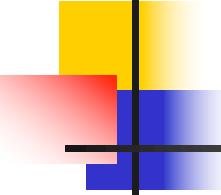


CMSC 510 – L05

Regularization Methods for Machine Learning



Instructor:
Dr. Tom Arodz



Big picture...

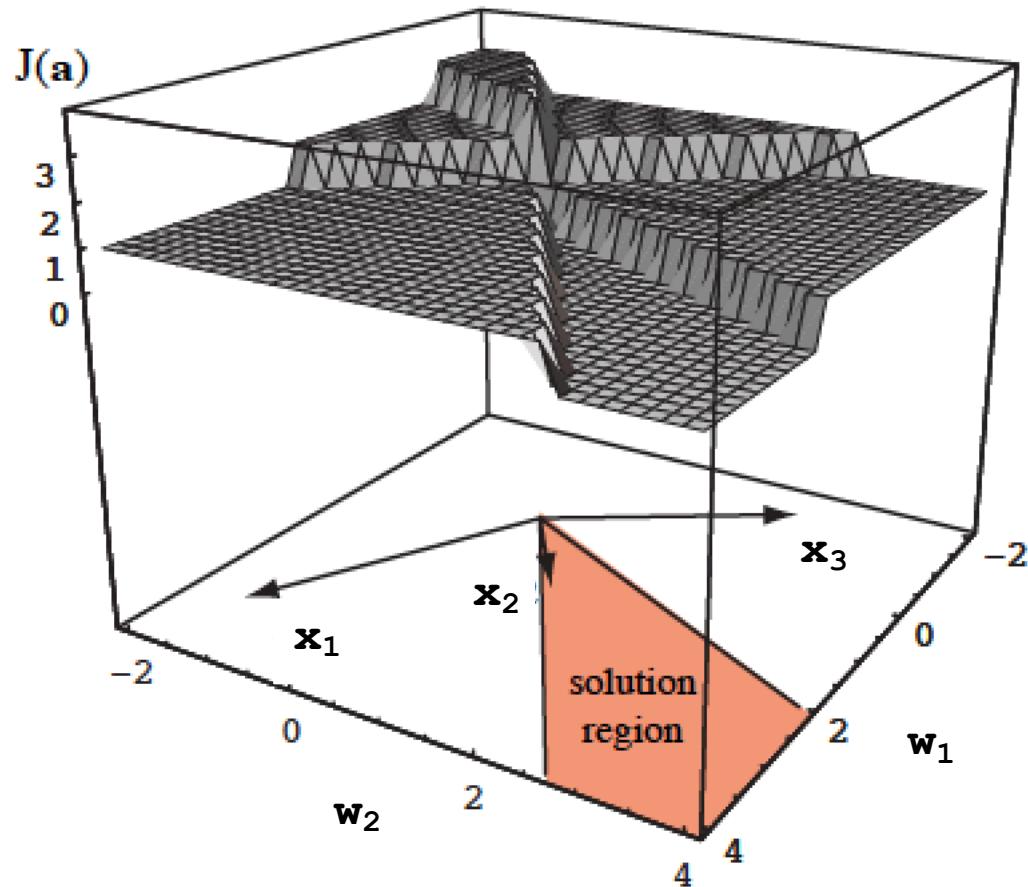
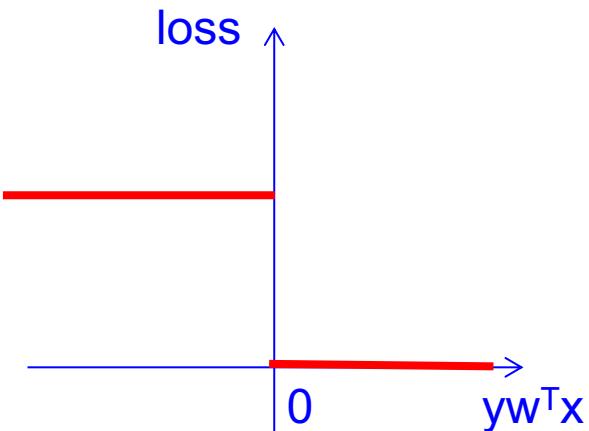
- Choose predictive model
 - We're focusing on $h(x) = wx + b$ (linear models)
- Choose loss function
 - E.g. $\text{loss} = (y - h(x))^2$
 - Loss depends on y, x – given to us (training set)
 - Loss also depends on w, b – we need to find those
- Choose the algorithm for finding model parameters (w, b for linear models) that lead to low avg. loss
 - Gradient descent is our default choice
 - $w_{\text{next}} = w_{\text{current}} - \frac{d \text{ loss}}{d w}$

Recap: 0/1 loss

- Present a sample x and predict:

$$h(x) = \text{sign}(w^T x - w_0)$$

- Why not use 0/1 loss?
- It's flat! We don't know in which direction to move to get a better solution!
- At each possible w , we have # of misclassified points as our empirical risk



Recap: Perceptron loss

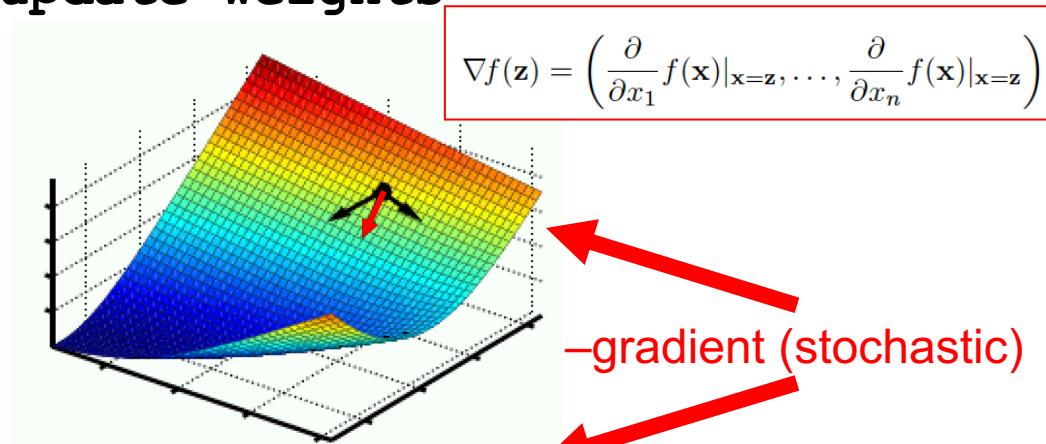
- Present a sample \mathbf{x} and predict:

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} - w_0)$$

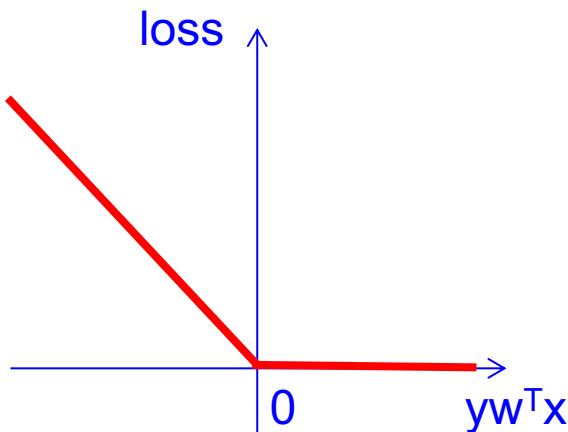
- Compare true class y with predicted class $h(\mathbf{x})$
- If prediction is wrong, update weights

$$\mathbf{w}^{\dagger}_{t+1} = \mathbf{w}^{\dagger}_t + 2c y \mathbf{x}^{\dagger}$$

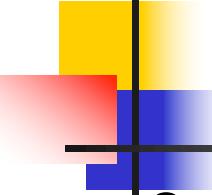
- Minimize empirical risk of perceptron loss:
- By modifying w after seeing each sample z



$$\mathbf{w}^{\dagger}_{t+1} = \mathbf{w}^{\dagger}_t - 2c \left. \frac{\partial \ell(h, \mathbf{z})}{\partial \mathbf{w}^{\dagger}} \right|_{\mathbf{w}^{\dagger}=\mathbf{w}^{\dagger}_t}$$



$$\ell(h, \mathbf{z}) = \begin{cases} 0 & \text{for } y \mathbf{w}^{\dagger T} \mathbf{x}^{\dagger} \geq 0 \\ -y \mathbf{w}^{\dagger T} \mathbf{x}^{\dagger} & \text{for } y \mathbf{w}^{\dagger T} \mathbf{x}^{\dagger} < 0 \end{cases}$$



(stochastic) gradient descent

- Our function to be minimized is (proportional to) empirical risk:
$$f(w) = \sum_i f_i(w) = \sum_i \text{loss}(y_i, x_i, w)$$
- Gradient descent (batch learning):
 - in a loop, modify w to minimize $f(w)$
 - that means we need to evaluate the sum over all samples
 - we use gradient of $\sum_i \text{loss}(y_i, x_i, w)$ w.r.t. w
 - We present all samples, and only then we update the weights
- Stochastic gradient descent (online learning):
 - in a loop, randomly pick a single i , modify w to minimize $f_i(w)$
 - that means we need to evaluate only one sample at a time
 - we use gradient of $\text{loss}(y_i, x_i, w)$ w.r.t. w
 - We present one sample, and immediately update the weights
 - Each step is then a “noisy” gradient, but the noise should average itself out over many iterations
- mini-batch stochastic gradient descent :
 - In between the above: $\sum_i \text{loss}(y_i, x_i, w)$ over e.g. randomly chosen 64 samples from the whole training set

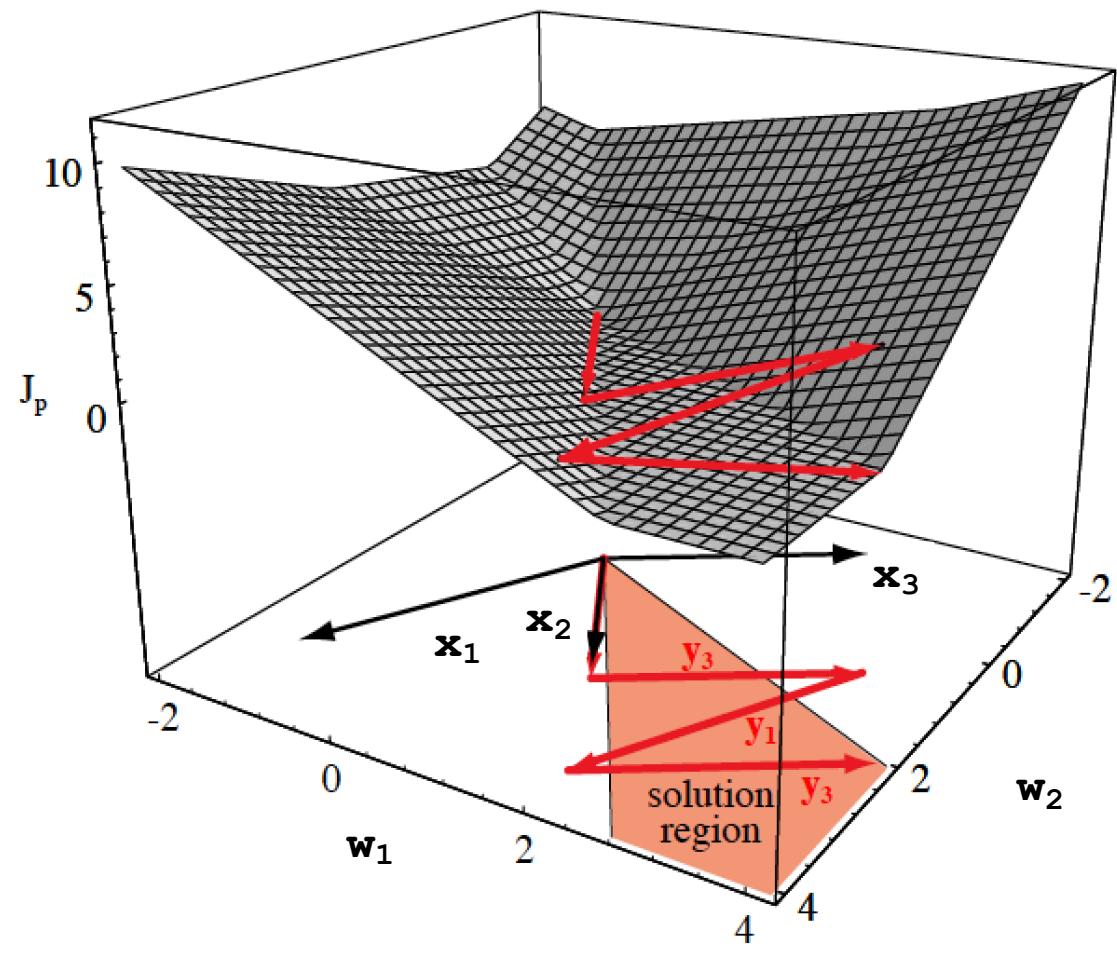
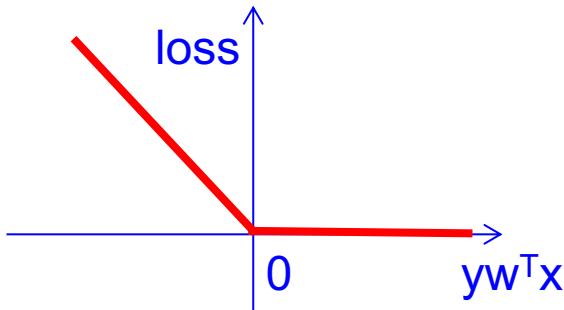
Perceptron

- Present a sample x and predict:

$$h(x) = \text{sign}(\mathbf{w}^T \mathbf{x} - w_0)$$

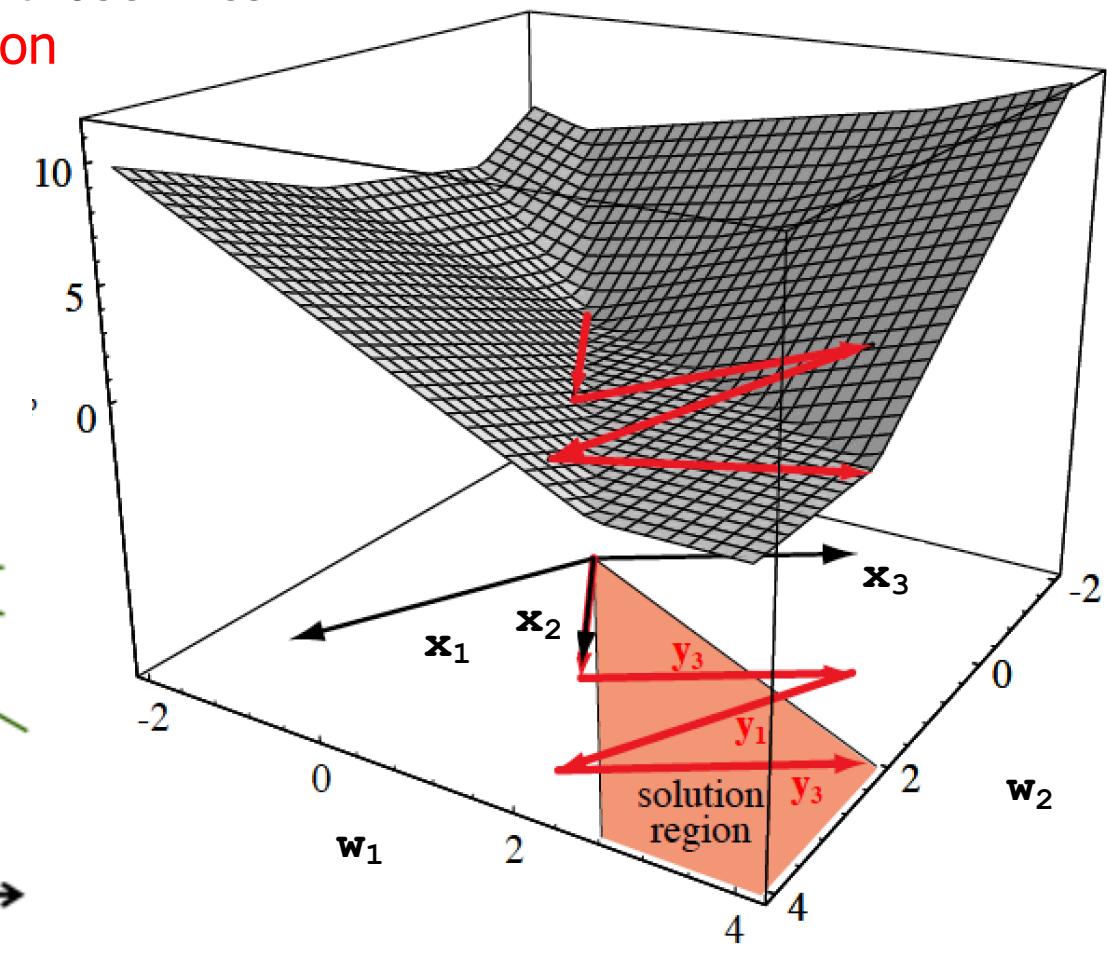
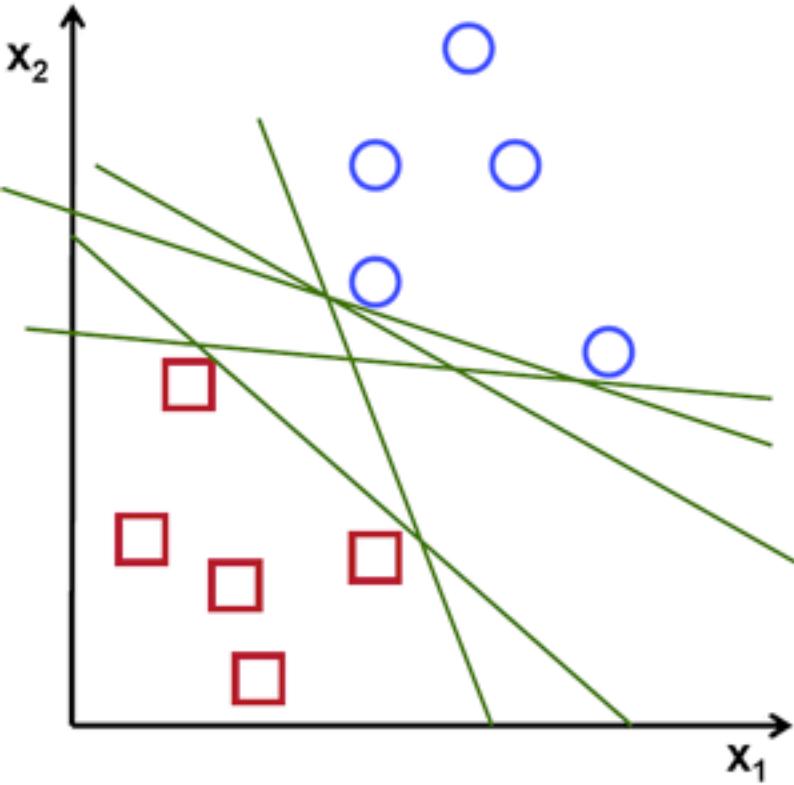
- Example: three points
- Two parameters
- Surface plot of risk:

$$\hat{R}_{S_m}(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, z_i)$$



Solution region

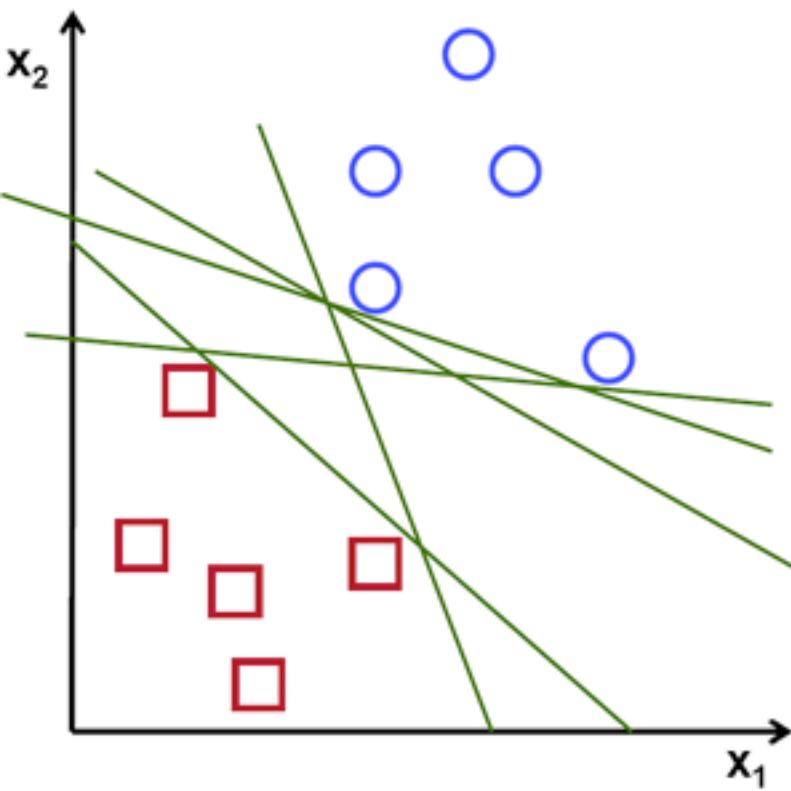
- If the two classes can be separated by a straight line (**problem/dataset is linearly separable**)
 - Then there's infinite number of lines that separate the classes
 - The coefficients w for those lines form the **solution region**



Separable problems

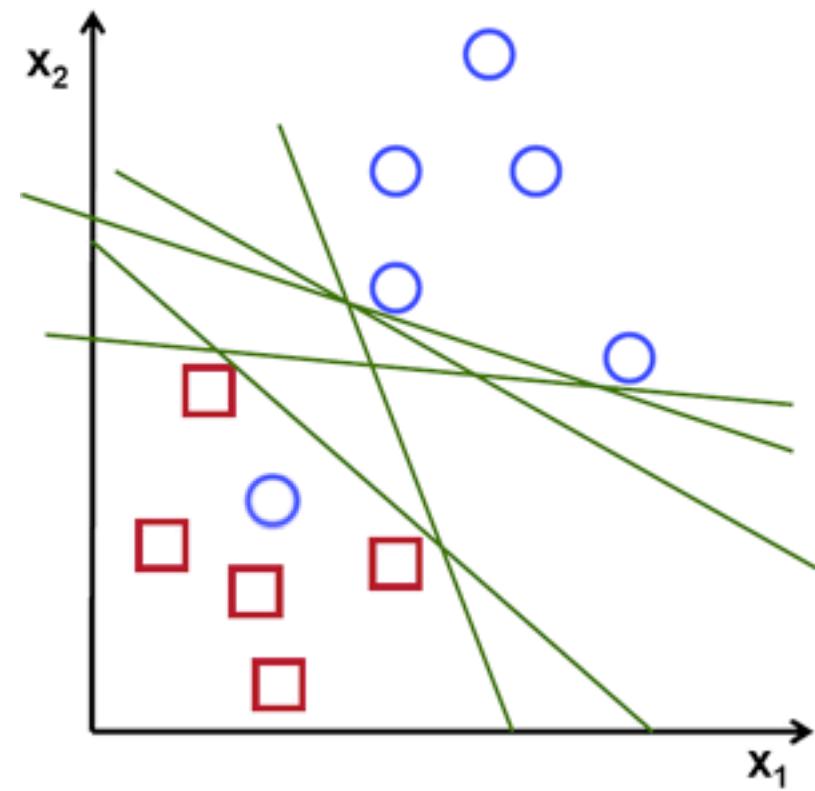
- Some problems/datasets are linearly non-separable

Separable problem



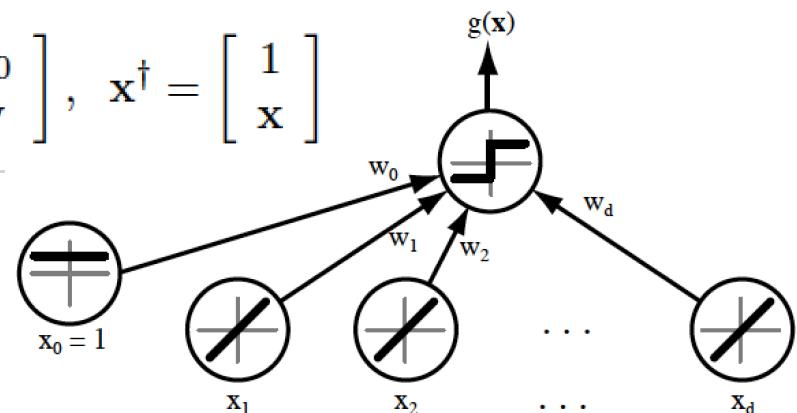
Non-separable problem

none of the lines works
no other line would work either



Perceptron

$$\mathbf{w}^\dagger = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}, \quad \mathbf{x}^\dagger = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$$

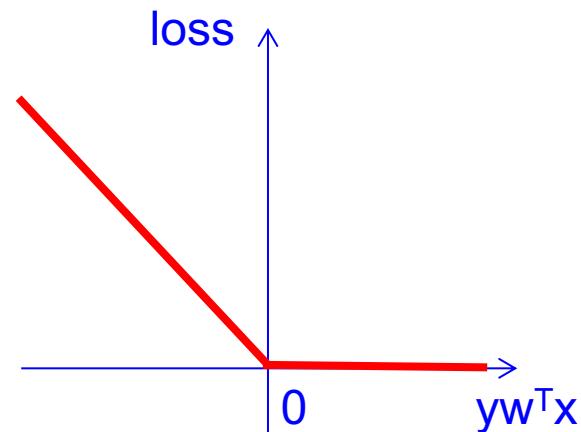


- Present a sample \mathbf{x} and predict:

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} - w_0)$$

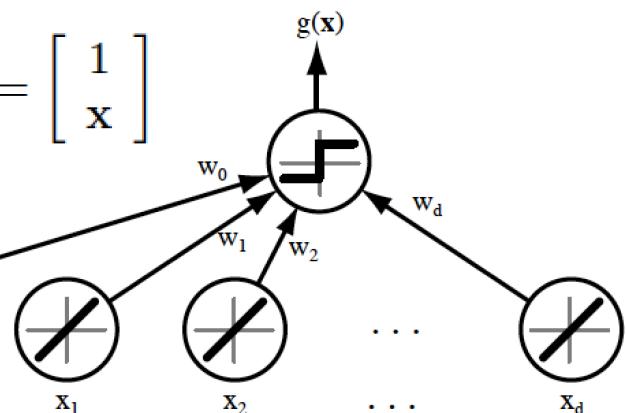
- Some mathematical problems:
 - Perceptron loss is not differentiable at 0
 - We would have to use subgradients instead of gradients
 - Loss is convex, so risk is convex
 - For separable problems, the algorithm eventually converges to optimal solution with 0 risk
 - For non-separable problems, we'll always have errors**
 - So the risk will always be >0 ???

$$\ell(h, z) = \begin{cases} 0 & \text{for } y\mathbf{w}^\dagger T \mathbf{x}^\dagger \geq 0 \\ -y\mathbf{w}^\dagger T \mathbf{x}^\dagger & \text{for } y\mathbf{w}^\dagger T \mathbf{x}^\dagger < 0 \end{cases}$$



Perceptron

$$\mathbf{w}^\dagger = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}, \quad \mathbf{x}^\dagger = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$$



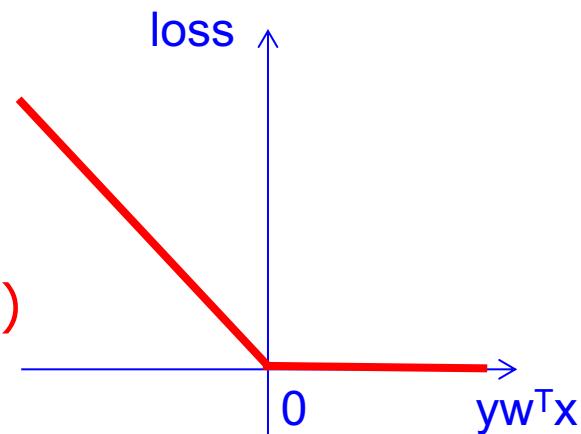
- Present a sample \mathbf{x} and predict:

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} - w_0)$$

- Some mathematical problems:

- Perceptron loss is not differentiable at 0
- We would have to use subgradients instead of gradients
- Loss is convex, so risk is convex
- For separable problems, the algorithm eventually converges to optimal solution with 0 risk
- For non-separable problems, risk has a single global minimum (with risk=0) at $\mathbf{w}=0$ (useless!!!)

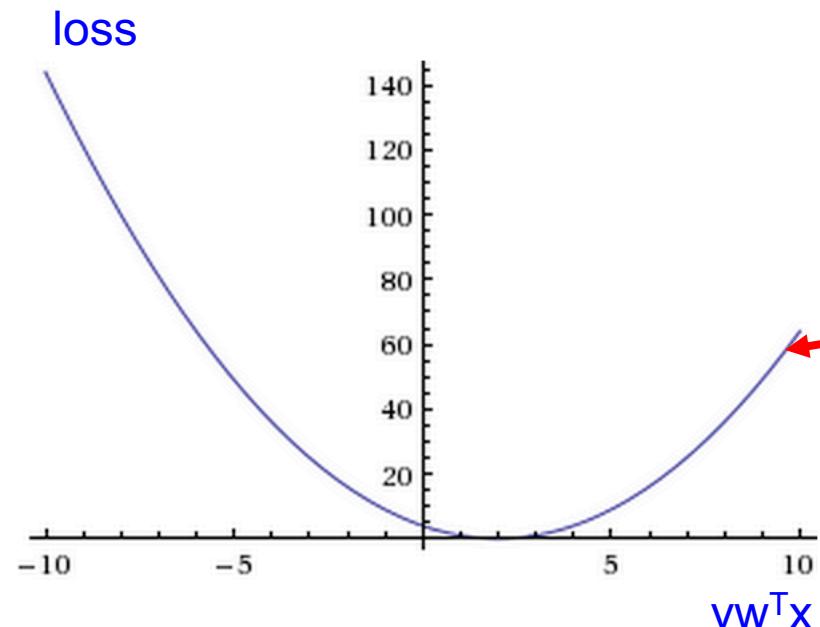
$$\ell(h, z) = \begin{cases} 0 & \text{for } y\mathbf{w}^\dagger T \mathbf{x}^\dagger \geq 0 \\ -y\mathbf{w}^\dagger T \mathbf{x}^\dagger & \text{for } y\mathbf{w}^\dagger T \mathbf{x}^\dagger < 0 \end{cases}$$



Losses so far

All have their problems:

- 0/1 loss: flat
- Perceptron loss: “no prediction”
- MSE for classification: penalty for correct predictions
- Let’s try to address that



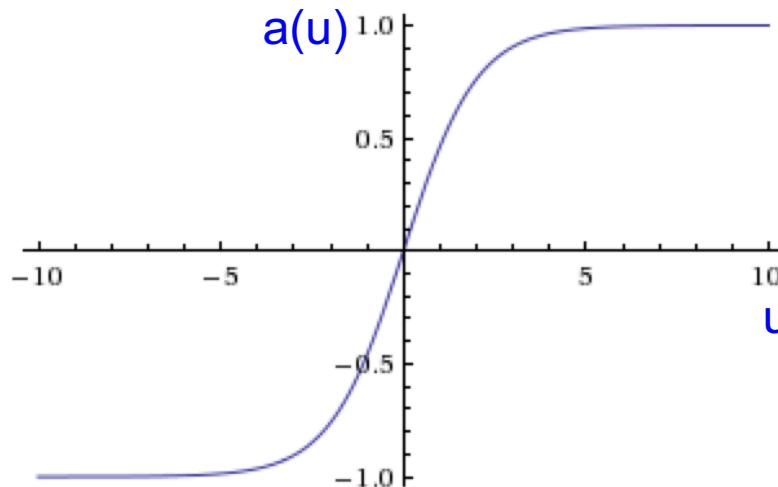
Differentiable perceptron ($y=-1/1$)

- Let's limit the predictions to $[-1,1]$:

$$h(x) = a(w^T x)$$

- If $h(x) > 0$ we predict +1
- If $h(x) < 0$ we predict -1

$$a(u) = \frac{2}{1 + e^{-u}} - 1$$



$a(u)$ is “activation function”

<- bipolar sigmoid
activation function

- Derivative of $a(u)$ with respect to u is (verify at home!):

$$a'(u) = \frac{2e^{-u}}{(1 + e^{-u})^2} = \frac{1}{2} (1 - a^2(u))$$

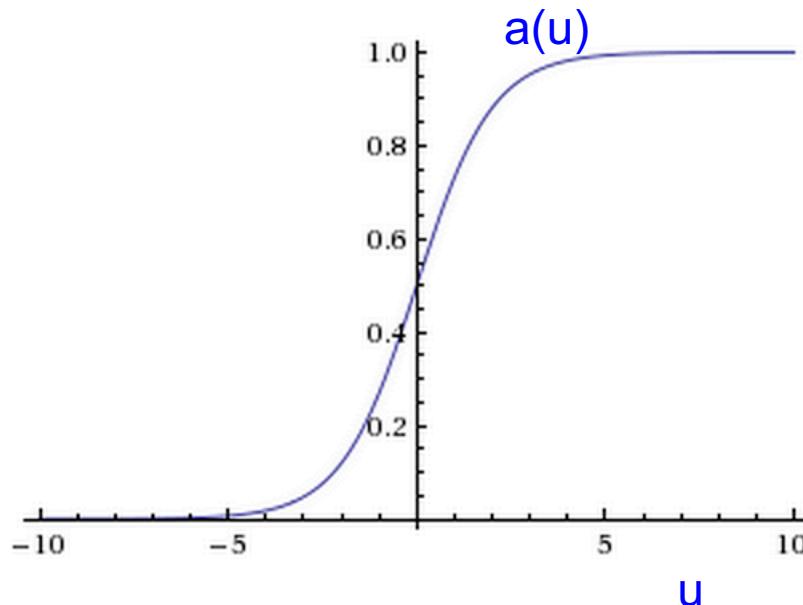
Differentiable perceptron ($y=0/1$)

- Another variant:

$$h(x) = a(w^T x)$$

- If $h(x) > 0.5$ we predict +1
- If $h(x) < 0.5$ we predict 0

$$a(u) = \frac{1}{1+e^{-u}}$$



$a(u)$ is “activation function”

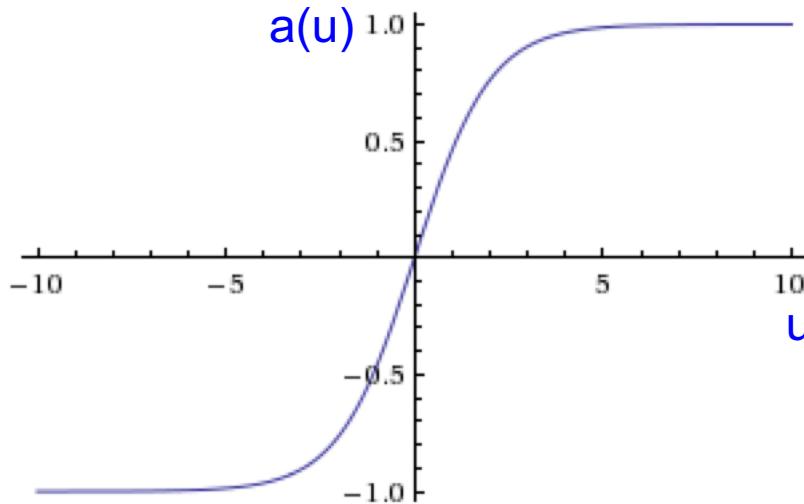
<- unipolar sigmoid
activation function

Differentiable perceptron

- Let's make h differentiable: $h(x) = a(w^T x)$

- If $h(x) > 0$ we predict +1
- If $h(x) < 0$ we predict -1

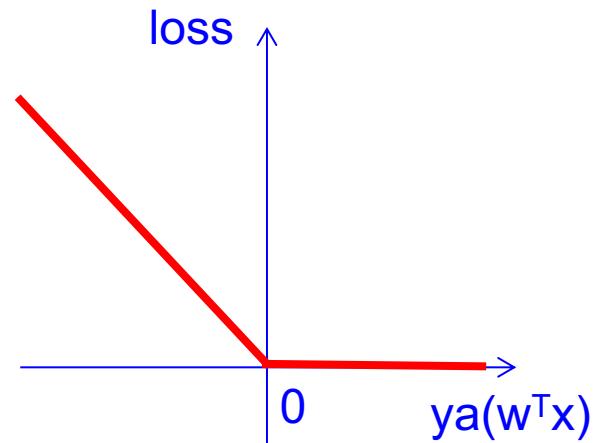
$$a(u) = \frac{2}{1 + e^{-u}} - 1$$



$a(u)$ is “activation function”

<- sigmoid activation function

- If we use the loss: $\max(0, -ya(w^T x))$
 - Still the same problem:
 - No prediction ($w=0$) is optimal



Differentiable perceptron

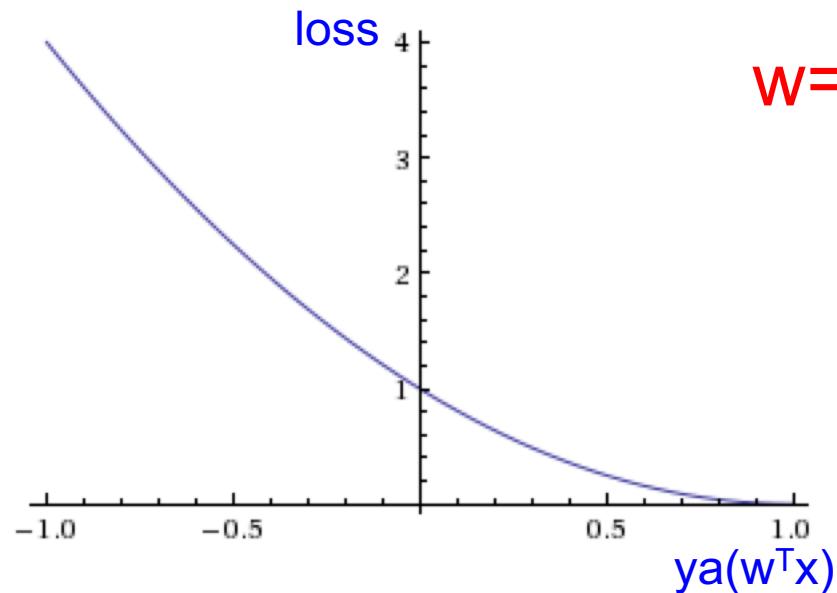
- $h(x) = a(w^T x)$

$$a(u) = \frac{2}{1 + e^{-u}} - 1$$

- Quadratic loss on sigmoid predictions:

$$\ell(h, z) = (y - a(w^T x)) ^ 2 = (1 - ya(w^T x)) ^ 2$$

How did we get this?



w=0 is no longer the best solution!

we won't be stuck with
“no prediction”
as a solution
that is always
a global minimum

Differentiable perceptron

- Assuming:

$$a(u) = \frac{2}{1 + e^{-u}} - 1$$

$$\ell(h, z) = (y - a(w^\dagger{}^T x^\dagger))^2 = (1 - ya(w^\dagger{}^T x^\dagger))^2$$

- We can start deriving the update of weights:

$$w^\dagger_{t+1} = w^\dagger_t - \frac{c}{2} \left. \frac{\partial \ell(h, z)}{\partial w^\dagger} \right|_{w^\dagger=w^\dagger_t}$$

$$w^\dagger_{t+1} = w^\dagger_t - \frac{c}{2} \left. \frac{\partial (1 - ya(w^\dagger{}^T x^\dagger))^2}{\partial w^\dagger} \right|_{w^\dagger=w^\dagger_t}$$

$$w^\dagger_{t+1} = w^\dagger_t - \frac{c}{2} \left. \frac{\partial (1 - ya(w^\dagger{}^T x^\dagger))^2}{\partial a(w^\dagger{}^T x^\dagger)} \frac{\partial a(w^\dagger{}^T x^\dagger)}{\partial w^\dagger} \right|_{w^\dagger=w^\dagger_t}$$

Differentiable perceptron

■ Assuming: $a(u) = \frac{2}{1 + e^{-u}} - 1$ $a'(u) = \frac{2e^{-u}}{(1 + e^{-u})^2} = \frac{1}{2} (1 - a^2(u))$

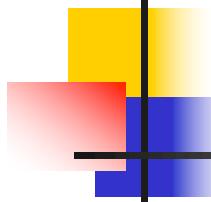
■ Continuing the derivation of the update of weights:

$$\mathbf{w}^\dagger_{t+1} = \mathbf{w}^\dagger_t - \frac{c}{2} \frac{\partial \left(1 - ya(\mathbf{w}^\dagger T \mathbf{x}^\dagger)\right)^2}{\partial a(\mathbf{w}^\dagger T \mathbf{x}^\dagger)} \frac{\partial a(\mathbf{w}^\dagger T \mathbf{x}^\dagger)}{\partial \mathbf{w}^\dagger} \Bigg|_{\mathbf{w}^\dagger = \mathbf{w}^\dagger_t}$$

$$\mathbf{w}^\dagger_{t+1} = \mathbf{w}^\dagger_t - \frac{c}{2} \left[2y^2 a(\mathbf{w}^\dagger_t T \mathbf{x}^\dagger) - 2y \right] \frac{\partial a(\mathbf{w}^\dagger T \mathbf{x}^\dagger)}{\partial \mathbf{w}^\dagger T \mathbf{x}^\dagger} \frac{\partial \mathbf{w}^\dagger T \mathbf{x}^\dagger}{\partial \mathbf{w}^\dagger} \Bigg|_{\mathbf{w}^\dagger = \mathbf{w}^\dagger_t}$$

$$\mathbf{w}^\dagger_{t+1} = \mathbf{w}^\dagger_t + c \left[y - a(\mathbf{w}^\dagger_t T \mathbf{x}^\dagger) \right] \frac{1}{2} \left[1 - a^2(\mathbf{w}^\dagger_t T \mathbf{x}^\dagger) \right] \mathbf{x}^\dagger$$

$$\mathbf{w}^\dagger_{t+1} = \mathbf{w}^\dagger_t + \frac{c}{2} \left[y - a(\mathbf{w}^\dagger_t T \mathbf{x}^\dagger) \right] \left[1 - a^2(\mathbf{w}^\dagger_t T \mathbf{x}^\dagger) \right] \mathbf{x}^\dagger$$



Differentiable perceptron

- a perceptron with differentiable activation function $a(w^T x)$:
- quadratic loss:

$$\ell(h, z) = \left(y - a(w^\dagger{}^T x^\dagger) \right)^2 = \left(1 - ya(w^\dagger{}^T x^\dagger) \right)^2$$

- This weight update rule is called **delta rule**:

$$w^\dagger_{t+1} = w^\dagger_t - \frac{c}{2} \frac{\partial \ell(h, z)}{\partial w^\dagger} \Big|_{w^\dagger=w^\dagger_t} = w^\dagger_t + c \left[y - a(w^\dagger_t{}^T x^\dagger) \right] a' \left(w^\dagger_t{}^T x^\dagger \right) x^\dagger.$$

- In our specific case of activation function $a(u) = \frac{2}{1 + e^{-u}} - 1$ we get:

$$w^\dagger_{t+1} = w^\dagger_t - \frac{c}{2} \frac{\partial \ell(h, z)}{\partial w^\dagger} \Big|_{w^\dagger=w^\dagger_t} = w^\dagger_t + \frac{c}{2} \left[y - a(w^\dagger_t{}^T x^\dagger) \right] \left[1 - a^2(w^\dagger_t{}^T x^\dagger) \right] x^\dagger$$

Differentiable perceptron $h(x)=a(w^T x)$

- $h(x)=a(w^T x)$

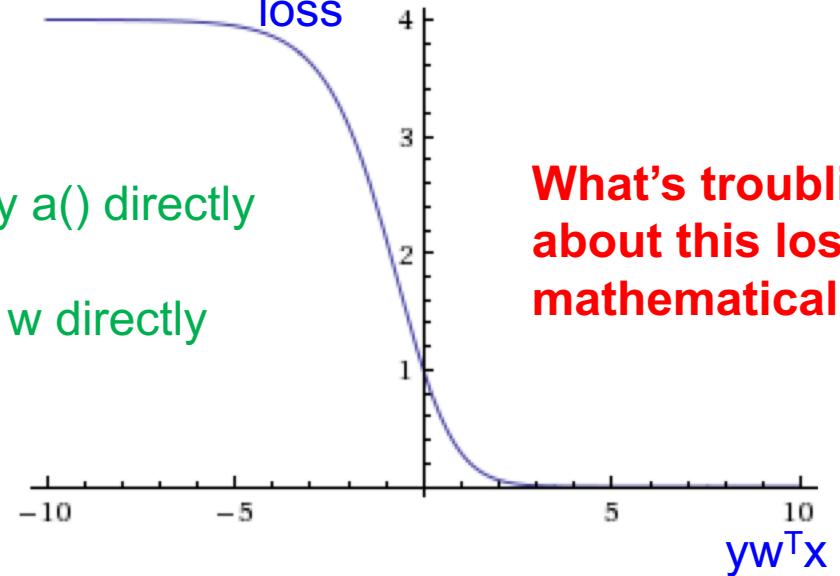
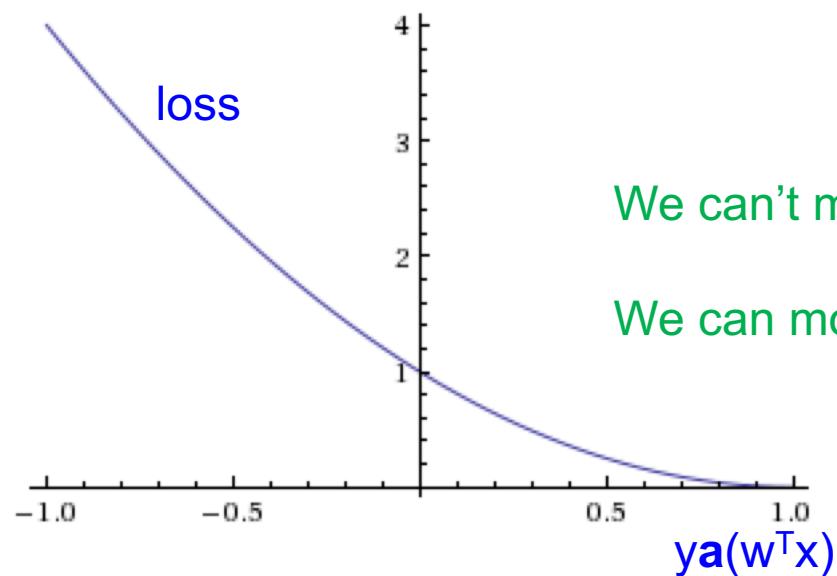
$$a(u) = \frac{2}{1 + e^{-u}} - 1.$$

- Loss as a function of y and of $a(w^T x)$ is:

$$\ell(h, z) = \left(y - a(w^T x) \right)^2 = \left(1 - ya(w^T x) \right)^2$$

- If we modify w for fixed sample (x, y) , how is loss affected?

- Loss should be expressed in function of w (or $yw^T x$, which is proportional to w for fixed training sample, and makes "correct" always on the right)



Differentiable perceptron

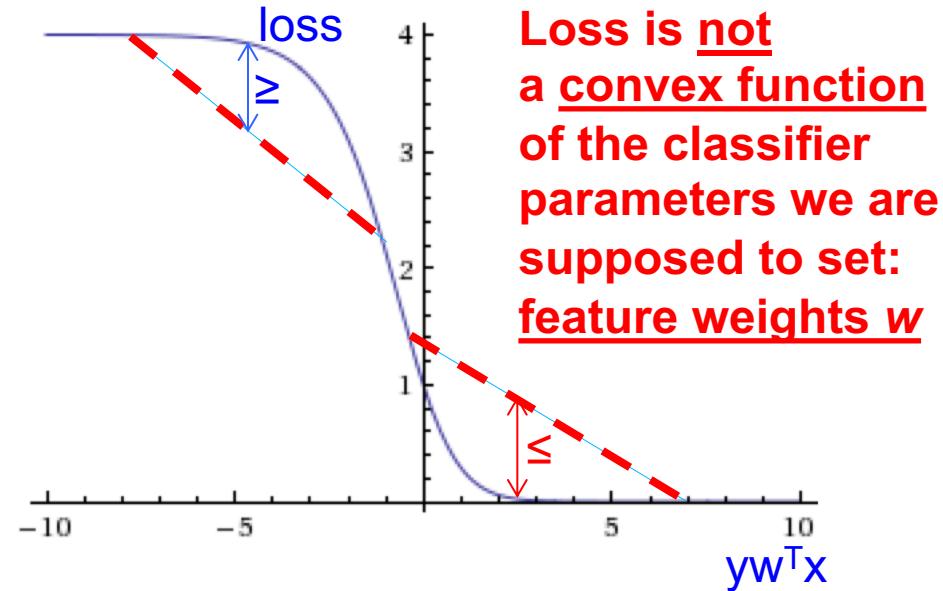
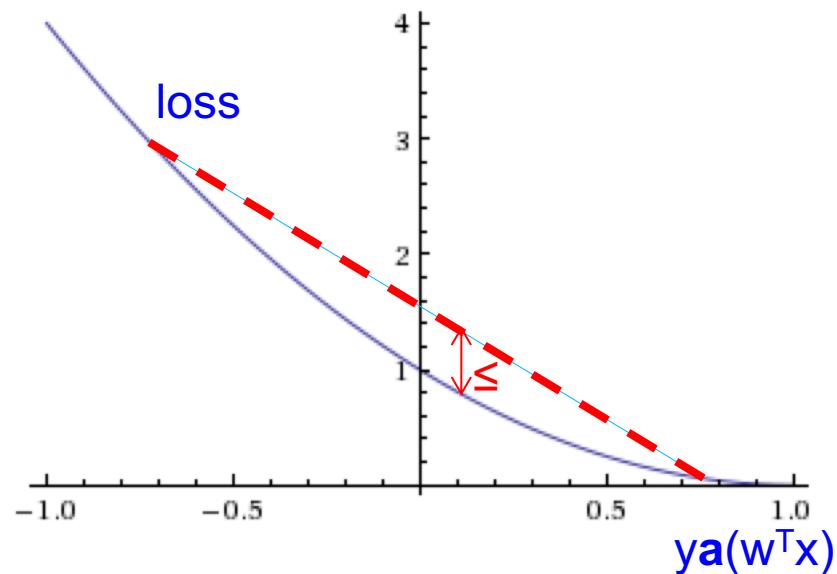
- $h(x) = a(w^T x)$

$$a(u) = \frac{2}{1 + e^{-u}} - 1$$

- Loss is:

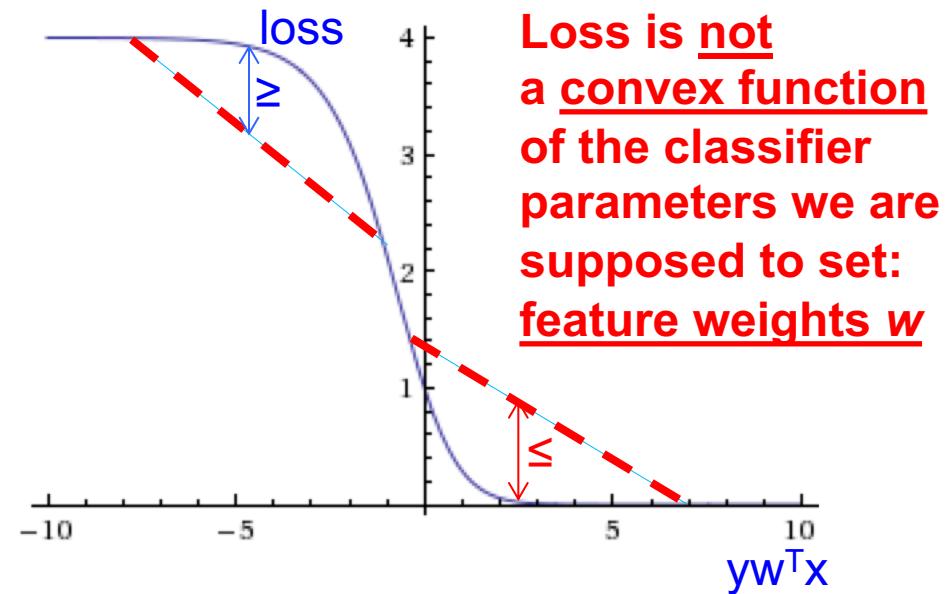
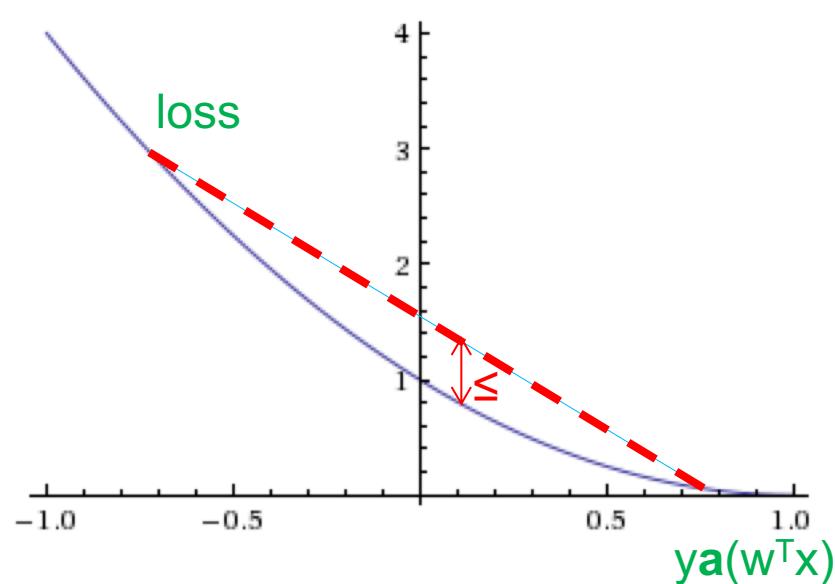
$$\ell(h, z) = (y - a(w^T x)) ^ 2 = (1 - ya(w^T x)) ^ 2$$

- If we modify w for fixed sample (x, y) , how is loss affected?
 - Loss should be expressed in function of w (or $yw^T x$ since y, x are constants coming from the dataset)



Differentiable perceptron

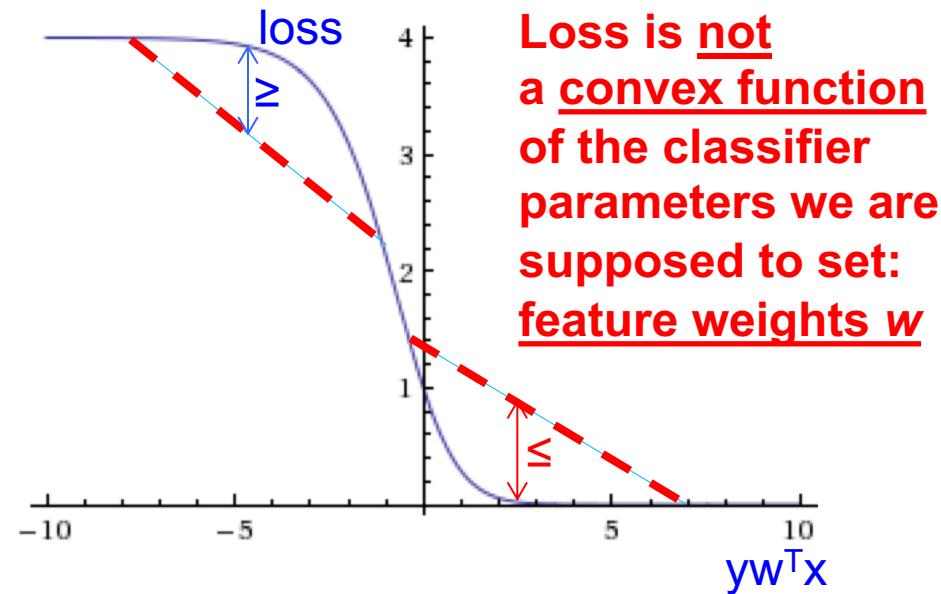
- Looking at loss related to parameters w as *a function of $ya(w^T x)$* vs. *as a function of $yw^T x$*
 - Let's say loss is: $\text{loss}(g(w))$
 - if $g(w)$ is a linear function of w (e.g. $g(w)=yw^T x$) then if $\text{loss}(g)$ is a convex function of g , also $\text{loss}(g(w))$ is a convex function of w
 - if $g(w)$ is not a linear function of w (e.g. $g(w)=ya(w^T x)$ where $a()$ is a sigmoid) then even if $\text{loss}(g)$ is a convex function of g , $\text{loss}(g(w))$ may not be a convex function of w



Loss is not a convex function of the classifier parameters we are supposed to set: feature weights w

Convexity

- Even if the loss is not convex, it still has one minimum (high positive $yw^T x$)
- So where is the problem?

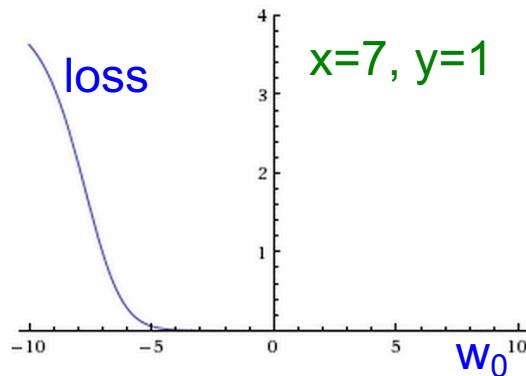


Example

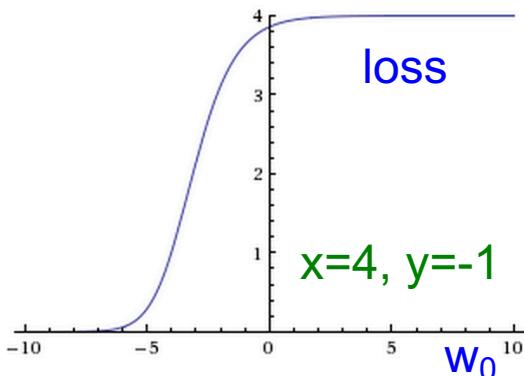
$$\ell(h, \mathbf{z}) = \left(y - a(\mathbf{w}^\dagger T \mathbf{x}^\dagger) \right)^2 \quad a(u) = \frac{2}{1 + e^{-u}} - 1$$

- Let's fix $w=1$, try changing w_0 and observe loss

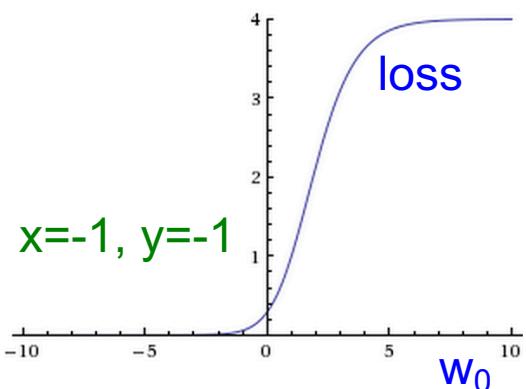
$$\left(1 - \left(\frac{2}{1 + \exp(-7 - w_0)} - 1 \right) \right)^2$$



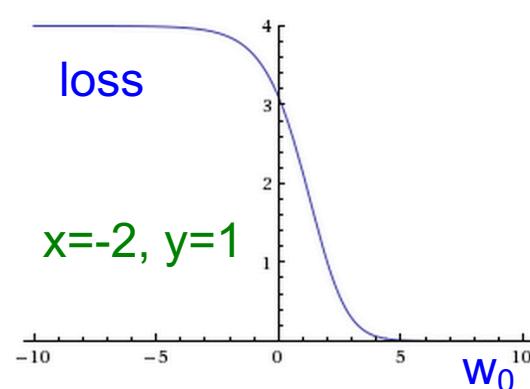
$$\left(-1 - \left(\frac{2}{1 + \exp(-4 - w_0)} - 1 \right) \right)^2$$



$$\left(-1 - \left(\frac{2}{1 + \exp(1 - w_0)} - 1 \right) \right)^2$$



$$\left(1 - \left(\frac{2}{1 + \exp(2 - w_0)} - 1 \right) \right)^2$$



- $h(x) = a(wx + w_0)$
 - $w=1: h(x) = a(x + w_0)$
- We have four samples:
 - $x=7, y=1$
 - $x=4, y=-1$
 - $x=-1, y=-1$
 - $x=-2, y=1$

+ - - +

- What is the plot of *empirical risk* for those four points, as we alter w_0 ?

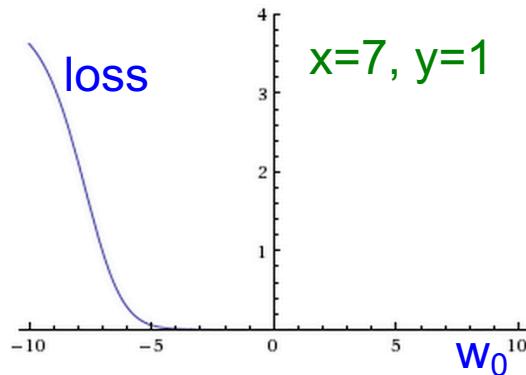
$$\widehat{R}_{S_m}(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, z_i)$$

Example

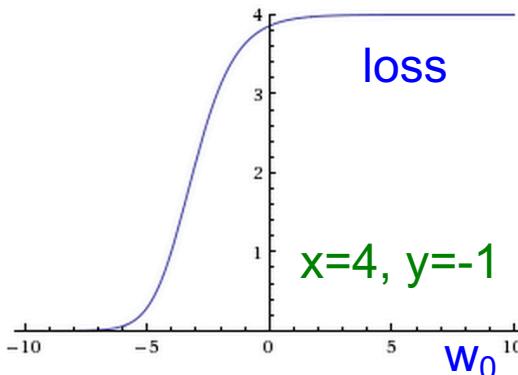
$$\ell(h, \mathbf{z}) = \left(y - a(\mathbf{w}^\dagger T \mathbf{x}^\dagger) \right)^2 \quad a(u) = \frac{2}{1 + e^{-u}} - 1$$

- Let's fix $w=1$, try changing w_0 and observe loss

$$\left(1 - \left(\frac{2}{1 + \exp(-7 - w_0)} - 1 \right) \right)^2$$

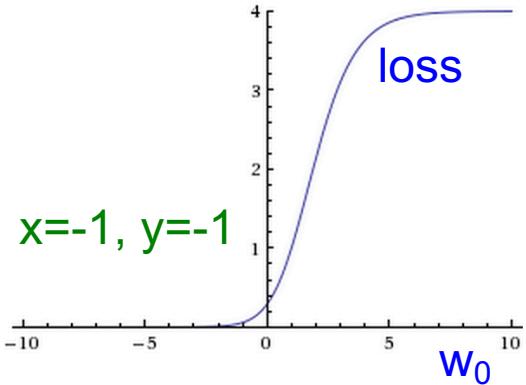


$$\left(-1 - \left(\frac{2}{1 + \exp(-4 - w_0)} - 1 \right) \right)^2$$

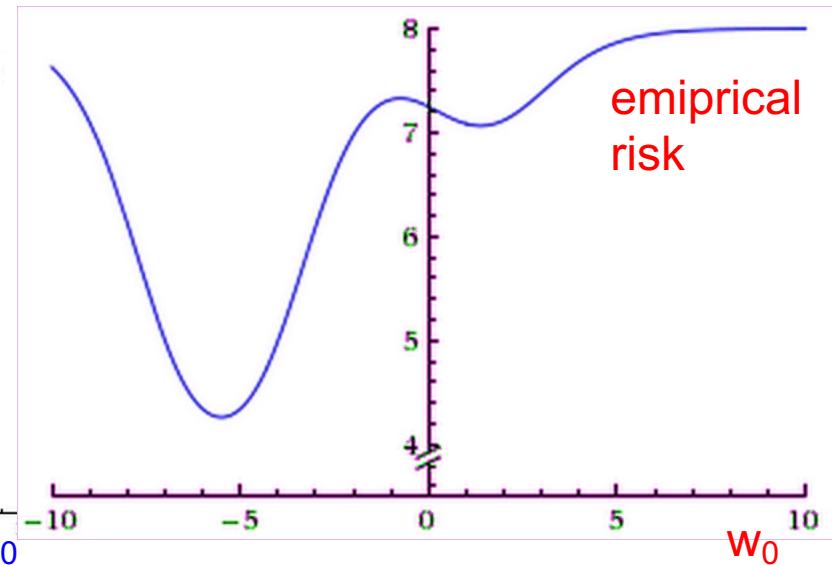
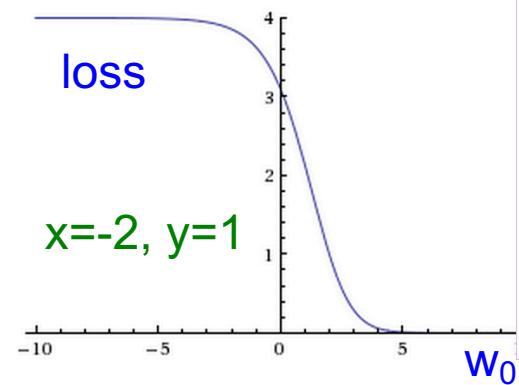


- $h(x) = a(wx + w_0)$
 - $w=1$: $h(x) = a(x + w_0)$
- We have four samples:
 - $x=7, y=1$
 - $x=4, y=-1$
 - $x=-1, y=-1$
 - $x=-2, y=1$

$$\left(-1 - \left(\frac{2}{1 + \exp(1 - w_0)} - 1 \right) \right)^2$$



$$\left(1 - \left(\frac{2}{1 + \exp(2 - w_0)} - 1 \right) \right)^2$$

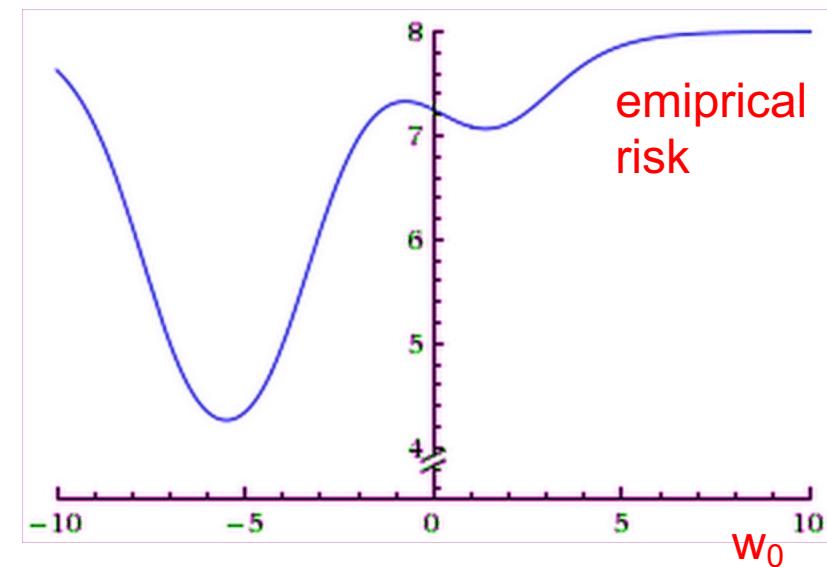


Example

$$\ell(h, \mathbf{z}) = \left(y - a(\mathbf{w}^\dagger T \mathbf{x}^\dagger) \right)^2 \quad a(u) = \frac{2}{1 + e^{-u}} - 1$$

- Let's fix $w=1$, try changing w_0 and observe loss
- Predictions for 4 samples:
 - If we start with $w_0=0$ and apply gradient descent we'll land at $w_0=2$
 - Solution $w_0=2$:
 - $7+2=9$ ok
 - $4+2=6$ wrong
 - $-1+2=1$ wrong
 - $-2+2=0$ wrong-ish
 - Solution $w_0=-5$:
 - $7-5=2$ ok
 - $4-5=-1$ ok
 - $-1-5=-6$ ok
 - $-2-5=-7$ wrong

- $h(x)=a(wx+w_0)$
 - $w=1: h(x)=a(x+w_0)$
- We have four samples:
 - $x=7, y=1$
 - $x=4, y=-1$
 - $x=-1, y=-1$
 - $x=-2, y=1$

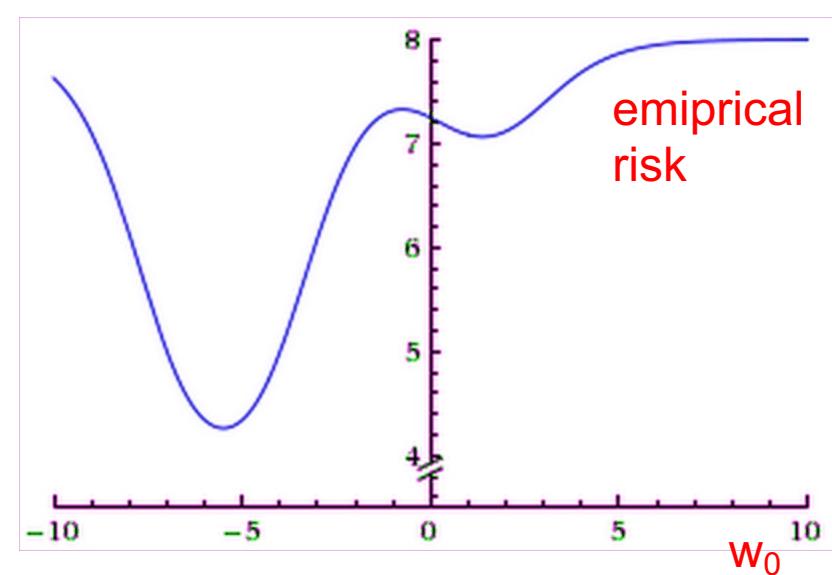
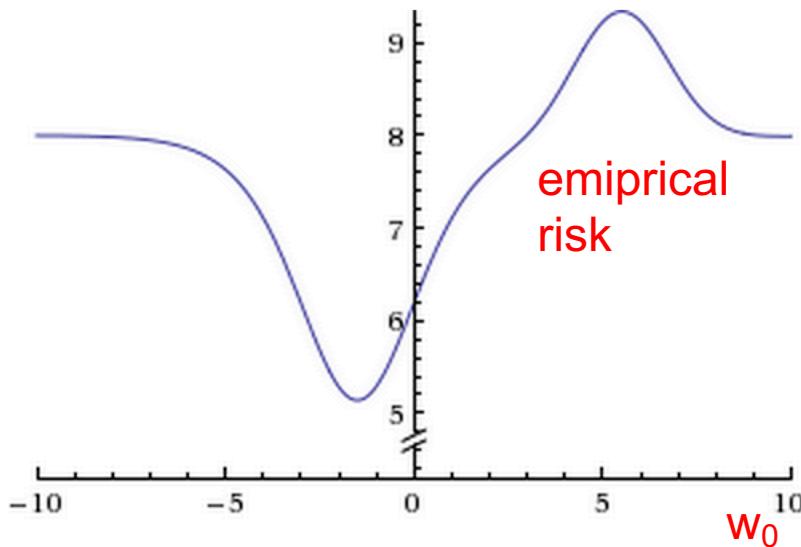


Example

$$\ell(h, \mathbf{z}) = \left(y - a(\mathbf{w}^\dagger T \mathbf{x}^\dagger) \right)^2 \quad a(u) = \frac{2}{1 + e^{-u}} - 1$$

- Let's fix \mathbf{w} , try changing w_0 and observe loss
- What if we fix $\mathbf{w} = -1$?
 - $w_0 = -1.5$ is best
 - again, one wrong prediction
 - $x = -2$ predicted as 1
 - $x = -1, 4, 7$ predicted as -1

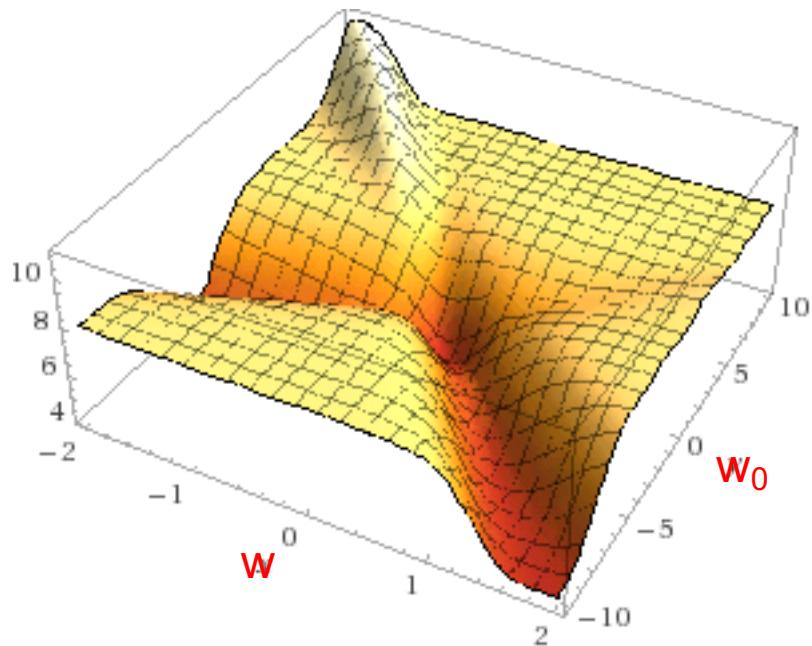
- $h(x) = a(wx + w_0)$
 - $w=1: h(x) = a(x + w_0)$
- We have four samples:
 - $x=7, y=1$
 - $x=4, y=-1$
 - $x=-1, y=-1$
 - $x=-2, y=1$



Example

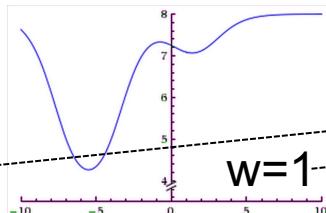
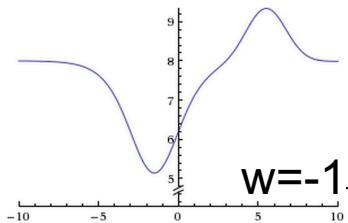
$$\ell(h, \mathbf{z}) = \left(y - a(\mathbf{w}^\dagger T \mathbf{x}^\dagger) \right)^2 \quad a(u) = \frac{2}{1 + e^{-u}} - 1$$

- Empirical risk for any w, w_0

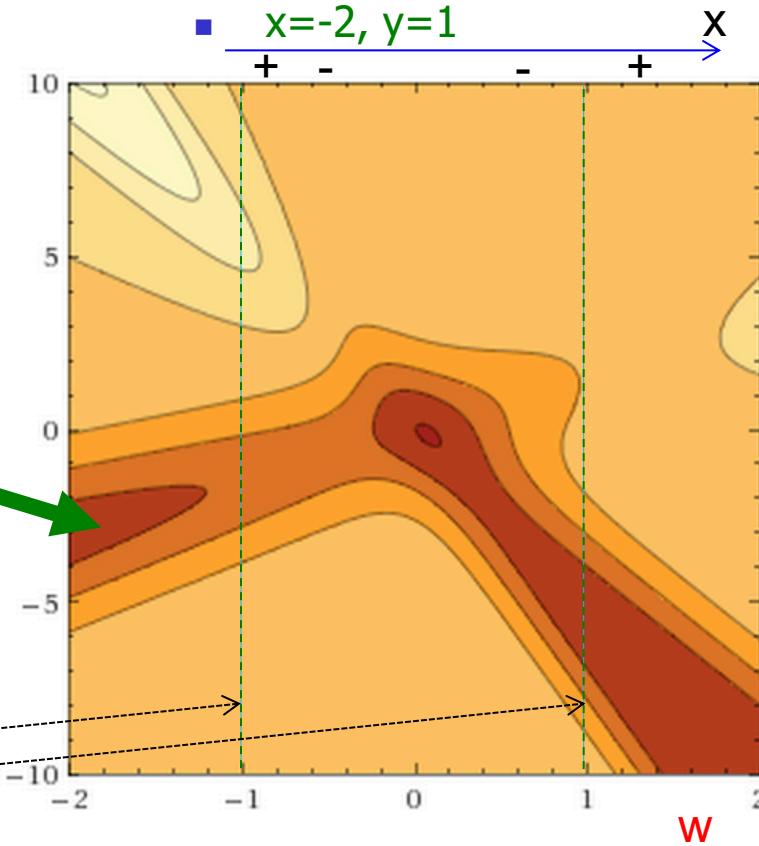


Where is global minimum?

Local minimum



- $h(x) = a(wx + w_0)$
- We have four samples:
 - $x=7, y=1$
 - $x=4, y=-1$
 - $x=-1, y=-1$
 - $x=-2, y=1$



Example

$$\ell(h, \mathbf{z}) = \left(y - a(\mathbf{w}^\dagger T \mathbf{x}^\dagger) \right)^2 \quad a(u) = \frac{2}{1 + e^{-u}} - 1$$

- “No prediction” still the best choice
- Solution ($w=1, w_0=-5.5$):
 - 2 correct-ish prediction => loss ≈ 0.1
 - 1 correct predictions => loss = 0
 - 1 error prediction => loss ≈ 4
- Solution ($w=0, w_0=0$):
 - 4 x “no prediction” => risk = 4

$$h(x) = a(wx + w_0)$$

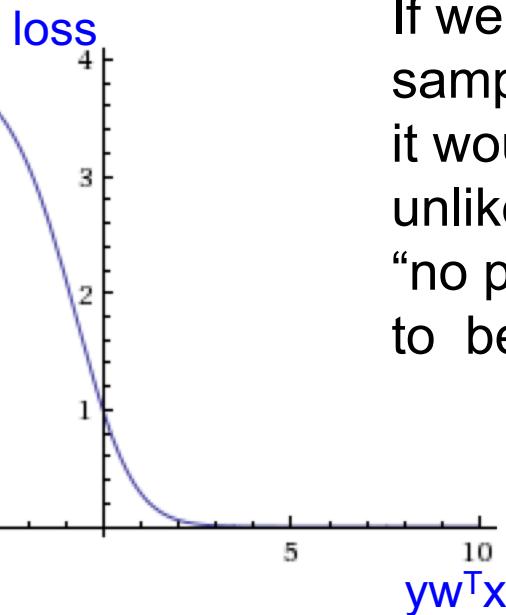
We have four samples:

$$x=7, y=1$$

$$x=4, y=-1$$

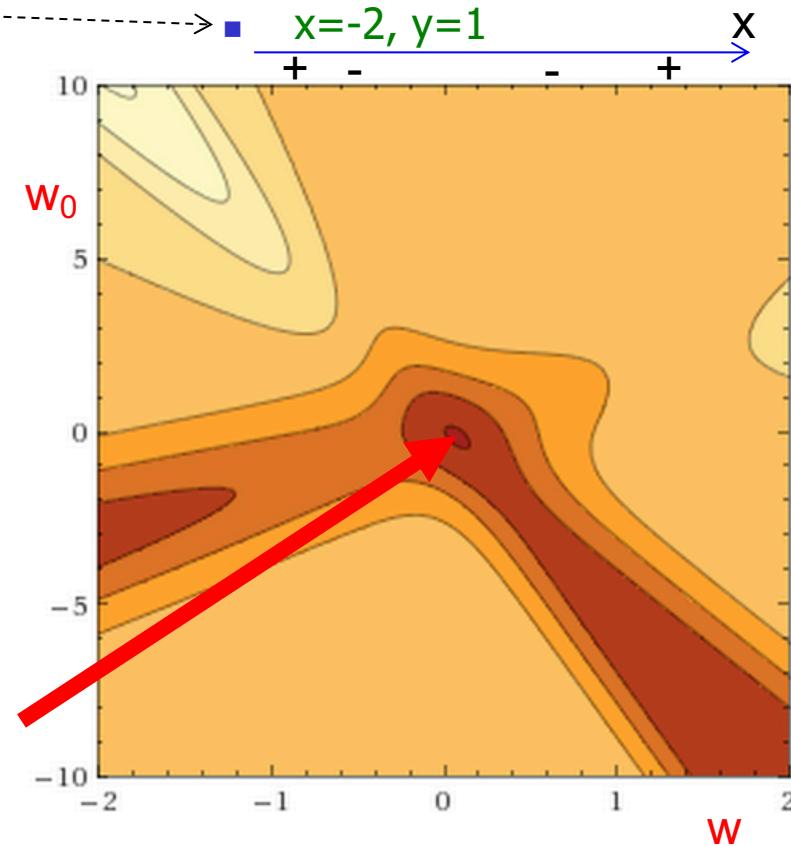
$$x=-1, y=-1$$

$$x=-2, y=1$$



If we had more samples,
it would be
unlikely for
“no prediction”
to be optimum

Global
minimum:
 $w=0, w_0=0$



Example

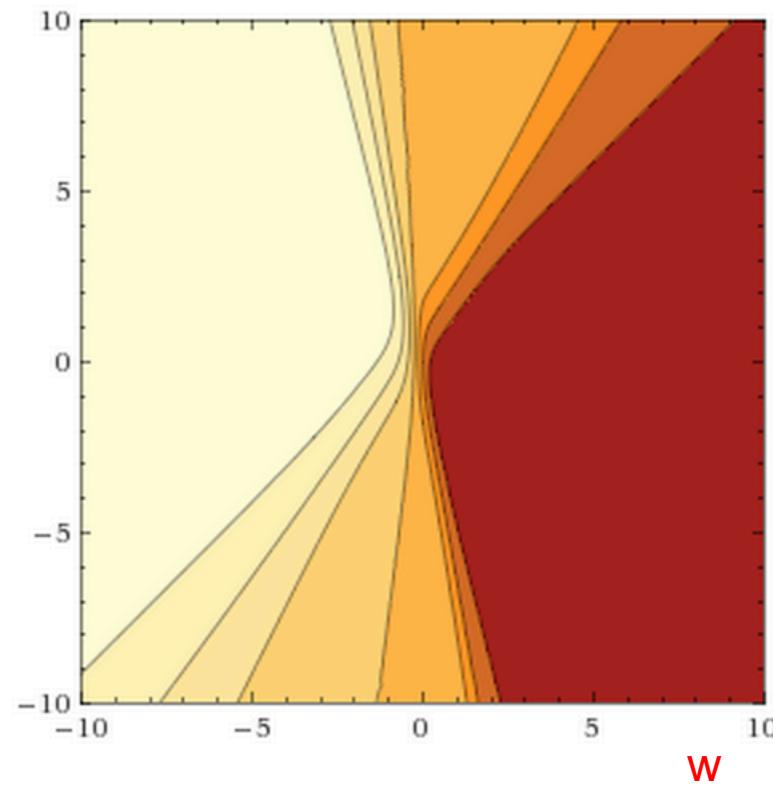
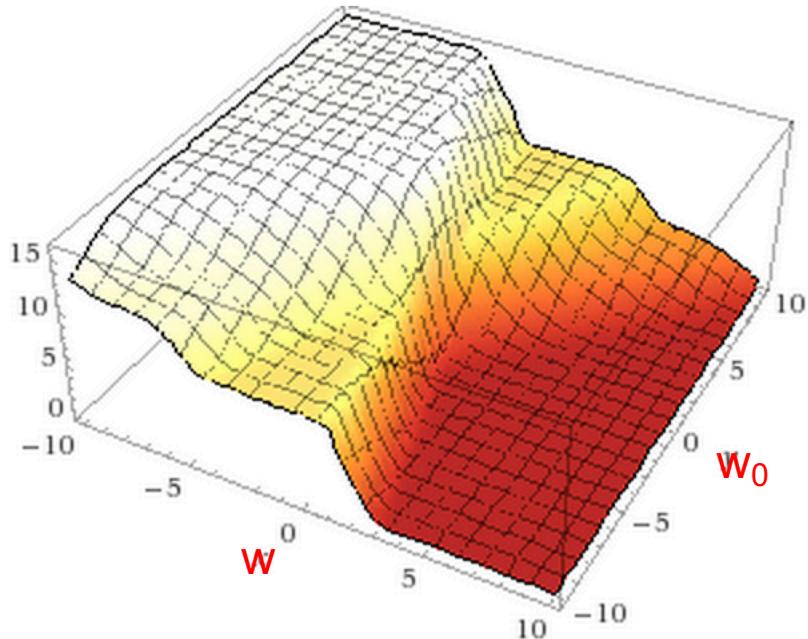
$$\ell(h, \mathbf{z}) = \left(y - a(\mathbf{w}^{\dagger T} \mathbf{x}^{\dagger}) \right)^2 \quad a(u) = \frac{2}{1 + e^{-u}} - 1$$

- Different, easier set of samples

- $x=7, y=1$
- $x=4, y=1$
- $x=-1, y=-1$
- $x=-2, y=-1$

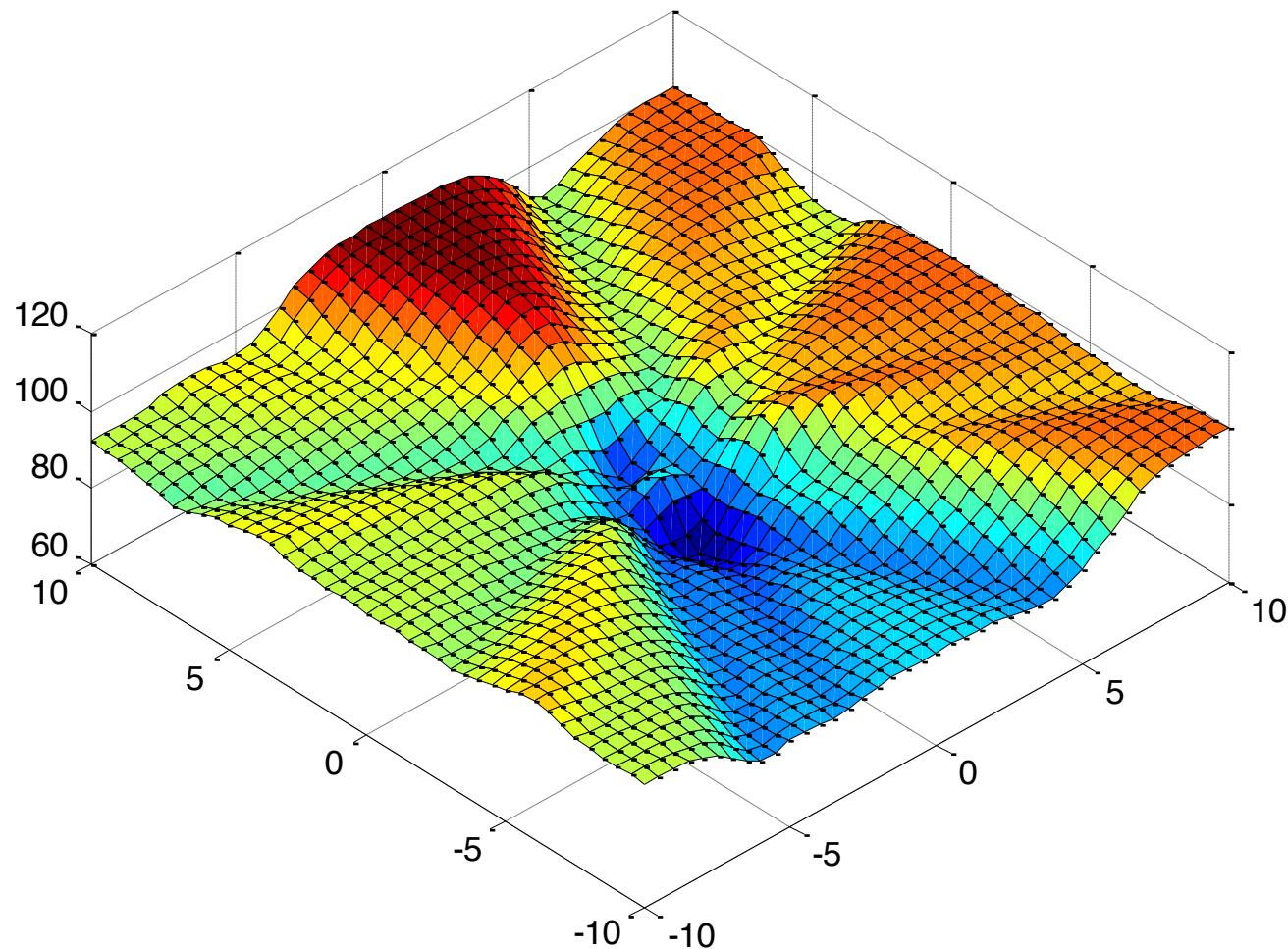


- Empirical risk for any w, w_0



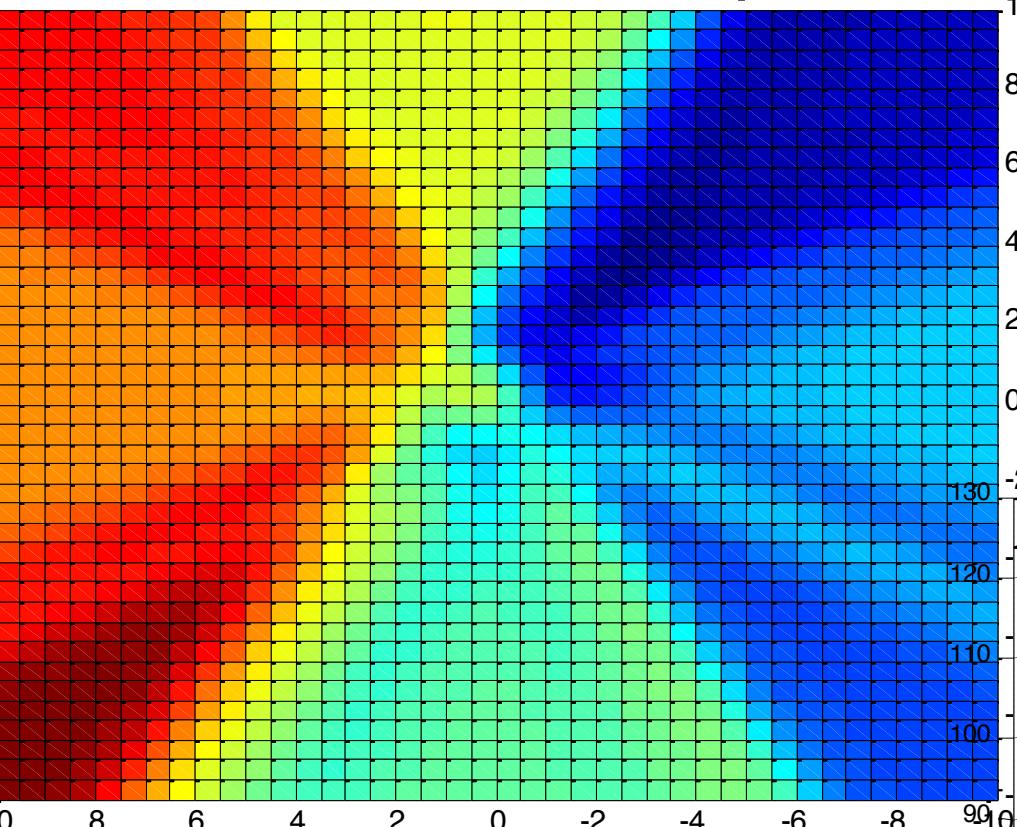
Differentiable perceptron

- Some other examples: empirical risk surface

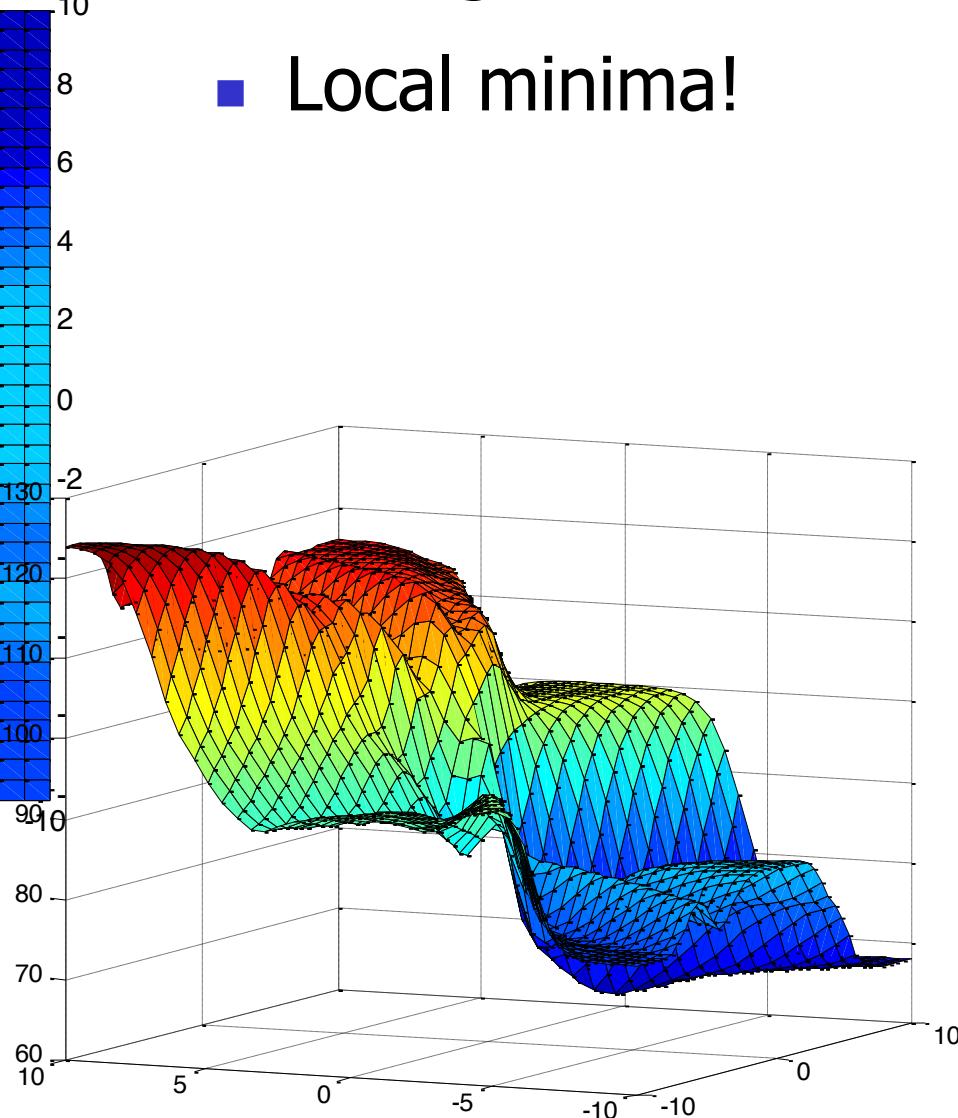


Differentiable perceptron

- Risk surface for quadratic loss over sigmoid

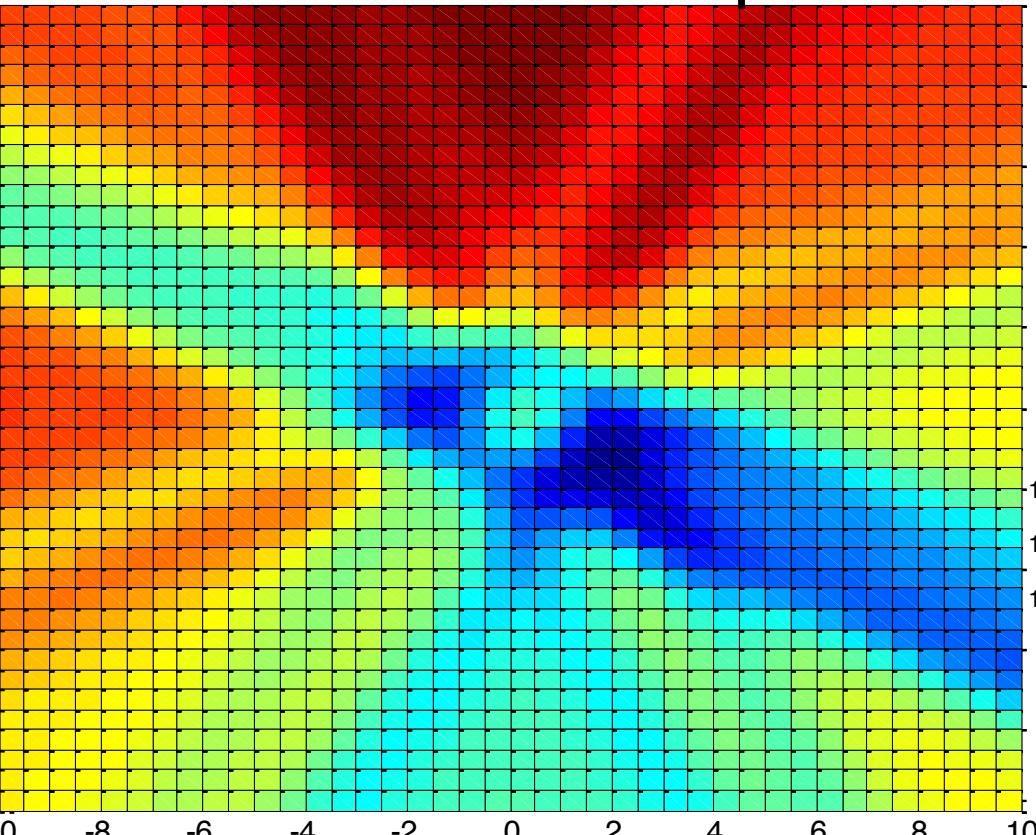


- Local minima!

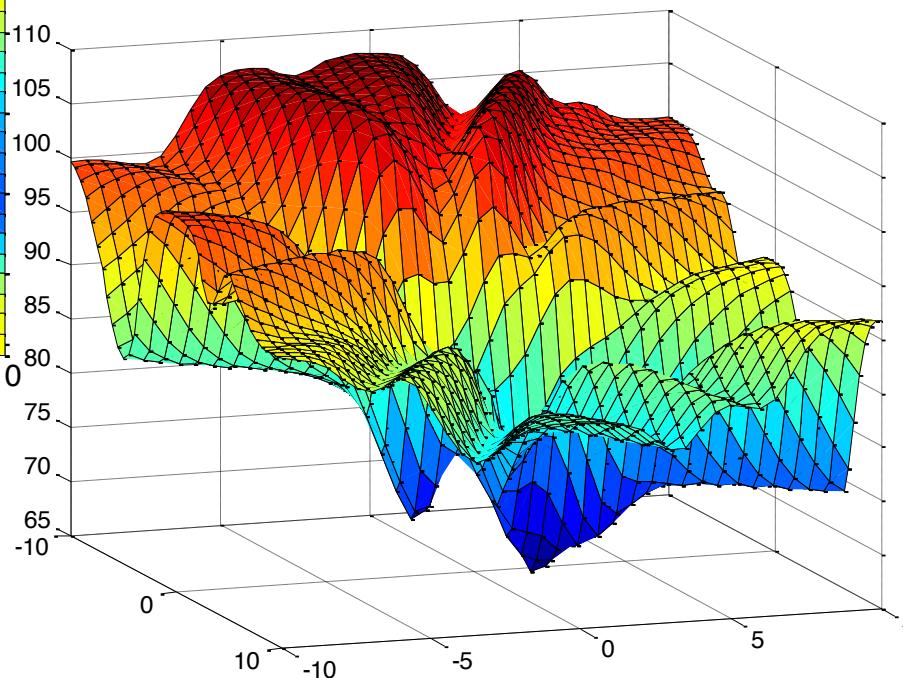


Differentiable perceptron

- Risk surface for quadratic loss over sigmoid

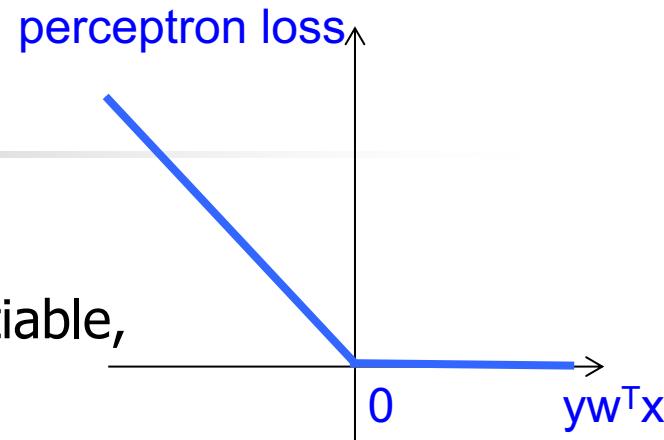


- We may get trapped in a shallow valley: and return a vector \mathbf{w} that's far from the best one!

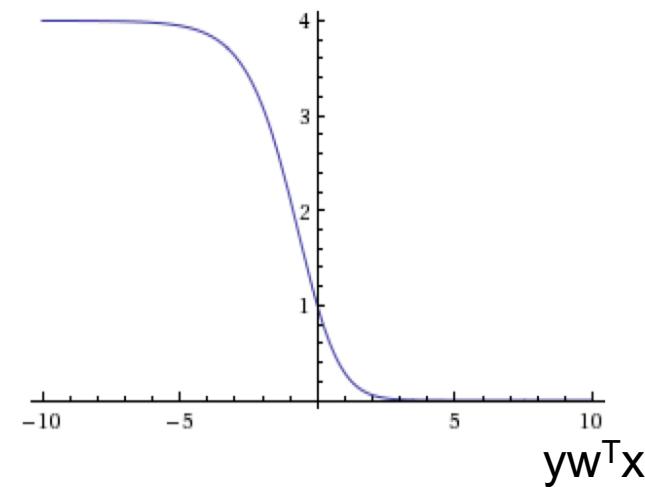
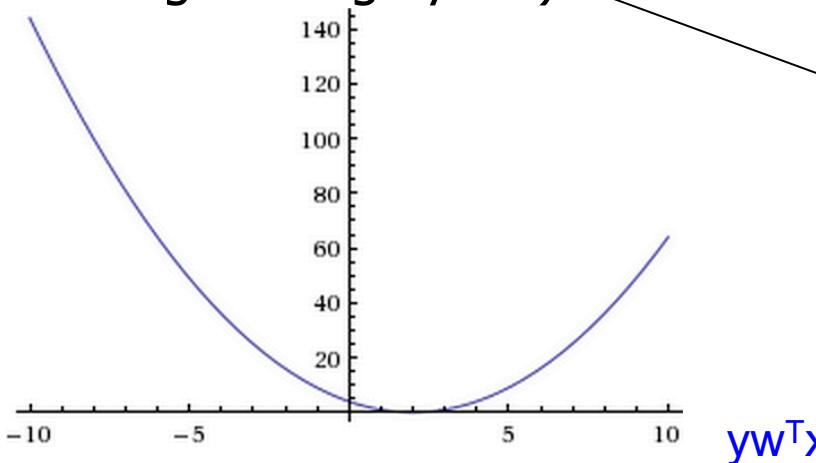


Recap

- Perceptron loss
 - convex, non-increasing, but non-differentiable, and $w=0$ is always a global minimum
- Quadratic loss on sigmoid activation delta rule loss
 - differentiable, non-increasing but non-convex => local minima
- LDA: quadratic loss directly on $yw^T x$
 - convex, differentiable, but non-monotonic
(increasing for large $yw^T x$)



Square loss



This wasn't a problem
for quadratic over sigmoid

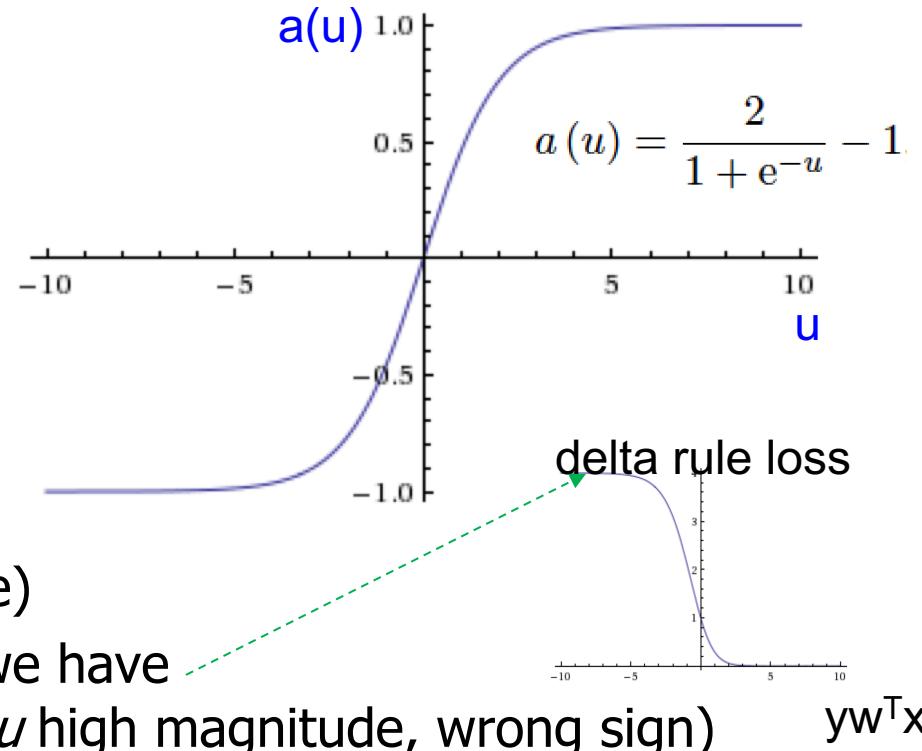
Differentiable perceptron $h(x) = a(w^T x)$

- Another problem with delta rule:

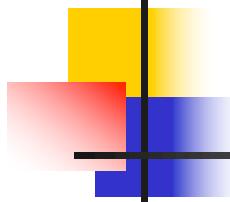
$$w^{\dagger}_{t+1} = w^{\dagger}_t - \frac{c}{2} \frac{\partial \ell(h, z)}{\partial w^{\dagger}} \Big|_{w^{\dagger}=w^{\dagger}_t} = w^{\dagger}_t + c \left[y - a(w_t^T x^{\dagger}) \right] a'(w_t^T x^{\dagger}) x^{\dagger}.$$

- The update depends on:

- $(y - a(u))$
 - Smaller as we approach correct prediction
- $a'(u)$
 - Much smaller as we approach correct prediction (u v.large)
 - And also very small when we have really incorrect prediction (u high magnitude, wrong sign)



- Very slow learning for large $|u|$



For next time

Please review:

Expected value

$$\mathbb{E}[X] = \sum_{i=1}^k x_i p_i = x_1 p_1 + x_2 p_2 + \cdots + x_k p_k.$$