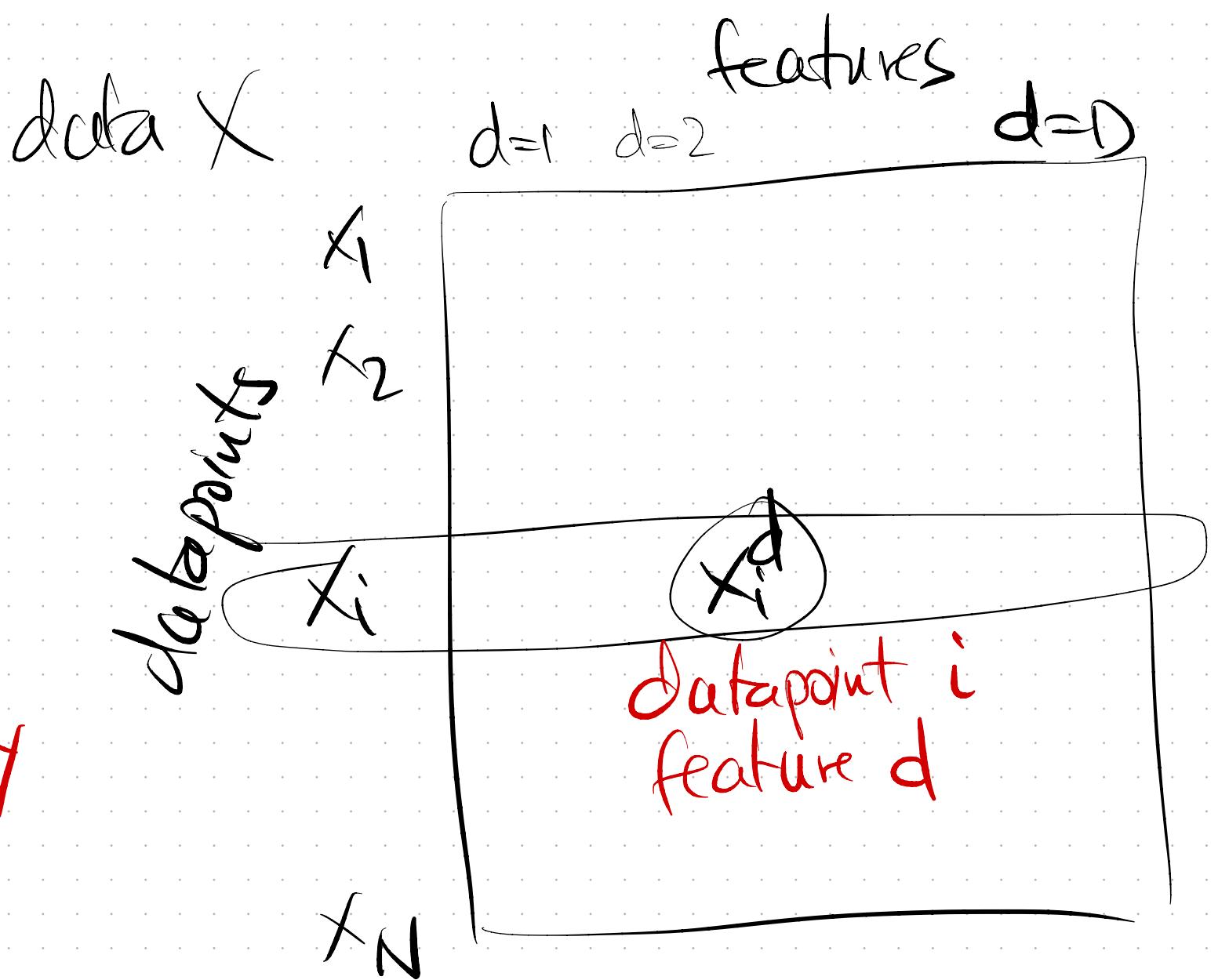


Lecture 5/21 - 23

- Gradient Descent
- Linear Regression with GD
- Logistic Regression
- HW2A (released later) due 6/3
 today
- ROC curve \approx AUC geometrically
- Perceptron w/GD or geometry



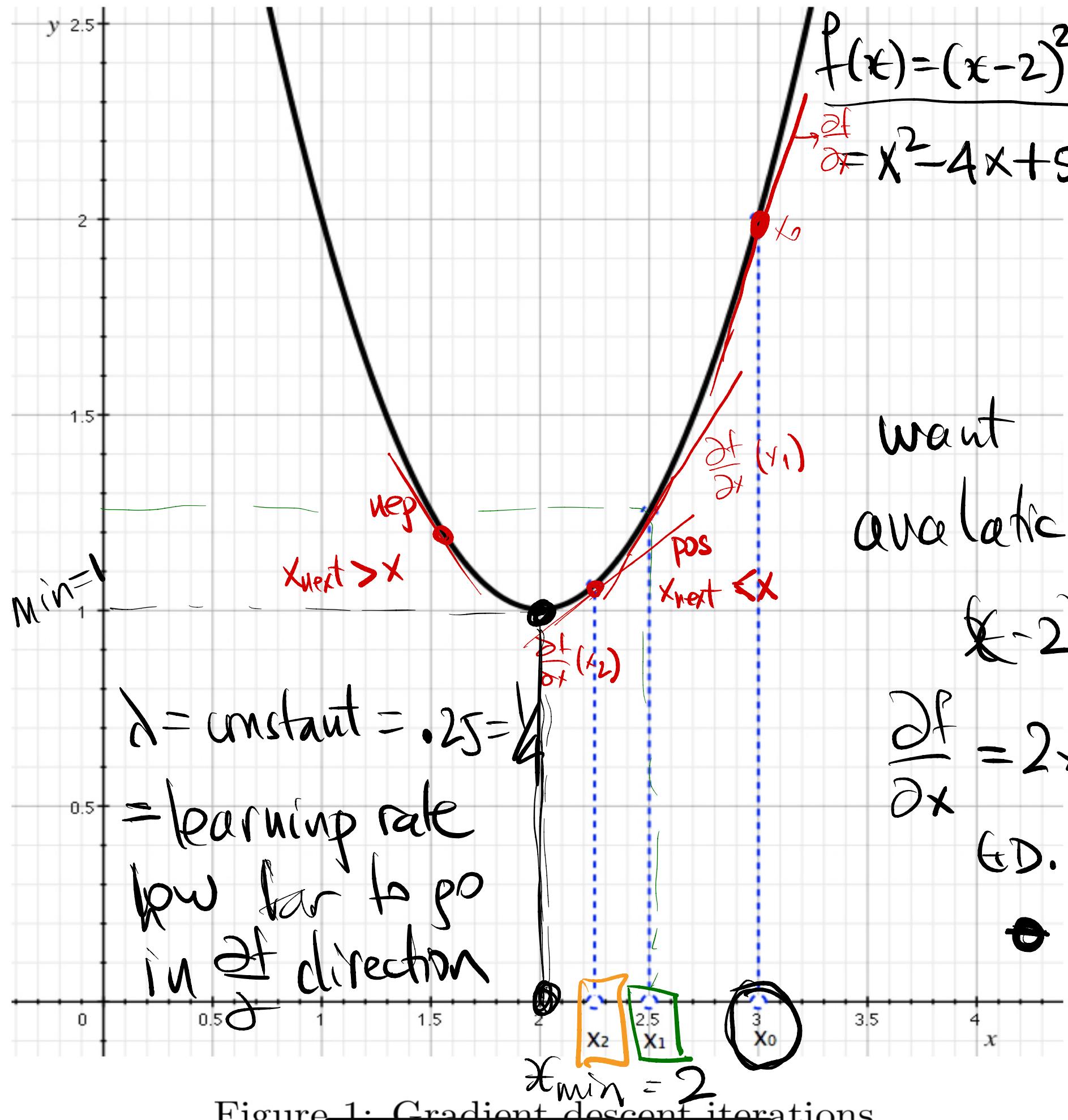


Figure 1: Gradient descent iterations

(+D): iteratively find
 by moving x in direction
 of gradient = differential.

want $\min_x f(x)$ $\arg\min = x_{min}$
 analytic solution $x_{min} = 2$
 $(x-2)^2 + 1 \geq 1$ equality if $x = 2$

$\frac{\partial f}{\partial x} = 2x - 4$ gradient. Start at $x_0 = 3$

(+D. iterations:

$$x_1 = x_0 - \lambda \frac{\partial f}{\partial x}(x_0)$$
 $= 3 - \frac{1}{4}(2 \cdot 3 - 4) = 2.5$

$$x_2 = x_1 - \lambda \frac{\partial f}{\partial x}(x_1) = 2.5 - \frac{1}{4} (3 \times 2.5 - 4) = \boxed{2.25}$$

i want $(x_1, x_2, x_3, \dots) \rightarrow x_m = 2$

- What if I swap bofar on the other side? $x_m < 2$

$$\frac{\partial f}{\partial x}(x_m) = 2 \cdot x_m - 4 < 0 \text{ negative}$$

$$\text{update } x_{m+1} = x_m - \frac{\partial f}{\partial x}(x_m) \text{ increasing } x_{m+1} > x_m$$

its moving x_{next} in right direction

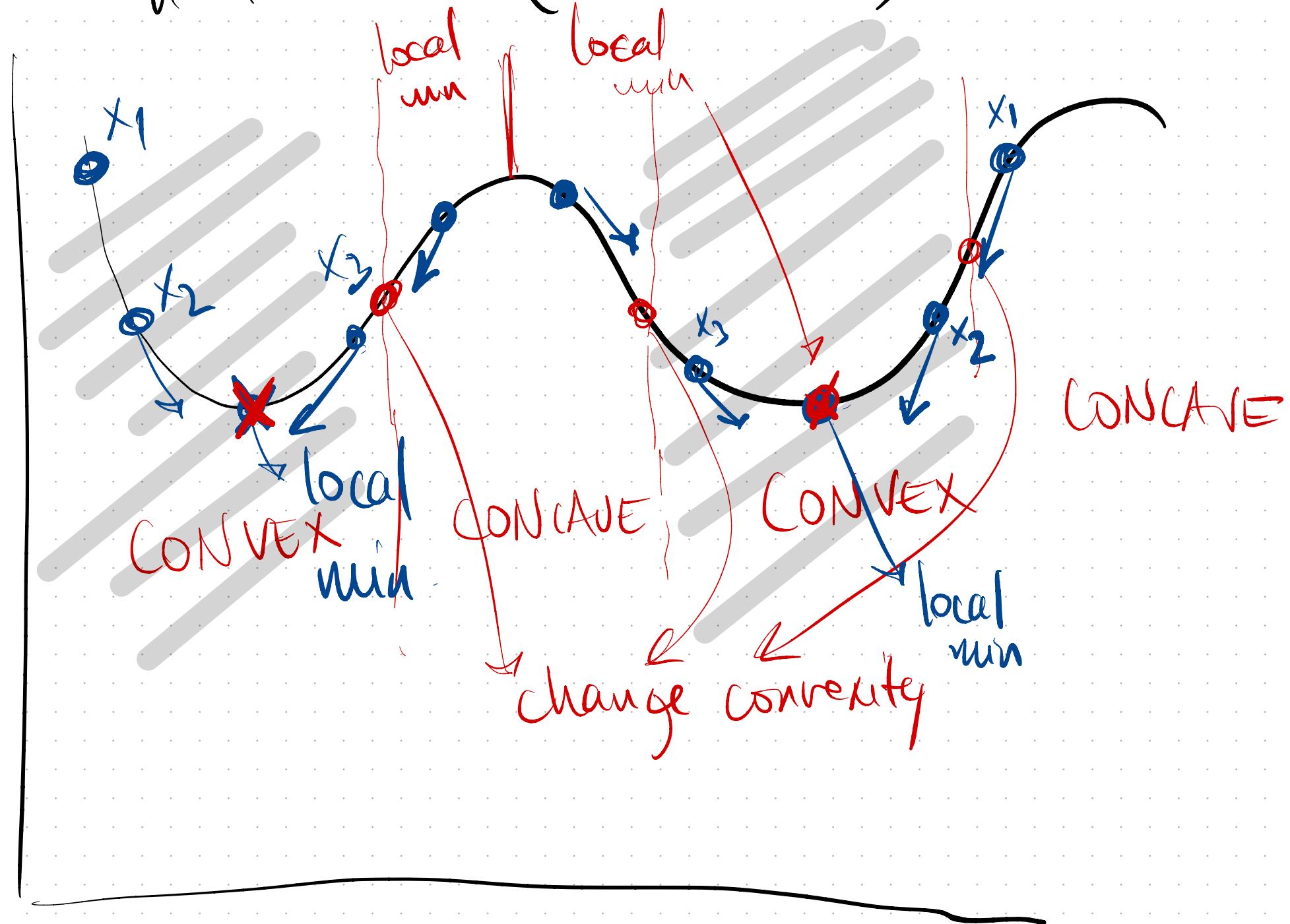
- What if $x_m = x_{\text{min}}$?

$$\boxed{\frac{\partial f}{\partial x}(x_{\text{min}}) \geq 0}$$

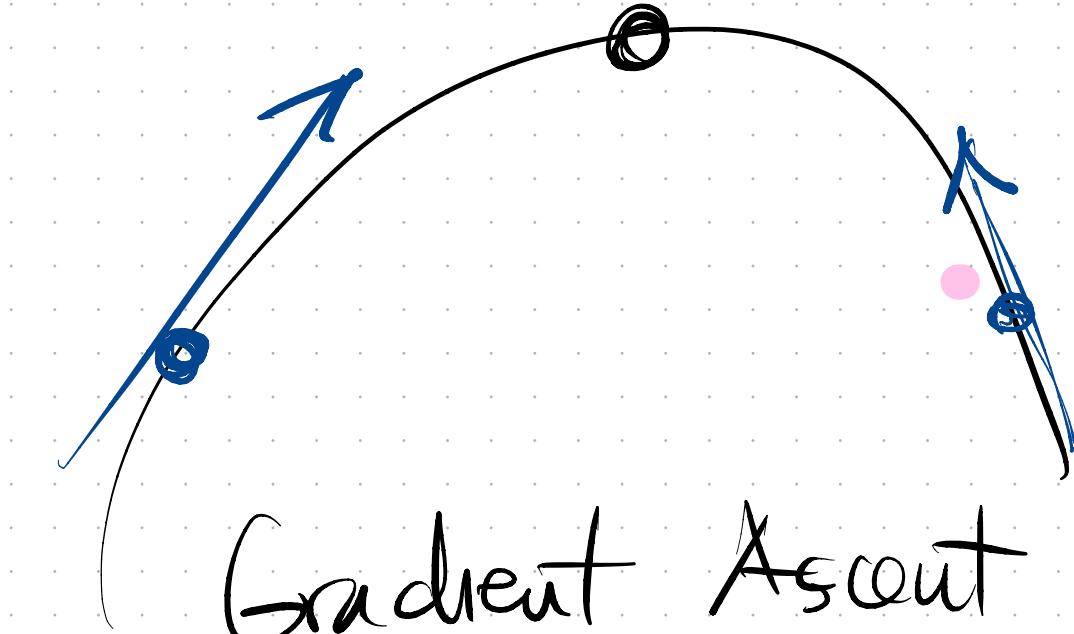
$$\Rightarrow \text{update } x_{m+1} = x_m - \boxed{\frac{\partial f}{\partial x}(x_m)} \\ = \text{save } x_m (\text{no change})$$

GD stays here.

Non convex (all the time) BUMPY



CONCAVE Region



Gradient Ascent

$$x_{\text{new}} = x + \lambda \frac{\partial F}{\partial x}(x)$$

λ = trial and error

- too low : many steps to converge

- too high : overshoot (other side) DIVERGE

Regression w/ GD. Error = Obj = $J(w) = \frac{1}{2} \sum_{i=1}^N (x_i w - y_i)^2$

$$= \frac{1}{2} \sum_{i=1}^N \left(\sum_{d=1}^D x_i^d \cdot w_d - y_i \right)^2$$

$$\frac{\partial J}{\partial w_d} = \frac{1}{2} \cdot 2(x_w - y) \frac{\partial (x_w - y)}{\partial w_d} = (x_w - y) \boxed{\frac{\partial (\sum_d x_i^d w_d - y)}{\partial w_d}} = (x_w - y) x_d$$

w
comp
d

GD update rule
(iterative)

$$w_{\text{new}}^d = w^d - \lambda [(x_w - y) x_d]$$

$$\text{for all } d \text{ and one datapoint } x_i \\ w_{\text{new}} = w - \lambda [(\cdot x_w - y) x_i]$$

vector vector

batch: keep $w = \text{fixed}$ (From previous iteration)

$$\text{update } w_{\text{new}} = w - \lambda \sum_{i=1}^N [(x_i w - y) x_i] \text{ for all } i=1:N$$

More stable
batch size \sqrt{N}

$$\Rightarrow w_{\text{new}} = w - \lambda \sum_{i=1}^N [(x_i w - y_i) x_i] \quad \text{(prev)}$$

sum of GD
updates for all
datapoints

• Stochastic: change w (current) with every datap update

faster

convergence
batch of size 1

for $i = 1 : N$

$$w_{\text{new}} = w \rightarrow [(x_i w - y_i) x_i]$$

$w = w_{\text{new}}$ → before moving to next datapoint x_{i+1}

• Batch: batch of size $T = 20$

batch (Keep w fixed for T datapoints)
prev

update $w = \dots$

move to next batch with updated w

TERMINATION

• $\frac{\partial J}{\partial w}$ close to zero

• $|w_{\text{new}} - w|$ very low

• Job good enough
(not for regression)

• performance (acc)
good enough

$$\text{Reg + } L_2 \text{ regularization } J(w) = \frac{1}{2} \sum_{i=1}^N (x_i w - y_i)^2 + \frac{\lambda}{2} \|w\|^2 \text{ 'Ridge'}$$

$$\frac{\partial J}{\partial w^d} = -(x_w - y) x^d + L_2 w^d$$

GD update $w^d_{\text{new}} = w^d - \lambda [-(x_w - y) x^d + L_2 w^d]$

for datap x_i

- High L_2 penalty $\Rightarrow \|w\|$ small (constrained)
 $J(w)$ not nec. optimal

- Low L_2 penalty $\Rightarrow J(w)$ minimize (optimal)
 $|w|$ large possibly.

$$L_1\text{-regularization } J(w) = \frac{1}{2} \sum_{i=1}^N (x_i w - y_i)^2 + L_1 \cdot |w|$$

"Lasso"

$$\text{both } L_1 + L_2 : \quad J(w) = \frac{1}{2} \sum_{i=1}^N (x_i w - y_i)^2 + L_1 |w| + \frac{1}{2} L_2 \|w\|^2$$

"Elastic Net"

derivation L_1 -reg

$$\frac{\partial J}{\partial w} = ? \text{ difficult}$$

Sparse
some $w_d \approx 0$

→ spread w over
features

Get toy L_1 -reg scikit-learn
call regression

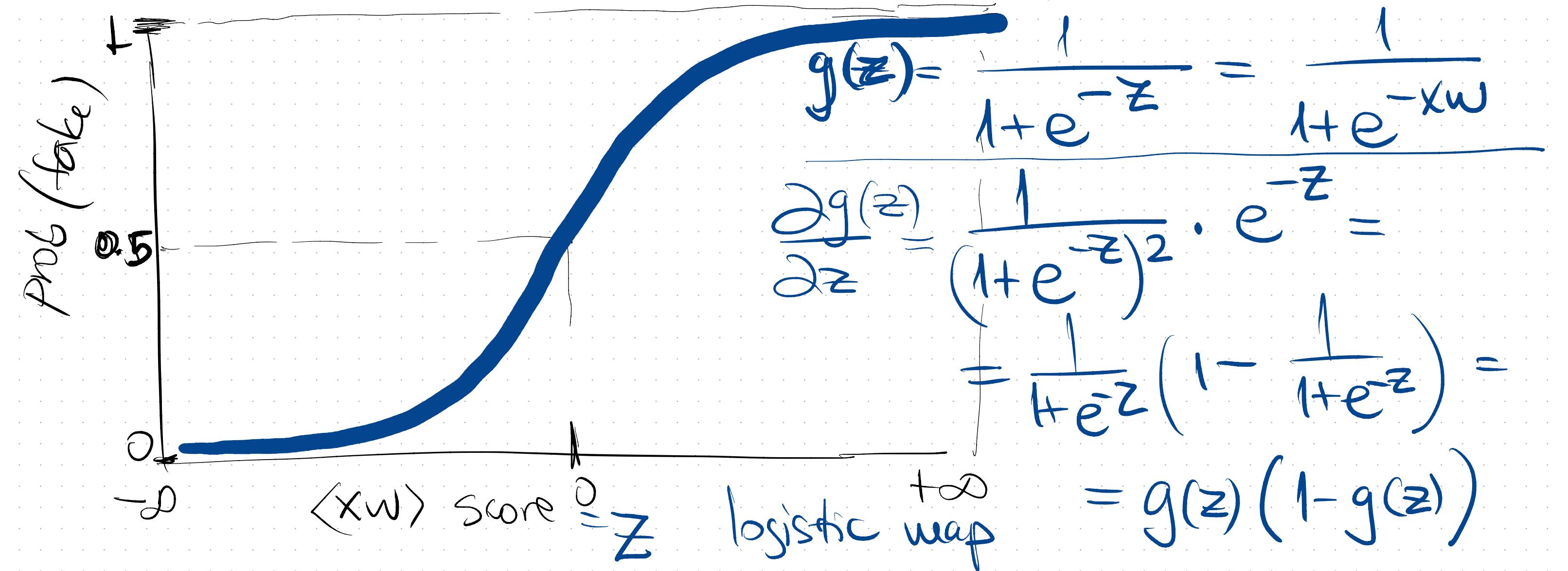
Logistic Regression : deal with classification classes $y \xrightarrow{\text{no}} \text{yes}$
 e.g. Spambase dataset

Lin. Regression $h(x) = xw \cong y = \text{quantity label (ex. house price)}$

Idea:

- xw = linear regression score
- $h(x) = \boxed{\text{map}}$ xw score \rightarrow probability of $y = \text{yes}$

logistic map $||g||$
 $Z = \text{Score} = xw$



- Key property: overshooting does not penalize OJ

Lin Regression

want $xw \approx y$

$$Y = \begin{cases} 0 \\ 1 \end{cases}$$

Logistic Regression

$$\text{want } \frac{1}{1+e^{-xw}} \approx Y$$

If $xw < 0$: error $(xw - 0)^2$

: xw very small $\rightarrow -\infty$
 $\frac{1}{1+e^0} \approx 0$

if $xw \gg 1$ error $(xw - 1)^2$

xw very large $\rightarrow +\infty$

$$\frac{1}{1+e^\infty} \approx 1$$

logistic Regression uses log-likelihood OBJECTIVE (not sq error)

WANT $h(x_i) = g(x_i w) = \frac{1}{1+e^{-x_i w}} \approx \text{prob}(y_i = 1)$ 2 classes

$$\begin{aligned} P(y_i = 1) &\approx h(x_i) \\ P(y_i = 0) &\approx 1 - h(x_i) \end{aligned} \Rightarrow P(y_i | x_i) = h(x_i) \cdot (1 - h(x_i))$$

$$h(x_i) \cdot (1 - h(x_i)) = \begin{cases} h(x_i) & \text{if } y_i = 1 \\ 1 - h(x_i) & \text{if } y_i = 0 \end{cases}$$

if $y_i = 1$ (yes) \Rightarrow want $h(x_i) = \text{high}$; $(1 - h(x_i)) = \text{low}$

$y_i = 0$ (no) \Rightarrow want $h(x_i) = \text{low}$; $(1 - h(x_i)) = \text{high}$

log likelihood

"LL"

$$\text{LL}(w) = \log \left[\prod_{i=1}^N P(y_i | x_i) \right]$$

want MAX LL(w)

product of prod of
correct prediction

datapoints indep. of each
other.

$$\begin{aligned}
 &= \log \left(\prod_{i=1}^N h(x_i)^{y_i} (1-h(x_i))^{1-y_i} \right) \\
 &= \sum_{i=1}^N \log(h(x_i)^{y_i}) + \log((1-h(x_i))^{1-y_i}) \\
 &= \sum_{i=1}^N [y_i \log(h(x_i)) + (1-y_i) \log(1-h(x_i))]
 \end{aligned}$$

gradient

$$\frac{\partial L_L}{\partial w} =$$

• GD update $w_{\text{new}} = w + \lambda \frac{\partial L_L}{\partial w}(w)$ ("ascent")