uniformly at random}$
for $1, 2, \dots, n$

A. cyclic $i_k = k \bmod n$

if $n=3$, i_k 's : 1, 2, 3, 1, 2, 3, ...

B. permutations

i_k 's : $\underbrace{3, 1, 2}_{\text{epoch 1}}, \underbrace{1, 3, 2}_{\text{epoch 2}}, \underbrace{2, 1, 3}_{\text{...}}$

C. i_k 's uniform at random

i_k 's : $\underbrace{1, 1, 1, 1, 3, 1, 2, 2, 3, 1}_{\text{...}}$

if i_k 's are uniform at random,

$$E[f_{i_k}] = f/n$$

Ex:

$$f(w) = \|y - Xw\|_2^2 + \lambda \|w\|_2^2$$
$$= \sum_{i=1}^n \left[(y_i - x_i^T w)^2 + \frac{\lambda}{n} \|w\|_2^2 \right]$$

$f_i(w)$

$$\nabla f_i(w) = -2(y_i - x_i^T w)x_i + \frac{2\lambda}{n} w$$

SGD:

$$\hat{w}^{(k+1)} = \hat{w}^{(k)} - \frac{\eta}{2} \tau \left[-(y_{i_k} - x_{i_k}^T \hat{w}^{(k)})x_{i_k} + \frac{2\lambda}{n} \hat{w}^{(k)} \right]$$

can replace gradients with subgradients:

Recall: if f is convex and differentiable:

$$f(u) \geq f(w) + (\underline{u} - w)^T \nabla f(w)$$

if f is convex but not differentiable,

then v is a subgradient

of f at w if

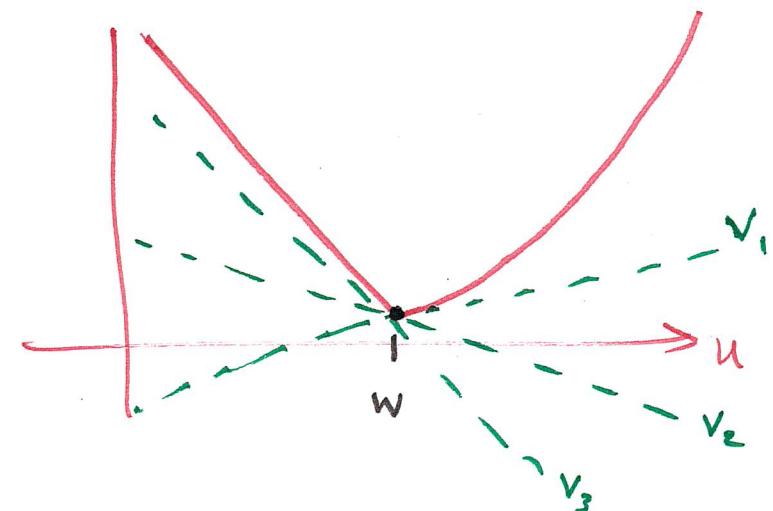
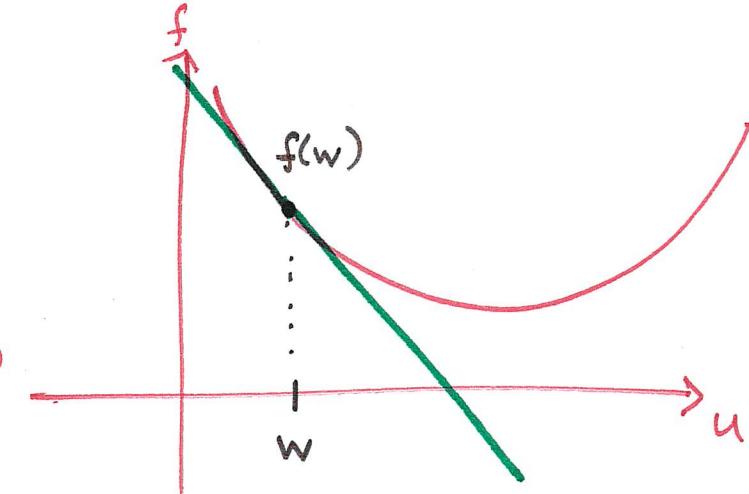
$$f(u) \geq f(w) + (\underline{u} - w)^T v$$

set of subgradients @ w

is called "differential set"

denoted $\partial f(w)$

write $v \in \partial f(w)$



e.g. $r(w) = \|w\|_1 = \sum_{j=1}^p |w_j|$

for $w_j \neq 0$, $|w_j|$ is differentiable

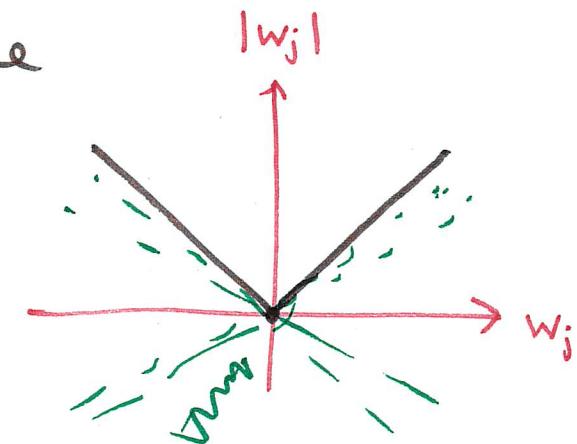
deriv is $= \text{Sign}(w_j)$

for $w_j = 0$, then $v_j \in [-1, +1]$

for $\underline{v} \in \partial r(w)$

then $v_j = \begin{cases} = \text{sign}(w_j) & \text{if } w_j \neq 0 \\ \in [-1, +1] & \text{if } w_j = 0 \end{cases}$

popular choice : $v_j = \begin{cases} \text{sign}(w_j) & \text{if } w_j \neq 0 \\ 0 & \text{if } w_j = 0 \end{cases} = " \text{sign}(w) "$



(6)

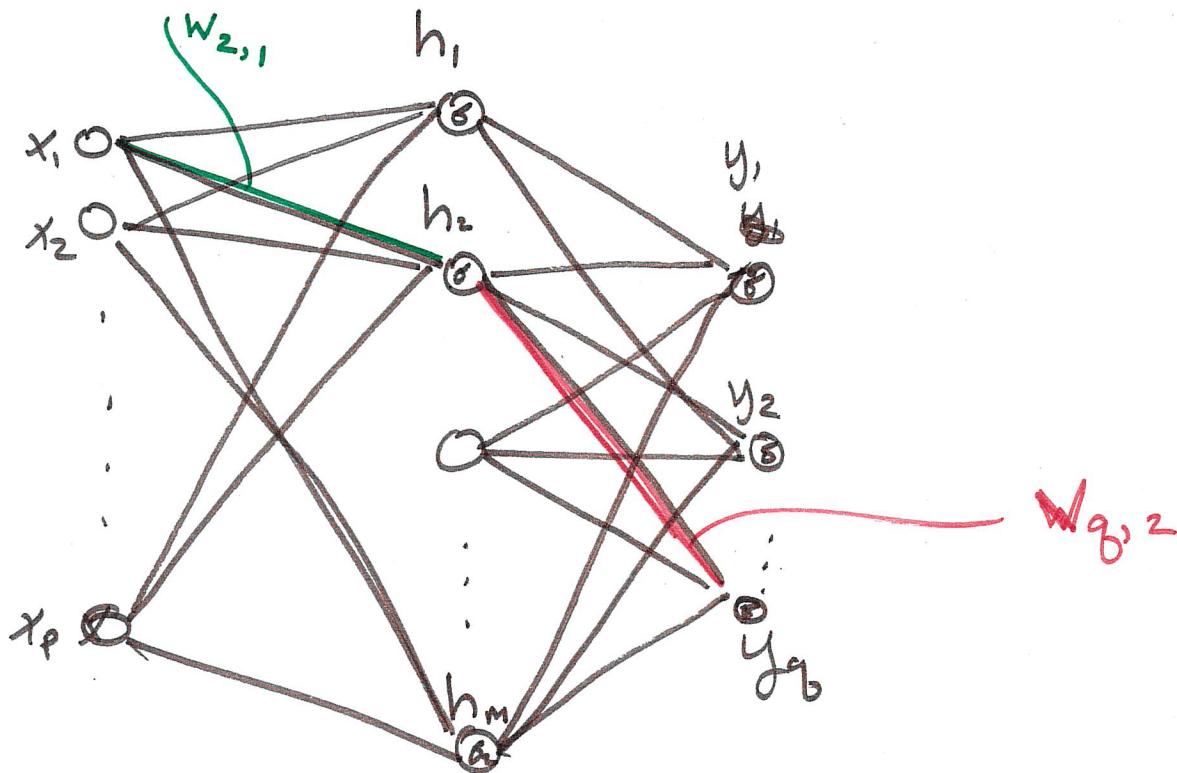
E_x 2 SGD : $f(w) = \|y - Xw\|_2^2 + \lambda \|w\|_1$ (LASSO)

$$= \sum_{i=1}^n \underbrace{\left[(y_i - x_i^T w)^2 + \frac{\lambda}{n} \|w\|_1 \right]}_{f_i(w)}$$

let $v = -2(y_i - x_i^T w)^* x_i + \frac{\lambda}{n} \text{sign}(w)$

SGD: $\hat{w}^{(k+1)} = \hat{w}^{(k)} - \frac{T}{2} (-v_i)$

$$= \hat{w}^{(k)} + T(y_{i_k} - x_{i_k}^T \hat{w}^{(k)}) - \frac{T\lambda}{2n} \text{sign}(\hat{w}^{(k)})$$

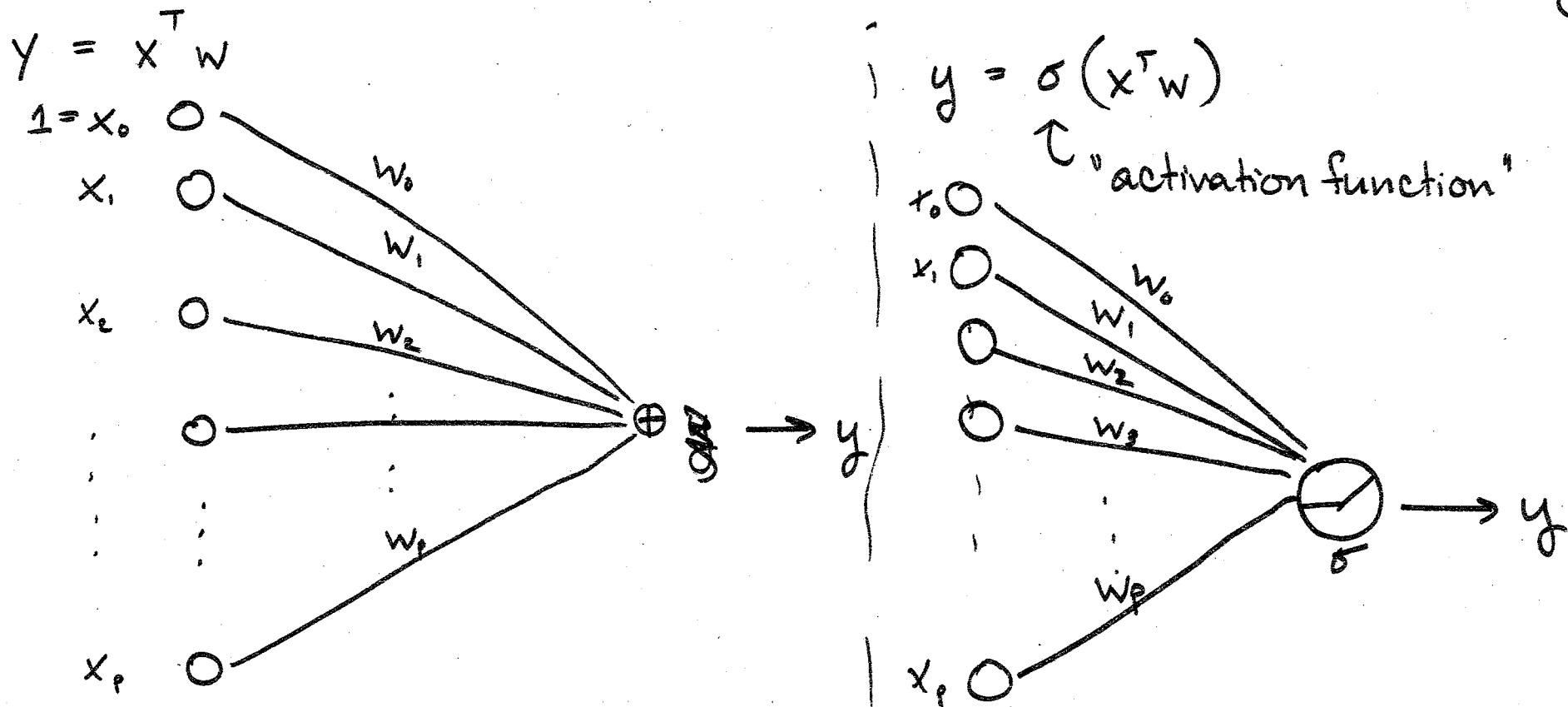


m -th hidden node's output

$$h_m = \sigma \left(\sum_{j=1}^p x_j w_{m,j} \right)$$

k^{th} output node

$$y_k = \sigma \left(\sum_{m=1}^M h_m v_{k,m} \right)$$



↳ $\sigma(z) = \max(0, z) = \text{ReLU}$

$\sigma(z) = \text{sign}(z) \in [-1, 0, +1]$

$\sigma(z) = \frac{1}{1+e^{-z}} \in [0, 1]$

