

today:

- logistic regression
- numerical optimization

Logistic regression:

setup: binary classification $Y = \{0, 1\}$ $X \in \mathbb{R}^d$

generative model motivation:

suppose only assumption is there exists pdf's (densities) in \mathbb{R}^d for each class conditional:

$$p(x|Y=1) \quad \downarrow \quad p(x|Y=0)$$

$$p(Y=1|X=x) = \frac{p(Y=1, X=x)}{p(Y=1, X=x) + p(Y=0, X=x)} \quad \left. \begin{array}{l} p(Y=1, X=x) \\ p(Y=0, X=x) \end{array} \right\} p(X=x)$$

$$= \frac{1}{1 + \frac{p(Y=0, X=x)}{p(Y=1, X=x)}} = \frac{1}{1 + \exp(-f(x))}$$

where $f(x) \triangleq \log \frac{p(X=x|Y=1)}{p(X=x|Y=0)}$

"log odd"

$\underbrace{\log \frac{p(X=x|Y=1)}{p(X=x|Y=0)}}_{\text{class-conditional ratio}}$

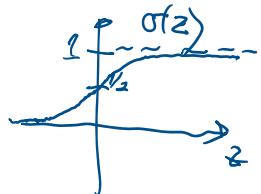
$\underbrace{\log \frac{p(Y=1)}{p(Y=0)}}_{\text{prior odds ratio}}$

In general

$$p(Y=1|X=x) = \sigma(f(x))$$

$$\text{where } \sigma(z) \triangleq \frac{1}{1 + \exp(-z)}$$

"sigmoid function"



Some properties of $\sigma(z)$: $\sigma(-z) = 1 - \sigma(z)$ $[\sigma(z) + \sigma(-z) = 1]$

$$\frac{d}{dz} \sigma(z) = \sigma(z) \sigma'(-z) = \sigma(z)(1 - \sigma(z))$$

* to motivate linear logistic regression

consider class conditional in the exponential family

$$p(x|m) \triangleq h(x) \exp(m^T T(x) - A(m))$$

"canonical parameter"

scalar fct. \rightarrow log partition fct.

"sufficient statistics"

normalize

These specify the "flat" exponential family

Gaussian:

$$\sim N(\mu, \Sigma) \Rightarrow -\frac{1}{2} \ln |\Sigma| - \frac{1}{2} (\mu - \mu_0)^T \Sigma^{-1} (\mu - \mu_0)$$

Gaussian:

$$\log p(x|\mu, \sigma^2) = -\frac{1}{2} \log 2\pi\sigma^2 - \frac{(x-\mu)^2}{2\sigma^2} - \left(\frac{x^2}{2\sigma^2} - \frac{x\mu}{\sigma^2} + \frac{\mu^2}{2\sigma^2} \right)$$

\downarrow
 $\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

$$\text{let } T(x) = \begin{bmatrix} -x^2/2 \\ x \end{bmatrix}$$

$$m(\mu, \sigma^2) = \begin{bmatrix} \sigma^2 \\ \mu/\sigma^2 \end{bmatrix}$$

$$A(\mu) = \frac{1}{2} \log(2\pi\sigma^2) + \frac{\mu^2}{2\sigma^2}$$

$$p(x|Y=1) = p(x|m_1)$$

$$p(x|Y=0) = p(x|m_0)$$

$$\text{log odds } f(x) = \log \frac{p(x|m_1)}{p(x|m_0)} + \log \frac{\pi}{1-\pi}$$

$$= (m_1 - m_0)^T T(x) + A(m_0) - A(m_1) + \log \left(\frac{\pi}{1-\pi} \right)$$

$$\triangleq w^T \phi(x)$$

$$\text{where } w = \begin{pmatrix} m_1 - m_0 \\ A(m_0) - A(m_1) + \log \left(\frac{\pi}{1-\pi} \right) \end{pmatrix} \quad \varphi(x) = \begin{pmatrix} T(x) \\ 1 \end{pmatrix}$$

get logistic regression model?

$$P_w(Y=1 | X=x) = \sigma(w^T \varphi(x))$$

"feature map"

decision boundary

$$\boxed{I \in \{w^T \varphi(x) > 0\}}$$

exercise to reader : try argument above with $p(x|y) = N(x|\mu_y, \Sigma_y)$

- if $\Sigma_y = \Sigma_1$, then $\varphi(x) = \begin{pmatrix} x \\ 1 \end{pmatrix} \rightarrow \text{"linear boundary"}$

- otherwise, $\varphi(x) = \begin{pmatrix} x \\ x^T \\ x \\ 1 \end{pmatrix}$

→ quadratic decision boundary

linear logistic regression model

linear logistic regression model

$$p(y=1|x, w) = \sigma(w^T x) \quad y \in \{0, 1\} \quad \text{decision rule: } \mathbb{1}\{w^T x > 0\}$$

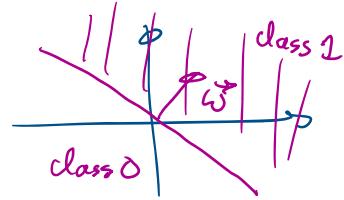
$$p(y=0|x, w) = 1 - \sigma(w^T x) = \sigma(-w^T x)$$

[if $y \in \{0, 1\}$ "Rademacher R.V."]

$$\text{encode } p(y|w) = \sigma(y w^T x)$$

$y | x=x$ is Bernoulli($\sigma(w^T x)$)

$$p(y|x) = \sigma(w^T x)^y (1 - \sigma(w^T x))^{1-y}$$



given $(x_i, y_i)_{i=1}^n$ iid

max. conditional log-likelihood to estimate \hat{w}_{ML}

$$l(w) = \sum_{i=1}^n \log p(y_i | x_i, w) = \sum_{i=1}^n y_i \log \sigma(w^T x_i) + (1-y_i) \log (1 - \sigma(w^T x_i))$$

$$\nabla_w \sigma(w^T x) = x \left[\sigma(w^T x_i) \sigma(-w^T x_i) \right] \quad \text{let } v_i \triangleq w^T x_i$$

$$\nabla_w l(w) = \sum_{i=1}^n \partial v_i \left[\sigma(v_i) \partial_w \sigma(v_i) \frac{y_i}{\sigma(v_i)} - \frac{(1-y_i)}{\sigma(v_i)} \right]$$

$$= \sum_{i=1}^n x_i \left[y_i \underbrace{\left[\sigma(-v_i) + \sigma(v_i) \right]}_{=1} - \sigma(v_i) \right]$$

$$\boxed{\nabla l(w) = \sum_{i=1}^n x_i [y_i - \sigma(w^T x_i)]}$$

need to use
numerical
methods

solve for $\nabla l(w) = 0 \Rightarrow$ need to solve transcendental equation

because $\frac{1}{1 + \exp(-w^T x_i)} \dots = 0$
solve for this

contrast to

linear regression

$$\nabla l(w) = \sum_{i=1}^n x_i [y_i - w^T x_i]$$

(non-linear w)

numerical optimization

want to minimize $f(w)$ (unconstrained) [suppose f is differentiable]
s.t. $w \in \mathbb{R}^d$

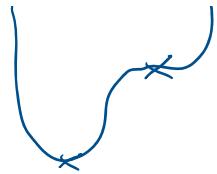
- 1) gradient descent (1st order method)

↓ ↓ ↓ ..



1) gradient descent (1st order method)

start w_0
 iterate : $w_{t+1} = w_t - \gamma_t \nabla f(w_t)$
 stopping criterion : if $\|\nabla f(w_t)\|^2 \leq \delta$
 then stop



$$\text{e.g. } \delta = 10^{-6}$$

note : if f is μ -strongly convex $\Rightarrow f(w_t) - \min_w f(w) \leq \frac{1}{2\mu} \|\nabla f(w_t)\|^2$
 $(\Leftrightarrow f(w) - \frac{1}{2} \|w\|^2 \text{ is convex in } w)$
 if $\|\nabla f(w_t)\|^2 \leq \delta \Rightarrow \text{subopt.}(w_t) \leq \frac{\delta}{2\mu}$

step-size rules

a) constant step-size : $\gamma_t = \frac{1}{L}$ L Lipschitz continuity constant for ∇f
 i.e. $\|\nabla f(w) - \nabla f(w')\| \leq L \|w - w'\|_2$

b) decreasing step-size rule
 $\gamma_t = \underbrace{\gamma_0}_{\leftarrow \text{constant}} \cdot \alpha^t$ $\forall w, w' \in \mathbb{R}^d$
 This is more common for stochastic optimization

usually want $\sum_t \gamma_t = +\infty$. e.g. $f(w) = \mathbb{E}_\xi g(w, \xi)$
 $\sum_t \gamma_t^2 < \infty$ $w_{t+1} = w_t - \gamma_t \nabla_w g(w, \xi_t)$
 e.g. $g(w, \xi) = \ell(y_i, h_w(x_i))$ for ERM
 $\xi = (x_i, y_i)$

c) choose γ_t by "line search" : $\min_{\gamma \in \mathbb{R}} f(w_t + \gamma d_t)$
 costly in general direction of update (e.g. $\nabla f(w_t)$)

\hookrightarrow instead do approximate search
 e.g. Armijo line search

[see Boyd's book]

15h38

Newton's method (2nd order method)

motive: minimizing a quadratic approximation

Truncate expansion $f(w) \approx f(w_t) + \nabla f(w_t)^T (w - w_t) + \frac{1}{2} (w - w_t)^T H(w_t) (w - w_t)$

Hessian $[H(w_t)]_{ij} = \frac{\partial^2 f(w_t)}{\partial w_i \partial w_j}$

Taylor expansion
at w_t

$$f(w) = f(w_t) + \nabla f(w_t)^T (w - w_t) + \frac{1}{2} (w - w_t)^T H(w_t) (w - w_t)$$

$$\stackrel{?}{=} Q_t(w)$$

$$= Q_t(w) + O(\|w - w_t\|^3)$$

+ $O(\|w - w_t\|^3)$
Taylor's remainder

quadratic model approximation

w_{t+1} ~ obtain by minimizing $Q_t(w)$

$$\nabla_w Q_t(w) = 0$$

$$\nabla f(w_t) + H(w_t)(w - w_t) \xrightarrow{\text{want}} 0$$

$$\Rightarrow w - w_t = H^{-1}(w_t) \nabla f(w_t)$$

$$w_{t+1} = w_t - H^{-1}(w_t) \nabla f(w_t)$$

inverse Hessian

$\sim O(d^3)$ time
and $O(d^2)$ space

Newton's update

$$d_t = H^{-1} \nabla f$$

$$H d_t = \nabla f$$

Note for homework 2: solve for $H_t d_t = \nabla f(w_t)$

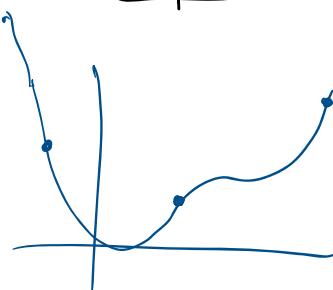
$$\min_d \|H_t d - \nabla f(w_t)\|^2$$

use numpy.linalg.lstsq to solve this

→ much more numerically stable

Clamped Newton: you add a step-size to stabilize Newton's method

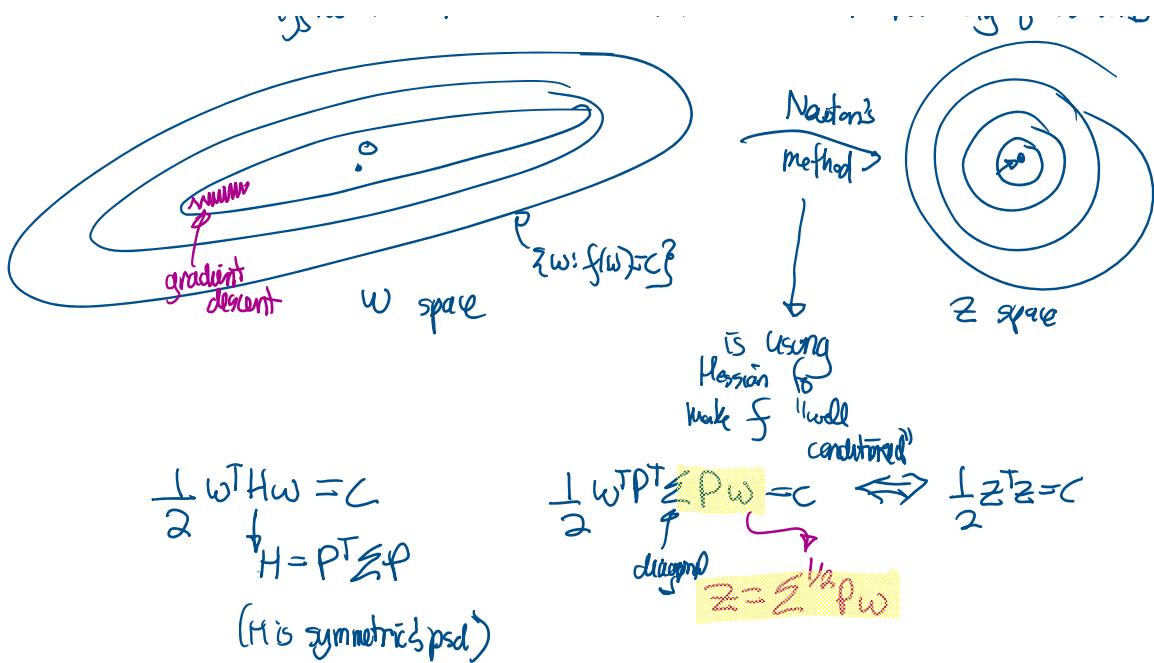
$$w_{t+1} = w_t - \underbrace{\frac{\nabla f(w_t)}{H(w_t)}}_{\text{step size}} \nabla f(w_t)$$



why Newton's method?

- much faster convergence in # iterations vs gradient descent
- affine invariant \Rightarrow method is invariant to rescaling of variables





Exercise to the reader:

$$z_{t+1} = z_t - \gamma \nabla f(z_t) \Leftrightarrow w_{t+1} = w_t - \gamma H^{-1} \nabla f(w_t)$$

Big data Logistic regression

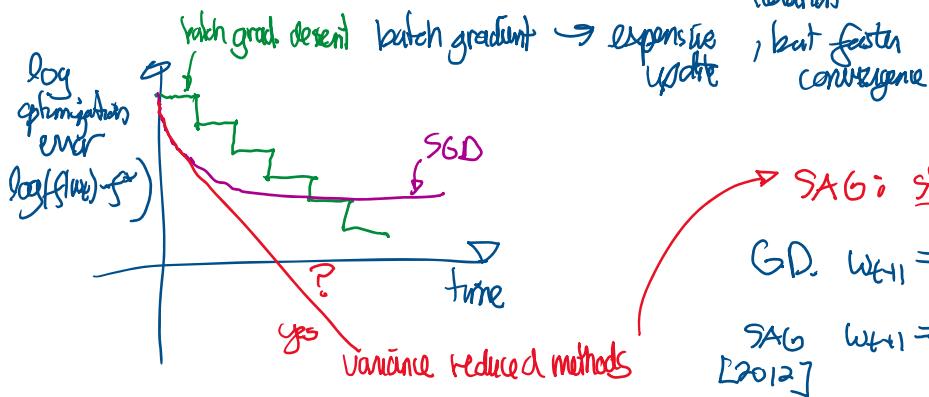
- big $d \rightarrow$ cannot do $O(d^2)$ or $O(d^3)$ operations \Rightarrow first order methods
- if n is large, you cannot use batch methods $\nabla f(w) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(w)$ (gradient of loss on some data pt)
(full) batch gradient

instead you can use "incremental gradient method"

e.g. stochastic gradient descent (SGD) $\therefore w_{t+1} = w_t - \gamma_t \nabla f_{i_t}(w_t)$

where i_t is picked uniformly at random

SGD \rightarrow cheap updates, but slower convergence per iteration



\rightarrow SAG; stochastic average gradient

GD. $w_{t+1} = w_t - \gamma_t \sum_{i=1}^n \nabla f_i(w_t)$

SAG [2012] $w_{t+1} = w_t - \gamma_t \frac{1}{n} \sum_{i=1}^n V_P$ memory

yes \ Variance Reduced methods

SAG [2012]

$$w_{t+1} = w_t - \gamma_t \frac{1}{n} \sum_{i=1}^n V_i$$

where $V_i = \nabla f_i(w_t)$

at each t, update $V_{it} = \nabla f_i(w_t)$

SAGA [2014]

(default method
for large scale logistic
regression in Scikit Learn)

SVRG

$$w_{t+1} = w_t - \gamma_t \left(\nabla f_i(w_t) + \frac{1}{n} \sum_{j=1}^n V_j - \frac{V_{i,t}}{\nabla f_i(w_t)} \right)$$

Variance Reduction Correction

Skipped part: Newton's method on logistic regression = IRLS:

Newton's method for logistic regression: IRLS

$$\text{recall for } l(w) : \nabla l(w) = \sum_{i=1}^n x_i [y_i - \sigma(w^T x_i)]$$

$$H(l(w)) = -\sum_{i=1}^n x_i x_i^T \sigma'(w^T x_i) \sigma(w^T x_i)$$

$$V^T H V = -\sum_{i=1}^n (\underbrace{x_i^T V}_{(x_i^T V)^2 \geq 0}) \underbrace{(\sigma'(-) \sigma(-))}_{\geq 0} \quad V^T H V \leq 0 \quad H \text{ is I.R.Q}$$

i.e. HJO
i.e. concave fct.

Notation: recall

$$X = \begin{pmatrix} -x_1^T & - \\ \vdots & \end{pmatrix}$$

Newton is maximizing
instead of minimizing

$$\text{let } \mu_i \triangleq \sigma(w^T x_i) \in]0, 1[$$

$$\nabla l(w) = \sum_i x_i [y_i - \mu_i] = X^T (y - \mu)$$

$$\text{Hessian} = -\sum_i x_i x_i^T \mu_i (1 - \mu_i) = -X^T D(w) X$$

$D_{ii} \triangleq \mu_i (1 - \mu_i)$

diagonal matrix

$$\text{Newton's update: } w_{t+1} = w_t - (X^T D_t X)^{-1} X^T (y - \mu_t)$$

$$= (X^T D_t X)^{-1} \left[(X^T D_t X) w_t + X^T (y - \mu_t) \right]$$

$$w_{t+1} = (X^T D_t X)^{-1} [X^T D_t z_t]$$

where $z_t \triangleq X w_t + D_t^{-1} (y - \mu_t)$

$w_{t+1} \leftarrow \dots \leftarrow z_t \leftarrow \dots \leftarrow \text{initialized } w_{t-1} \leftarrow \dots \leftarrow w_0 \leftarrow \dots \leftarrow w_0$

this is a solution to a "weighted least square problem"

$$\min_w \left\| D^{1/2} (z - x^w) \right\|_2^2$$

$\xrightarrow{\text{weights}}$

$$\sum_i (z_i - w^T x_i)^2 / d_{ii}^{-1}$$

$\xrightarrow{\text{new target}}$

compare with
Gaussian noise
model for least
square

$$\sum_i (y_i - w^T x_i)^2 / \sigma_i^2$$

Newton's method for logistic regression

= iterated reweighted least squares (IRLS)