

today : • logistic regression
• numerical optimization

logistic regression model :

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

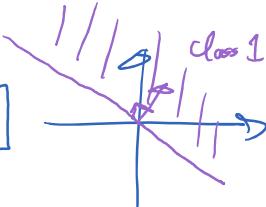
$$p(y=0|x, w) = 1 - \sigma(w^T x) = \sigma(-w^T x)$$

[if $Y = \{-1, 1\}$ "Rademacher R.v."]

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

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

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



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

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

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

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

$$\begin{aligned} \nabla_w l(w) &= \sum_{i=1}^n x_i \left[\underbrace{\sigma(v_i) \sigma(-v_i)}_{\sigma(v_i)} y_i - \frac{(1-y_i)}{\sigma(-v_i)} \right] \\ &= \sum_{i=1}^n x_i \left[y_i \underbrace{[\sigma(v_i) + \sigma(v_i)]}_{1} - \sigma(v_i) \right] \end{aligned}$$

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

need to use numerical methods

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

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

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

\uparrow linear w

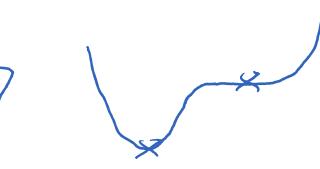
Numerical optimization

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

i) gradient descent (1st order method)

i) gradient descent (1st order method)

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



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

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

step size rules:

a) constant step-size : $\gamma_t = \frac{1}{L} \leftarrow$ Lipschitz continuity constant for ∇f

$$\text{f.e. } \|\nabla f(w) - \nabla f(w')\| \leq L \|w - w'\| \quad \forall w, w' \in \mathbb{R}^d$$

b) decreasing step-size rule [This is more common for stochastic optimization]

$$\gamma_t = \frac{C}{t} \leftarrow \text{constant}$$

$$\begin{aligned} \text{usually want : } & \sum_t \gamma_t = \infty \\ & \sum_t \gamma_t^2 < \infty \end{aligned}$$

$$\text{e.g. } f(w) \triangleq \mathbb{E}_{\xi} g(w, \xi)$$

$$w_{t+1} = w_t - \gamma_t \nabla_w g(w, \xi_t)$$

$$\text{e.g. } g(w, \xi) = f(y_i, h_w(x_i)) \text{ for ERM}$$

$$\xi = (x_i, y_i)$$

c) choose γ_t by "line search" :

$$\min_{\gamma \in \mathbb{R}} f(w_t + \gamma d_t)$$

costly in general

direction of update
 $\rightarrow \nabla f(w_t)$

instead do approximate search

e.g. Armijo line search

[see Boyd's book]

15h50

Newton's method (2nd order method)

motivate : minimizing a quadratic approximation

$$\begin{aligned}
 \text{Taylor expansion at } w_t : f(w) &= f(w_t) + \nabla f(w_t)^T (w - w_t) + \frac{1}{2} (w - w_t)^T H(w_t) (w - w_t) \\
 &\triangleq Q_t(w) \\
 &= Q_t(w) + O(\|w - w_t\|^3)
 \end{aligned}$$

$$\begin{aligned}
 \triangleq \text{Hessian} \\
 [H(w_t)]_{ij} &= \frac{\partial^2 f(w_t)}{\partial w_i \partial w_j}
 \end{aligned}$$

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

\mathcal{L} quadratic model approximation

$w_{t+1} \rightsquigarrow$ obtain by minimizing $\mathcal{Q}_t(w)$

$$\nabla_w \mathcal{Q}_t(w) = 0 \quad \nabla f(w_t) + H(w_t)(w - w_t) \xrightarrow{\text{want}} 0$$

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

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

Newton's update
inverse Hessian $d_t = H^{-1} \nabla f$
 \rightsquigarrow
 $\rightsquigarrow O(d^3)$ time
and $O(d^2)$ space $H d_t = \nabla f$

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

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

use numpy, linalg. linalg. lsq to solve this ↗

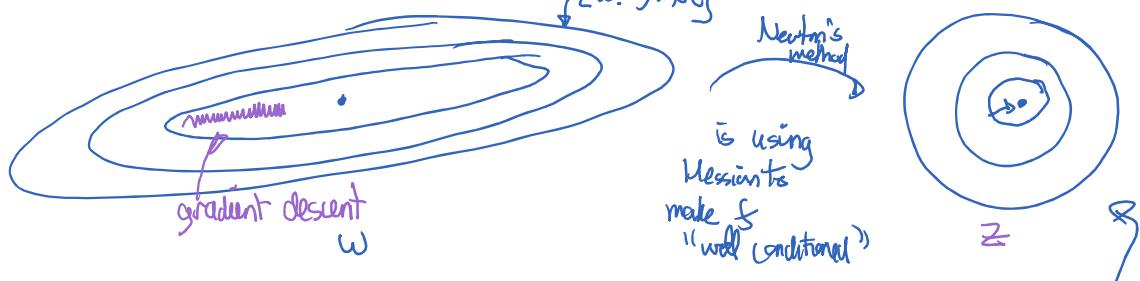
↗ much more numerically stable

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

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

why Newton's method?

- much faster convergence in # iterations vs gradient descent
- affine invariant \Rightarrow method is invariant to rescaling of variables
 $\Sigma w: f(w) \propto \frac{1}{2} w^T H w$



$$\frac{1}{2} w^T H w = C$$

\downarrow
 $H = P^T \Sigma P$
(H is symmetric)
 \downarrow psd.

$$\frac{1}{2} w^T P^T \Sigma P w = C \text{ and } \frac{1}{2} z^T z = C$$

\downarrow
diagonal
 $z = P^{-1} w$

exercise to recall: $z_{t+1} = z_t - \gamma \nabla \tilde{f}(z_t)$

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

Newton's method for logistic regression: IR LS

recall for $\ell(w)$: $\nabla \ell(w) = \sum_{i=1}^n x_i [y_i - \sigma(w^T x_i)]$
 $H(\ell(w)) = -\sum_{i=1}^n x_i x_i^T \sigma'(w^T x_i) \sigma(w^T x_i)$

$$V^T H V = -\sum_{i=1}^n \underbrace{[V^T x_i] (x_i^T V)}_{(x_i^T V)^2 \geq 0} \underbrace{\sigma'(-)}_{\geq 0} \underbrace{\sigma'(+)}_{\geq 0} \quad V^T H V \leq 0 \quad \forall v \in \mathbb{R}^d$$

Notation: recall

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

i.e. HJO
i.e. concave fn.
Newton is maximizing instead of minimizing

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

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

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

↳ diagonal matrix

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

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

X
 $w_{t+1} = (X^T D_t X)^{-1} [X^T D_t z_t]$
where $z_t \triangleq X w_t + D_t^{-1} (y - \mu_t)$

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

$$\min_w \|D^{1/2} (z_t - X w)\|_2^2$$

\uparrow new target
 \downarrow weights

$$\rightarrow \sum_i (z_i - w^T x_i)^2$$

d_{ii}^{-1}

compare with Gaussian noise model for least square

$$\sum_i \frac{(y_i - w^T x_i)^2}{2\sigma^2}$$

Newton's method for logistic regression

= iterated reweighted least squares (IR LS)

Big data logistic regression

• big d \Rightarrow cannot do $O(d^2)$ or $O(d^3)$ operations \Rightarrow first order methods

gradient descent, ...

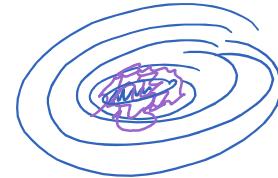
- big $d \Rightarrow$ cannot do $O(d^2)$ or $O(d^3)$ operations \Rightarrow first order methods
- if n is large, you cannot use batch methods $\nabla f(w) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(w)$ (gradient of some data point)
 $\underbrace{\nabla f_i(w)}_{\text{batch gradient}}$ $\Theta(nd)$

Instead you use "incremental gradient methods"

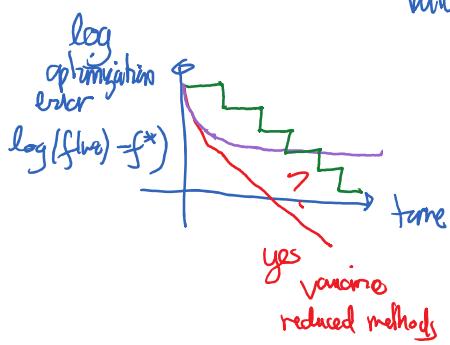
e.g. stochastic gradient descent (SGD): $w_{t+1} = w_t - \gamma_t \nabla f_{i_t}(w_t)$

where i_t is
picked uniformly at random $\mathcal{O}(d)$

SGD \rightarrow cheap updates, but slower convergence per iteration



batch gradient \rightarrow expensive updates, but faster convergence



SGD: stochastic average gradient

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

$$SAG: w_{t+1} = w_t - \gamma_t + \frac{1}{n} \sum_{i=1}^n v_i \leftarrow \text{memory}$$

where $v_i = \nabla f_i(w_t)$

at each t , update $v_{i,t} \triangleq \nabla f_i(w_t)$.

$$SAGA: w_{t+1} = w_t - \gamma_t \left(\nabla f_t(w_t) + \frac{1}{n} \sum_{j=1}^{t-1} v_j \right)$$

(default method
for large scale logistic regression in Scikit-learn)

$$\frac{1}{n} \sum_{j=1}^n v_j - v_{i,t} \nabla f_i(w_t)$$

Variance reduction correction

SVRG