

today - Logistic regression
Fisher LDA and multivariate Gaussians
(generative)

Logistic regression

$Y \in \{0, 1\}, X \in \mathbb{R}^d$

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



$p(x|y)$ is in exp. family for each y

then $p(y=1|x) = \sigma(w^T \phi(x))$

parameter related to sufficient statistics $T(x)$ of exp family

logistic regression model

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

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

[if had encoded $y \in \{-1, 1\}$, $p(y|x) = \sigma(y w^T x)$]

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

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

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

maximum conditional likelihood:

$L(w) = \sum_{i=1}^n \log p(y_i|x_i; w) = \sum_{i=1}^n [y_i \log \sigma(w^T x_i) + (1-y_i) \log \sigma(-w^T x_i)]$

$\sigma'(z) = \sigma(z)\sigma(-z)$

$\nabla \sigma(w^T x_i) = x_i [\sigma(w^T x_i) \sigma(-w^T x_i)]$ $v_i \triangleq w^T x_i$

$\nabla L(w) = \sum_{i=1}^n x_i \left[\frac{y_i \sigma(v_i) \sigma(-v_i)}{\sigma(v_i)} - \frac{(1-y_i) \sigma(-v_i) \sigma(v_i)}{\sigma(-v_i)} \right]$
 $\left[y_i [\sigma(v_i) + \sigma(-v_i)] - \sigma(v_i) \right]$

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

need numerical optimization

$\nabla L(w) = 0 \Rightarrow$ need to solve transcendental eq. because $\frac{1}{1 + \exp(w^T x_i)} (\dots) = 0$
 solve for this

(contrast with least squares where get linear in w)

$$\sum_{i=1}^n x_i [y_i - w^T x_i] = 0$$

brief recap on optimization:

want to minimize $\min_{w \in \mathbb{R}^d} f(w)$

1) gradient descent start at w_0

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

step-size rules:

- constant step-size $\delta_t = \frac{1}{L}$ $L \leftarrow$ Lipschitz constant of ∇f
- decreasing step-size rule: $\delta_t = \frac{c}{t}$ $c \leftarrow$ constant
usually want $\sum_t \delta_t = \infty$, $\sum_t \delta_t^2 < \infty$
- choose δ_t by "line-search"; $\min_{\delta \in \mathbb{R}} f(w_t + \delta d_t)$ $d_t \leftarrow$ direction you were exploring $[-\nabla f(w_t)]$
 \downarrow costly in general
 \rightarrow can use instead approximate search
eg. Armijo line search

2) Newton's method (2nd order method)

motivation: minimize quadratic approximation:

Taylor expansion: $f(w) = f(w_t) + \nabla f(w_t)^T (w - w_t) + \frac{1}{2} (w - w_t)^T H(w_t) (w - w_t) + O(\|w - w_t\|^3)$

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

$[H(w)]_{ij} \triangleq \frac{\partial^2 f(w)}{\partial w_i \partial w_j}$ quadratic function

want

$$w_{t+1} = \min_w Q_t(w)$$

$$\nabla Q_t(w)$$

$$= \nabla f(w_t) + H(w_t)(w - w_t)$$

$$\Rightarrow 0$$



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

this rule is

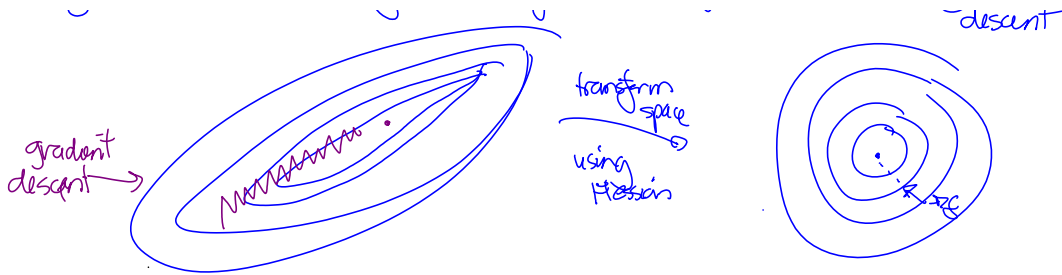
$$w_{t+1} = w_t - \underbrace{H^{-1}(w_t)}_{\text{inverse Hessian}} \nabla f(w_t)$$

Newton's update

damped Newton: you add a step-size

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

why Newton: • much faster convergence in # of iterations vs gradient



• affine invariant

but Newton need Hessian $\rightarrow O(d^2)$ space
 $O(d^3)$ time

Newton for logistic regression: IRLS:

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

$H(L(w)) = -\sum_{i=1}^n x_i x_i^T \sigma(w^T x_i) \sigma(-w^T x_i)$

$v^T H v = -\sum_{i=1}^n \underbrace{(v^T x_i)^2}_{\geq 0} \underbrace{\sigma(-) \sigma(-)}_{\geq 0}$

$v^T H v \leq 0 \quad \forall v$
 i.e. $H \leq 0$

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

let's look at minimizing the negative log-likelihood

recall design $X = \begin{pmatrix} \vdots \\ -x_i^T \\ \vdots \end{pmatrix}$
 $n \times d$

then Hessian = $\sum_{i=1}^n x_i x_i^T [\mu_i (1 - \mu_i)]$
 gradient = $\sum_{i=1}^n x_i [\mu_i - y_i] = X^T (\mu - y)$

$\rightarrow = X^T (DX)$ where $D_{ii} = \mu_i (1 - \mu_i)$] depends on w
 $\begin{bmatrix} \dots & x_i & \dots \end{bmatrix} \begin{bmatrix} \vdots \\ D_{ii} x_i^T \\ \vdots \end{bmatrix}$

Newton's update: $w_{t+1} = w_t - (X^T D_t X)^{-1} X^T (\mu - y)$
 $= (X^T D X)^{-1} [X^T D X w_t + X^T (y - \mu)]$
 $= (X^T D X)^{-1} [X^T D z_t]$ where $z_t \triangleq X w_t + D^{-1} (y - \mu)$

$(z_t)_i = x_i^T w_t + \frac{y_i - \mu_i}{\mu_i (1 - \mu_i)}$

this is solution to "weighted least square problem"

$\min_w \| D^{1/2} (z - Xw) \|^2$ new target
 $\hookrightarrow \sum_{i=1}^n \frac{(z_i - x_i^T w)^2}{d_{ii}}$

compare with Gaussian noise model for least square
 $\sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{\sigma^2}$
 noise

Newton's method for logistic regression = iterative reweighted least squares

Newton's method for logistic regression = iterative reweighted least squares

big data logistic regression

cannot do $O(d^2)$, $O(d^3)$ operations

⇒ first order methods

if n is huge, you cannot use batch method

$$\nabla f = \sum_{i=1}^n \nabla f_i(w)$$

$O(n)$

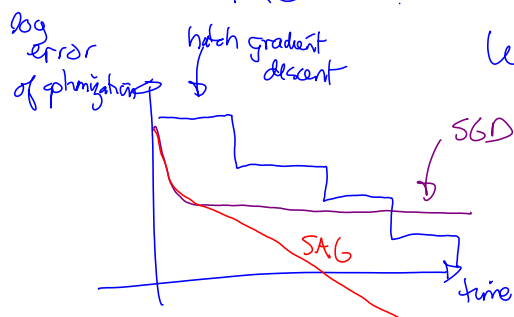
suppose $f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$

instead you can use SGD: stochastic gradient descent

SGD: $w_{t+1} = w_t - \gamma_t \nabla f_i(w_t)$ (much cheaper update but slow convergence)

batch gradient $w_{t+1} = w_t - \gamma_t \frac{1}{n} \sum_{i=1}^n \nabla f_i(w_t)$ (expensive update, fast convergence)

SAG: stochastic average gradient:



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

where $V_i = \nabla f_i(w_t)$

and at each t , update only one $V_i = \nabla f_i(w_t)$

kernels

recall for least-square, gradient $\sum_{i=1}^n x_i (y_i - w^T x_i)$

log-regression $\sum_{i=1}^n x_i (y_i - \sigma(w^T x_i))$

you always have $w_t = \sum_{i=1}^n \alpha_i x_i$

eg least-square:

$$w = (X^T X + \lambda I)^{-1} X^T y$$

$$= X^T \alpha \quad \text{where } \alpha = \underbrace{(X X^T + \lambda I)^{-1}}_{\text{Gram matrix } K} y$$

for prediction $w^T x = \sum_{i=1}^n \alpha_i (x_i^T x)$

Gram matrix K

$$K_{ij} = x_i^T x_j$$

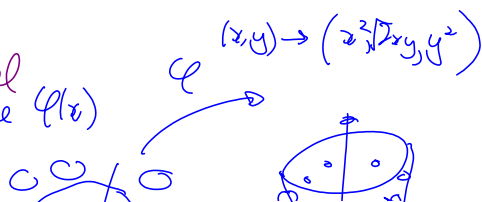
suppose map $x \xrightarrow{\phi} \phi(x)$

$$w^T \phi(x) = \sum_{i=1}^n \alpha_i \underbrace{\phi(x_i)^T \phi(x)}_{K(x_i, x)}$$

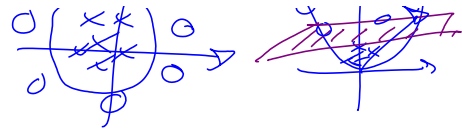
"kernel trick"

→ you can work implicitly in high dimensional space $\phi(x)$

only using $k(\cdot, \cdot)$ kernel



u u



$$z = (x, y)$$

$$\langle \ell(z), \ell(z') \rangle$$

$$= \left\langle \begin{pmatrix} x^2 \\ \sqrt{2}xy \\ y^2 \end{pmatrix}, \begin{pmatrix} x'^2 \\ \sqrt{2}x'y' \\ y'^2 \end{pmatrix} \right\rangle = (xx')^2 + 2(xx')(yy') + y^2y'^2$$
$$= (xx' + yy')^2$$
$$= \langle z, z' \rangle^2$$