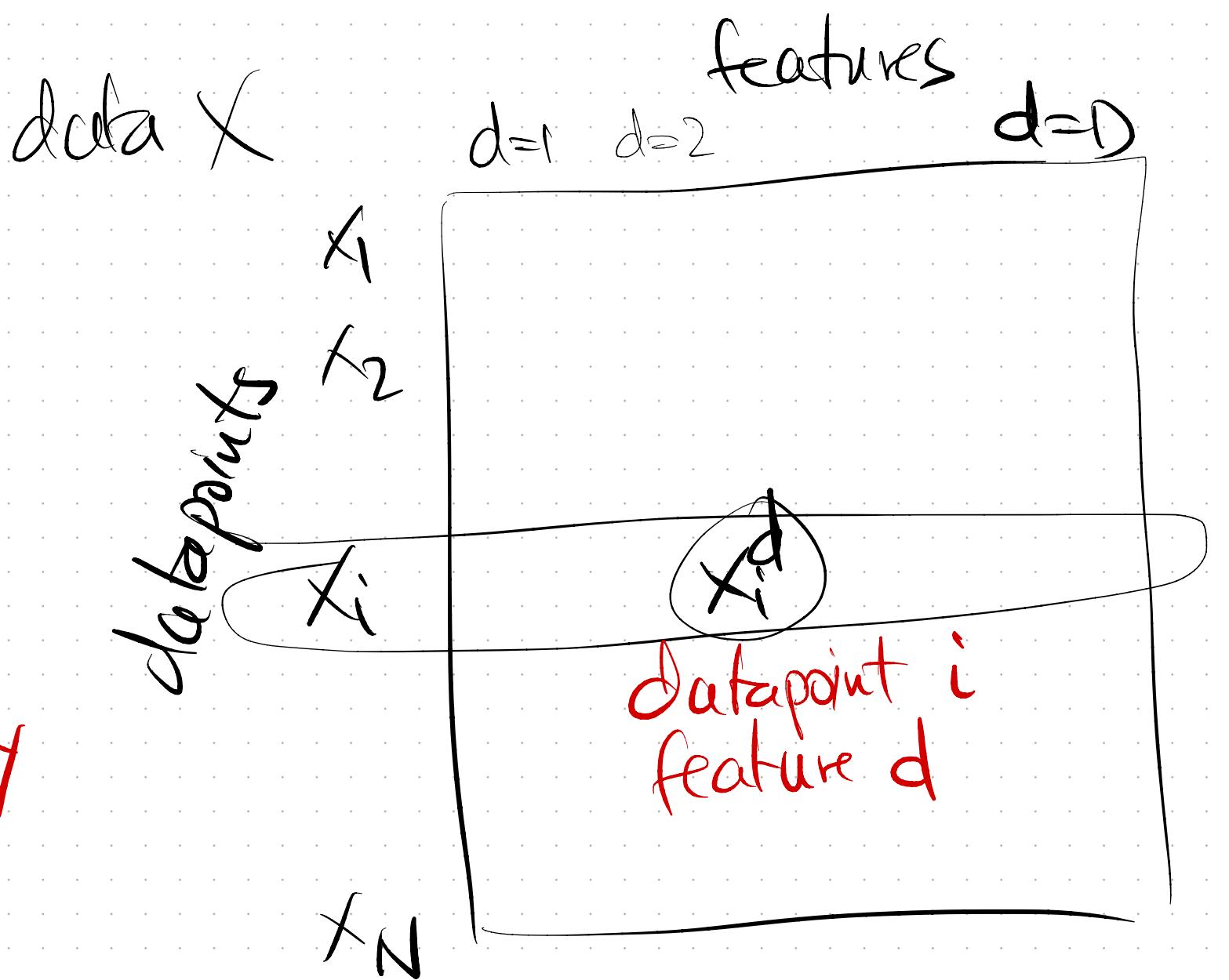


Lecture 5/21 - 23

- 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



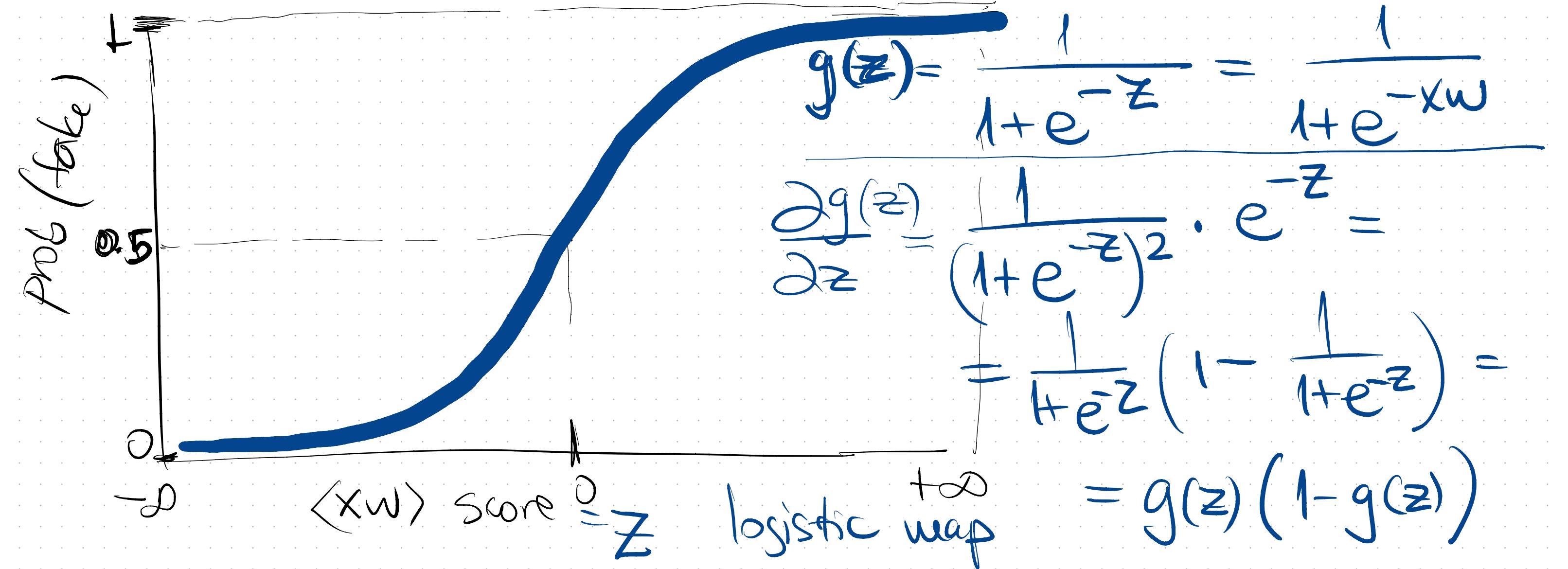
Logistic Regression : deal with classification classes $y \xrightarrow{\text{no}} \text{yes}$
 e.g. Spambase dataset

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

Idea:

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

probs (label = yes)
 logistic map $||g||$
 $Z = \text{Score} = xw$



- Key property: overshooting does not penalize OJ

Lin Regression

want $xw \approx y$

$$Y = \underbrace{1}_{\text{want}} + \underbrace{\omega}_x$$

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 (coef w)

$$h(x) = \text{logistic}(\text{Regscore}(x)) = \text{logistic}(xw)$$

$$= \frac{1}{1+e^{-xw}} \underset{\text{interpret}}{\approx} \underset{\text{probab } (y \rightarrow \text{class yes})}{\approx}$$

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

$P(y_i = 1) \approx h(x_i)$
 $P(y_i = 0) \approx 1 - h(x_i)$ $\Rightarrow P(y_i | x_i) = h(x_i) \cdot (1 - h(x_i))$

y_i 0/1 filters
correct pred

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

Chance of being correct
if $y_i = 1$ (yes) \Rightarrow want $h(x_i)$ = high ; $(1 - h(x_i))$ low
 $y_i = 0$ (no) \Rightarrow want $h(x_i)$ = low ; $(1 - h(x_i))$ high

log likelihood "LL" $\log [\prod_{i=1}^N P(y_i | x_i)]$ product of prod of correct prediction

Want MAX LL(w)
all datapoints indep. of each other.

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

all train
data \leftarrow gradient

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

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

$$\begin{aligned}
 \log(a \cdot b \cdot c) &= \log(a) + \\
 &\quad \log(b) + \log(c) \\
 \log(a^n) &= n \log(a)
 \end{aligned}$$

$$L = \sum_{i=1}^n [y_i \log(h(x_i)) + (1-y_i) \log(1-h(x_i))] = \sum_{i=1}^n [y_i \log(g(x_iw)) + (1-y_i) \log(1-g(x_iw))]$$

one datapoint x_i :

$$\frac{\partial L}{\partial w_d} = \left[y_i \frac{1}{g(x_iw)} - (1-y_i) \frac{1}{1-g(x_iw)} \right] \frac{\partial}{\partial w_d} g(x_iw)$$

$$= \left[y_i \frac{1}{g(x_iw)} - (1-y_i) \frac{1}{1-g(x_iw)} \right] \cdot \frac{\partial g}{\partial (x_iw)} \cdot \frac{\partial (x_iw)}{\partial w_d}$$

$$= \left[y_i \frac{1}{g(x_iw)} - (1-y_i) \frac{1}{1-g(x_iw)} \right] \cdot [g(x_iw)(1-g(x_iw))] \frac{\partial}{\partial w_d}$$

$$= [y_i(1-g(x_iw)) - (1-y_i) \cdot g(x_iw)] \frac{\partial}{\partial w_d}$$

~~$$= [y_i - y_i g(x_iw) - g(x_iw) + y_i g(x_iw)] \frac{\partial}{\partial w_d}$$~~

$$= [y - g(x_iw)] x_i^d = [y - h(x_i)] x_i^d$$

rate of compound
differentials

GD update for datap. x_i : $w_d^{\text{new}} = w_d + \lambda(y_i - h(x_i))x_i^d$

GD update for all datap.: $w_d^{\text{new}} = w_d + \lambda \sum_{T=1}^N (y_i - h(x_i))x_i^d$

GD update all $\text{amps}(d)$
all $\text{datap}(i)$ $w^{\text{new}} = w + \lambda [y - h(x)]x$

Regularization Logistic Regression

$$OBJ = \text{LogLikelihood}(\omega) + L_2 \sum_{d=1}^D (\omega_d)^2$$

easy

$$\sum_{i=1}^N h(x_i)y_i (1-h(x_i))^{(y_i)}$$

$$OBJ = \text{log likelihood} + L_1 \|\omega\|$$

$$L_1 \sum_{d=1}^D |\omega_d| \rightarrow \text{not easy.}$$

ML training:

- $\frac{\partial OBJ}{\partial \omega} = \nabla_{\omega} \text{ (gradient)}$

- GD update $\omega_{\text{new}} = \omega + \lambda \cdot \nabla_{\omega}$

Multiple classes (multinomial) $y \in \{1, 2, 3, \dots, K\}$ K classes.

- train a linear score (lin regression) for each class

$$f_K(x) = \langle x \cdot w_K \rangle = \sum_{d=1}^D x^d \cdot w_K^d \quad w_K = \text{vector of coef for class } K.$$

binary \rightarrow k classes

• logistic \rightarrow Softmax

$$\text{prob}(Y=y_K) = \frac{e^{f_K(x)}}{\text{normalized}} = \frac{\exp(f_K(x))}{\sum_{t=1}^k \exp(f_t(x))}$$

↓
distribution over possibilities

$$Y = y_1 y_2 \dots y_K$$

Newton Update instead of GD

unidim:

$$w_{\text{new}} = w - \frac{\nabla L(w)}{\nabla^2 L(w)}$$

$\nabla L(w)$

$\nabla^2 L(w)$

1st diff

w
dif. (twice differential)

$$LL = \text{obj}(w)$$

$$\nabla LL = \text{gradient} \quad \frac{\partial LL}{\partial w}$$

$$\nabla^2 LL = \text{second grad} \quad \frac{\partial^2 LL}{\partial w^2}$$

Intuition: $f(w) = LL(w)$. Use Taylor series expansion for $f(w)$
at point w_n = current w at iteration n

$$f(w) = f(w_n) + \frac{\partial f(w_n)}{\partial w} \frac{(w-w_n)}{1!} + \frac{\partial^2 f(w_n)}{\partial w^2} \frac{(w-w_n)^2}{2!} + \frac{\partial^3 f(w_n)}{\partial w^3} \frac{(w-w_n)^3}{3!}$$

$$\approx f(w_n) + \frac{\partial f(w_n)}{\partial w} (w-w_n) + \frac{1}{2} \frac{\partial^2 f(w_n)}{\partial w^2} (w-w_n)^2 \quad (\text{two terms})$$

Solve for w
assume $f(w) \approx f(w_n) = \frac{\partial f(w_n)}{\partial w} (w-w_n) + \frac{1}{2} \cdot 2 \frac{\partial^2 f(w_n)}{\partial w^2} (w-w_n)^2$

$$-\frac{\partial f}{\partial w} = \frac{\partial^2 f}{\partial w^2} (w-w_n)$$

$$w = w_n - \frac{\frac{\partial f}{\partial w}(w_n)}{\frac{\partial^2 f}{\partial w^2}(w_n)}$$

multidim $w = (w^1 w^2 \dots w^D)$:
 first differential is $\nabla f = \frac{\partial f}{\partial w} = \left(\frac{\partial f}{\partial w^1}, \frac{\partial f}{\partial w^2}, \dots, \frac{\partial f}{\partial w^D} \right)$
 vector

2nd differential is Hessian Matrix partial second deriv.

$$D^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial w^1 \partial w^1} & \frac{\partial^2 f}{\partial w^1 \partial w^2} & \dots & \frac{\partial^2 f}{\partial w^1 \partial w^D} \\ \vdots & \ddots & & \vdots \\ \frac{\partial^2 f}{\partial w^m \partial w^1} & \dots & \frac{\partial^2 f}{\partial w^m \partial w^D} \\ \vdots & & \ddots & \frac{\partial^2 f}{\partial w^D \partial w^D} \end{bmatrix}$$

GD update: $w_{\text{new}} = w - \alpha \nabla f(w) - \frac{1}{2} \alpha^2 H^{-1} \nabla^2 f(w)$

learn rate α

Intuition: $\frac{\partial f}{\partial w}$ 1st diff
 $\frac{\partial^2 f}{\partial w^2}$ 2nd diff

Newton-Hessian for LL=OBJ of Logistic Regression

$$H = X^T \underline{\text{DIAG}} \cdot X$$

Diag:

