

Quiz on Wednesday

Project proposal due Nov 3rd

HW3 due Nov 8th

**TODAY**

Logistic Regression

Gradient descent and Newton's Method

Separating hyperplanes (support vector machine)

In LR we assume  $\eta(x) = P(Y=1 | X=x)$

$$= \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}}$$

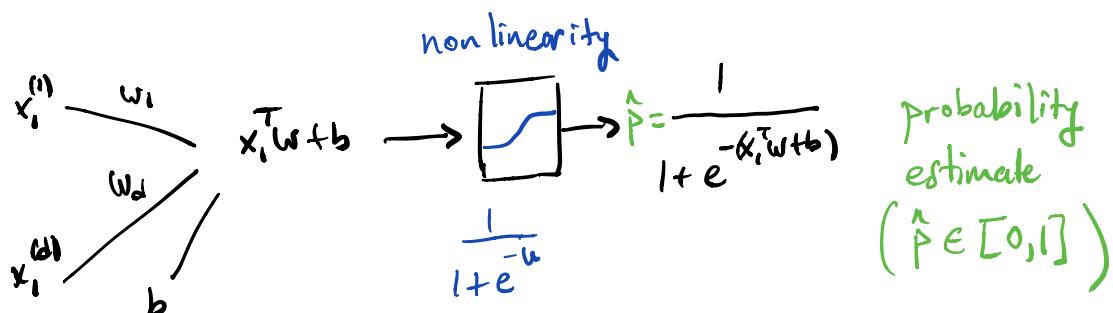
Given data  $(x_1, y_1), \dots, (x_n, y_n)$  we estimate  $w$  and  $b$   
using the method of Maximum Likelihood

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n p(y_i | x_i, \theta) \\ &= \prod_{i=1}^n p(x_i)^{y_i} \cdot (1 - p(x_i))^{1-y_i} \end{aligned}$$

where  $p(x_i) = P(Y=1 | X=x_i)$

$$-\log L(\theta) = - \sum_{i=1}^n y_i \cdot \log p(x_i) - (1 - y_i) \cdot \log (1 - p(x_i))$$

(cross-entropy loss)



$$-\log L(\theta) = -\sum_i \log(1-p(x_i)) - \sum_i y_i \log \frac{p(x_i)}{1-p(x_i)}$$

$$= -\sum_{i=1}^n \log \frac{1}{1+e^{w^T x_i + b}} - \sum_{i=1}^n y_i (w^T x_i + b)$$

Define:  $\tilde{x}_i = \begin{bmatrix} x_i \\ 1 \end{bmatrix}$      $\theta = \begin{bmatrix} w \\ b \end{bmatrix}$      $\tilde{x}_i^T \theta = x_i^T w + b$

Exercise: If we define new labels  $\tilde{y}_i \in \{-1, +1\}$ , then

$$-\log L(\theta) = \sum_{i=1}^n \log \left( 1 + e^{-\tilde{y}_i \theta^T \tilde{x}_i} \right) \quad (\tilde{y}_i = 2 \cdot y_i - 1)$$

$$y_i = 0 \Rightarrow \tilde{y}_i = -1$$

$-\log L(\theta)$ : convex function.

$$y_i = +1 \Rightarrow \tilde{y}_i = +1$$

$$\nabla -\log L(\theta) = 0 \quad \text{no closed form solution!}$$

$$\nabla -\log L(\theta) = X^T (y - \hat{p}) \quad \boxed{\text{HW 3}}$$

$$\nabla^2 -\log L(\theta) = X^T S X \quad S \text{ is diagonal.}$$

Gradient Descent

$$\min_{\theta} l(\theta)$$

$$\theta_{t+1} = \theta_t - \alpha_t \cdot \nabla l(\theta_t)$$

$\alpha_t$ : step size  
 $\exp \alpha = 0.1$

$$\begin{aligned} l(\theta_{t+1}) &= l(\theta_t + \alpha \cdot v) \\ &= l(\theta_t) + \alpha \cdot \langle \nabla l(\theta_t), v \rangle + O(\alpha^2) \end{aligned}$$

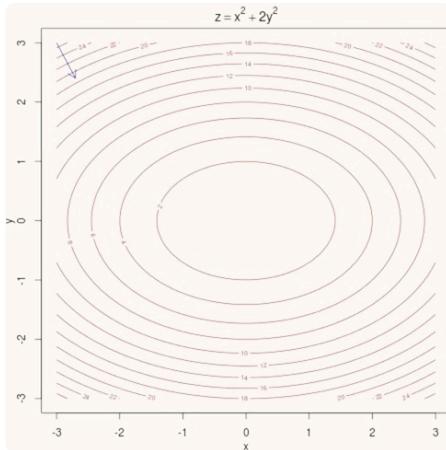
$v^* = \frac{-\nabla l(\theta_t)}{\|\nabla l(\theta_t)\|}$  will minimize the above Taylor's approximation

$$\|\nabla \ell(\theta_t)\|_2$$

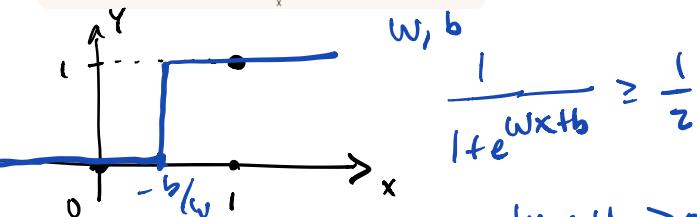
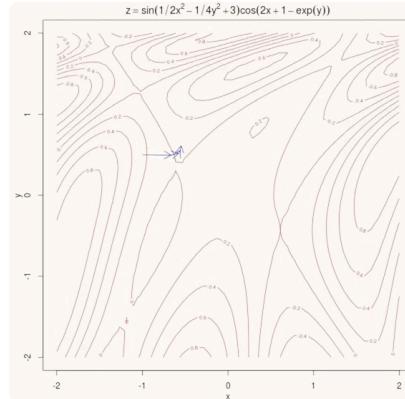
0 "

$$\text{Cauchy-Schwarz: } \langle \nabla \ell(\theta_t), v \rangle \geq -\|\nabla \ell(\theta_t)\|_2 \cdot \|v\|_2$$

Convex



Non-Convex



linearly separable:  $\hat{p} = y \Rightarrow$  stationary point  $\nabla L(\theta) = 0.$

Newton's Method

$$\theta_{th} = \theta_t - (\nabla^2 \ell(\theta_t))^{-1} \nabla \ell(\theta_t)$$

$$\frac{1}{1+e^{wx+b}} \geq \frac{1}{2}$$

$$wx+b \geq 0$$

$$x \geq -b/w$$

not unique!

$$\frac{-\alpha \cdot b}{\alpha w} = -\frac{b}{w}$$

Second order  
Taylor's approximation.

$$\ell(\theta) = \ell(\theta_t) + \langle \nabla \ell(\theta_t), \theta - \theta_t \rangle + \frac{1}{2} (\theta - \theta_t)^T \nabla^2 \ell(\theta_t) (\theta - \theta_t)$$

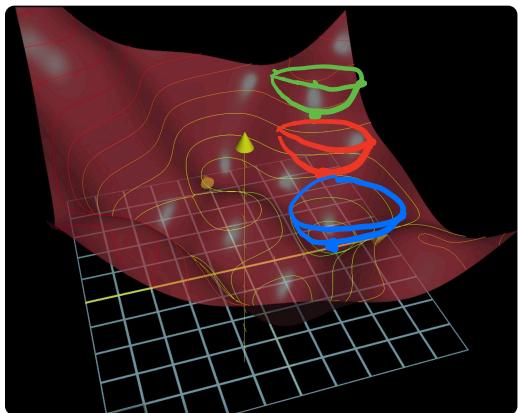
$$\nabla F(\theta) = 0$$

$$+ O(\|\theta - \theta_t\|^3)$$

$$\nabla \ell(\theta_t) + \nabla^2 \ell(\theta_t) (\theta - \theta_t) = 0$$

$$\theta^* = \theta_t - (\nabla^2 \ell(\theta_t))^{-1} \nabla \ell(\theta_t)$$

$= \theta_{\text{new}}$  (Newton's update)



GD

Complexity

$\mathcal{O}(nd)$

$$\begin{matrix} n \\ \times \\ d \end{matrix}$$

Newton  $\mathcal{O}(nd^2)$  Hessian  $X^T S X$

$$\begin{matrix} n \\ d \\ \times \\ d \end{matrix}$$

$\mathcal{O}(nd^2)$  (Hessian) Gradient  
direct methods

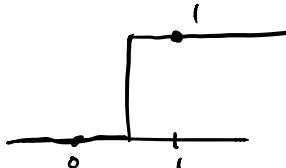
SVD

QR decomp.

## Regularized Logistic Regression

Problem 1 :  $n < d \Rightarrow (\text{Hessian})^{-1}$  doesn't exist!

Problem 2 :



linearly separable     $\theta \rightarrow \infty$   
no unique solution    unique solutions

$$\min L(\theta) + \lambda \|\theta\|_2^2 \quad \underbrace{\text{L2 regularization}}$$

$$\min L(\theta) + \lambda \|\theta\|_1 \quad \text{L1 regularization}$$

Exercise

$$\text{LDA assumptions} \Rightarrow p(Y=1 | X=x) = \frac{1}{1+e^{-(w^T x + b)}}$$

## Separating hyperplanes (support vector machine)

LDA / QDA / NB  $p(x)$  and  $p(Y|x)$

logistic  $- p(Y|x)$

T.7

-SVM

$$w^T x + b \geq 0$$

directly estimate  $\theta = [w]$   
 $b$

**Vapnik's Principle:**

"When solving a problem of interest, do not solve a more general problem as an intermediate step."

"Don't solve a harder problem than you have to"