# Lecture 4

_____

- weighted least squares:

$$\text{Loss}(\beta) = \frac{1}{2} \sum_{i=1}^{n} w_i (y_i - s_i)^2 = \frac{1}{2} \sum_{i=1}^{n} w_i (y_i - x_i^T \beta)^2$$

Continuous $\downarrow$ $s_i$

Approximate iteratively with quadratic terms

- Logistic regression:

$$p(y|x) = \frac{e^{ys}}{Z} = \frac{e^{ys}}{1 + e^s} \qquad y \in \{0,1\}$$

$$\text{Loglikelihood}(\beta) = \sum_{i=1}^{n} \left[ y_i s_i - \log(1 + e^{s_i}) \right]$$

$\ell(s_i)$

no analytical solution to maximize

- Iterative Reweighted Least Squares:

$$\beta_t \xrightarrow{\Delta\beta} \beta_{t+1}$$

$$x_i^T \downarrow \qquad \qquad \downarrow$$

$$\hat{s}_i \xrightarrow{\Delta s_i = x_i^T \Delta\beta} s_i = \hat{s}_i + \Delta s_i$$

$$\ell(s_i) \doteq \ell(\hat{s}_i) + \ell'(\hat{s}_i) \Delta s_i + \frac{1}{2} \ell''(\hat{s}_i) \Delta s_i^2$$

$$\ell'(s_i) = y_i - p_i = e_i$$
$$\ell''(s_i) = -p_i(1-p_i) = -w_i$$

$$\ell(s_i) \doteq \ell(\hat{s}_i) + \ell'(\hat{s}_i) \Delta s_i + \frac{1}{2} \ell''(\hat{s}_i) \Delta s_i^2$$

$\hat{e}_i$ $\qquad -\hat{w}_i$

surrogate

Taylor Expansion

2nd order approximation $\quad \ell(s_i)$
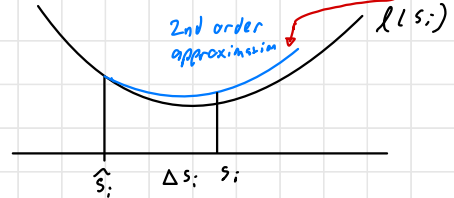
$\hat{s}_i \qquad \Delta s_i \quad s_i$

$$\doteq \text{const} + \hat{e}_i \Delta s_i - \frac{1}{2} \hat{w}_i \Delta s_i^2$$

$$-\ell(s_i) = \frac{1}{2} \hat{w}_i \left( \Delta s_i^2 - 2 \frac{\hat{e}_i}{\hat{w}_i} \Delta s_i \right) + \text{const}$$

$$-\ell(s_i) = \tfrac{1}{2}\,\hat{w_i}\left(\Delta s_i^2 - 2\,\frac{\hat{e_i}}{\hat{w_i}}\,\Delta s_i\right) + const$$

$$-\ell(s_i) \overset{\text{square form}}{=} \tfrac{1}{2}\,\hat{w_i}\left(\Delta s_i - \frac{\hat{e_i}}{\hat{v_i}}\right)^2 + const$$

$$\underbrace{\hphantom{\frac{\hat{e_i}}{\hat{v_i}}}}_{\hat{y_i}}$$

$$= \tfrac{1}{2}\,\hat{w_i}\left(\hat{y_i} - x_i^T \Delta\beta\right)^2$$

$$\text{So} \quad \min_{\Delta\beta} \ -\sum_{i=1}^{n}\ell(s_i)$$

$$\Rightarrow \Delta\beta = \left(\sum \hat{w_i}\, x_i\, x_i^T\right)^{-1}\left(\sum_{i=1}^{n}\hat{w_i}\, x_i\, \hat{y_i}\right)$$
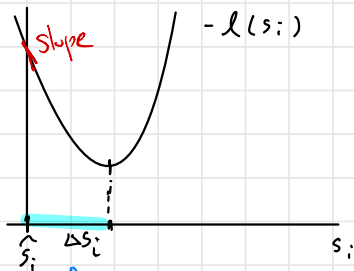
- This is a special case of the New-ton - Raphson method.

- Note $w_i$ measures the curvature since
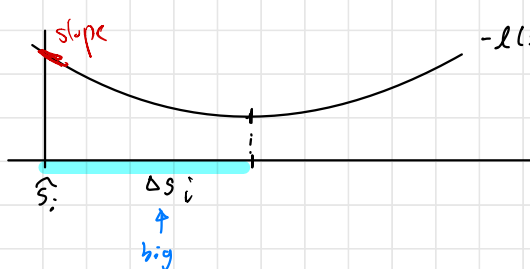
$$\ell''(s_i) = -P_i(1-P_i) = -w_i$$

and is maximized when $P_i = \tfrac{1}{2}$ (uncertain examples)

• $w_i$ big $\Rightarrow$



slope    $-\ell(s_i)$

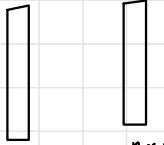$\hat{s_i}$   $\Delta s_i$       $s_i$

↑ small

golf last shot
care about accuracy

• $w_i$ small $\Rightarrow$



slope           $-\ell(s_i)$

$\hat{s_i}$    $\Delta s_i$

↑ big

golf first shot

- Review of Multivariate Calculus:

$$y = F(x)$$

$m \times 1$    $n \times 1$

$$F'(x) = \frac{\partial Y}{\partial x^T}\bigg|_{m \times n} = \left( \partial y_i \; \begin{matrix} 1 \\ \vdots \\ i \\ \vdots \\ m \end{matrix} \right)_{m \times 1} \left[ \begin{matrix} 1 & \cdots & j & \cdots & n \\ & & \frac{1}{\partial x_j} & & \end{matrix} \right]_{1 \times n}$$

$$= \begin{matrix} j \\ i \\ m \end{matrix} \left( \begin{matrix} & & n \\ & \frac{\partial y_i}{\partial x_j} & \\ & & \end{matrix} \right)$$

- example: $\quad y = Ax$

$y_i$   $m \times 1$    $a_{ij}$   $m$    $n$   $x_j$   $n \times 1$

$$y_i = \sum_{j=1}^{n} a_{ij} x_j$$

$$\frac{\partial y}{\partial x^T} = \begin{pmatrix} & j & \\ i & a_{ij} & \\ & & \end{pmatrix} = A$$
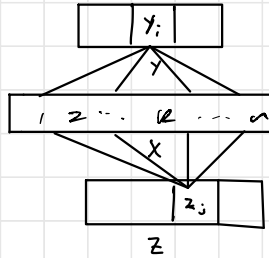
- Chain Rule:

$$y = F(x) \qquad x = G(z)$$
$$m \times 1 \qquad n \times 1 \qquad n \times 1 \qquad \ell \times 1$$

$$\frac{\partial y}{\partial z^T} = \frac{\partial y}{\partial x^T} \frac{\partial x}{\partial z^T}$$
$$m \times \ell \qquad m \times n \qquad n \times \ell$$



$$\frac{\partial y}{\partial z_j} = \sum_{k=1}^{n} \frac{\partial y_i}{\partial x_k} \frac{\partial x_k}{\partial z_j}$$

$$i \left( \frac{\partial y_i}{\partial z_j} \right) = i \left( \frac{\partial y_i}{\partial x_k} \right) \left( \frac{\partial x_k}{\partial z_j} \right)$$

- So consider: $L(\beta) = \frac{1}{2}|e|^2$

$$\frac{\partial L}{\partial \beta^T} = \frac{\partial L}{\partial e^T} \frac{\partial e}{\partial \beta^T}$$

Note: $\frac{\partial L}{\partial e} = \left( \frac{\partial L}{\partial e_i} \right) = \left( e_i \right) = e$

Note: $e = y - x\beta$ so $\frac{\partial L}{\partial \beta^T} = -x$

Thus $\frac{\partial L}{\partial \beta^T} = -e^T x$ 

$$\frac{\partial L}{\partial \beta} = -x^T e = 0$$
$$-x^T(y - x\beta) = 0$$
$$x^T y - x^T x \beta = 0$$
$$(x^T x)^{-1} x^T y = \hat{\beta}$$

○    $y \underset{1 \times 1}{} = F(x) \underset{n \times 1}{}$

<span style="color:red">↓ multidimensional</span>

$F' \underset{n \times 1}{} = F'(x) = \dfrac{\partial y}{\partial x} \underset{n \times 1}{}$    <span style="color:red">instead of</span>   $1 \times n$    }   <span style="color:blue">Gradient</span>

$F''(x) = \dfrac{\partial F'}{\partial x^T} \underset{n \times n}{} = \dfrac{\partial \frac{\partial y}{\partial x}}{\partial x^T} = \dfrac{\partial^2 y}{\partial x \, \partial x^T}$   }   <span style="color:blue">Hessian</span>

$\left( \dfrac{\partial}{\partial x} \right) \left( \dfrac{\partial}{\partial x^T} \right) \, y$

$$\begin{pmatrix} \vdots \\ \frac{\partial}{\partial x_i} \\ \vdots \end{pmatrix} \left( \cdots \quad \frac{\partial}{\partial x_j} \quad \cdots \right) y = \begin{matrix} j \\ i \end{matrix} \begin{pmatrix} & & \vdots & & \\ & & \vdots & & \\ - & - & \frac{\partial^2 y}{\partial x_i \partial x_j} & \cdots & - \\ & & \vdots & & \end{pmatrix}$$

● Consider quadratic form :

$$y = x^T \underset{n \times n}{A} \underset{n \times 1}{x}$$

$$= \sum_{i, j} a_{ij} \, x_i x_j \;=\; const \,+\, a_{ii} x_i^2 \,+\, \sum_{j \neq i} a_{ij} \, x_i x_j \,+\, \sum_{j \neq i} a_{ji} \, x_i x_j$$

$\dfrac{\partial y}{\partial x}$

$$\begin{matrix} 1 \\ \vdots \\ \vdots \\ n \end{matrix} \begin{pmatrix} \\ \frac{\partial y}{\partial x_i} \\ \\ \end{pmatrix} \qquad \dfrac{\partial y}{\partial x_i} = 2 a_{ii} x_i \,+\, \sum_{j \neq i} (a_{ij} + a_{ji}) \, x_j =$$

- Consider quadratic form :

$$y = x^T \underset{n \times n}{A} \underset{n \times 1}{x}$$

$$= \sum_{i,j} a_{ij} x_i x_j = const + a_{ii} x_i^2 + \sum_{j \neq i} a_{ij} x_i x_j + \sum_{j \neq i} a_{ji} x_i x_j$$

$$\frac{\partial y}{\partial x}$$

$$\begin{matrix} 1 \\ \vdots \\ \vdots \\ n \end{matrix} \left( \frac{\partial y}{\partial x_i} \right) \qquad \frac{\partial y}{\partial x_i} = 2 a_{ii} x_i + \sum_{j \neq i} (a_{ij} + a_{ji}) x_j =$$

$$= i \left( \overline{\phantom{aaa} a_{ij} + a_{ji} \phantom{aaa}} \right) \left( \begin{matrix} x_j \end{matrix} \right)$$

So $$\frac{\partial y}{\partial x} = (A + A^T) X$$
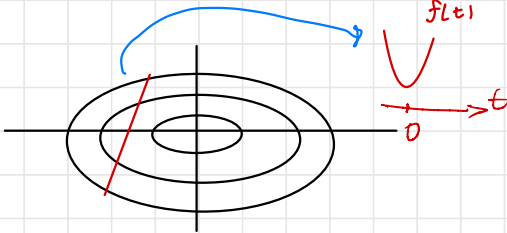
$$\&$$

$$\frac{\partial^2 y}{\partial x \, \partial x^T} = A + A^T$$

- Taylor Expansion

$$y_{1\times 1} = F(x) = F(x_0) + \langle F'(x_0), x-x_0 \rangle + \frac{1}{2}(x-x_0)^T F''(x_0)(x-x_0)$$

with dimension labels: $y$ is $1\times 1$, $F(x)$ is $n\times 1$, $F'(x_0)$ is $n\times 1$, $x-x_0$ is $n\times 1$, $(x-x_0)^T$ is $1\times n$, $F''(x_0)$ is $n\times n$, $(x-x_0)$ is $n\times 1$



$\widehat{F}(x_1, x_2)$

Contour Plot :



$f(t)$

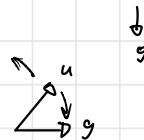$$X = x_0 + \vec{u}\, t \quad {}_{1\times 1}$$

$$|\vec{u}| = 1$$

$$f(t) = F(x) = F(x_0 + \vec{u}\,t)$$

$$\doteq f(0) + f'(0)\,t + \frac{1}{2}f''(0)\,t^2$$

note : $f'(t) = \dfrac{\partial y}{\partial t} = \dfrac{\partial y}{\partial x^T}\dfrac{\partial x}{\partial t} = F'(x)^T u$

$\qquad\qquad\qquad\qquad = \langle F'(x), u\rangle$

$$f'(0) = \langle F'(x_0), u\rangle = \langle g, u\rangle$$
$$\downarrow$$
$$g$$

Implies $f'(0)$ is maximized
when $u$ aligned with $g$
$g$ is the gradient
the steepest direction.

$$= |g||u|\cos(\theta)$$
$$= |g|\cos(\theta)$$
$$\overset{max}{\Rightarrow} u \propto g$$

- note: $f''(t) = \dfrac{\partial f'}{\partial t} = \dfrac{\partial}{\partial t} u^T F' = u^T \dfrac{\partial}{\partial t} F'$

$$= u^T \dfrac{\partial F'}{\partial x^T} \dfrac{\partial x}{\partial t} = u^T \dfrac{\partial^2 F}{\partial x \partial x^T} u$$

$$= u^T F''(x) u$$

$$f''(0) = u^T F''(x_0) u = u^T H u$$

Going back to $f(t) \doteq f(0) + f'(0) t + \frac{1}{2} f''(0) t^2$
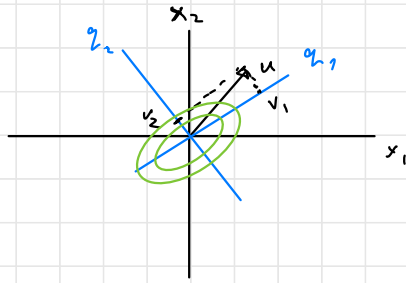so we have

$$f(t) \doteq F(x_0) + \langle F'(x_0), x - x_0 \rangle + \frac{1}{2}(x - x_0)^T F''(x_0)(x - x_0)$$

$$H = Q \Lambda Q^T$$

$$u^T H u = \underbrace{u^T Q}_{v^T} \Lambda \underbrace{Q^T u}_{v} = v^T \Lambda v = \sum_{i=1}^{n} \lambda_i v_i^2$$
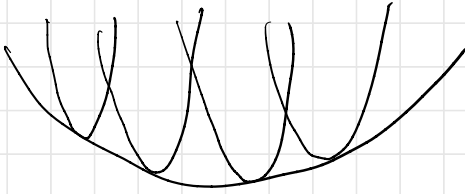
Change of Basis $\begin{cases} v = Q^T u \\ u = Q v \end{cases}$

$$u = \begin{pmatrix} \vec{q}_1 & \cdots & \vec{q}_i & \cdots & \vec{q}_n \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ \vdots \\ v_n \end{pmatrix} = \vec{q}_1 v_1 + \cdots + \vec{q}_i v_i + \cdots + \vec{q}_n v_n$$
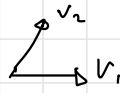
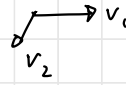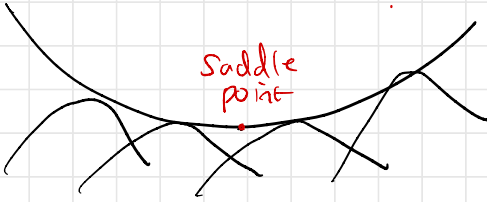$$\lambda_i \equiv \text{curvature along } \vec{q_i}$$

$$\begin{pmatrix} v_i \end{pmatrix} = \begin{pmatrix} q_i^T \, u \end{pmatrix} = \begin{pmatrix} \langle u, \vec{q_i} \rangle \end{pmatrix}$$

$$\lambda_1 v_1^2 + \lambda_2 v_2^2$$

$$\lambda_1 > 0 \qquad \lambda_2 > 0$$
$$(\lambda_1 = 1) \qquad (\lambda_2 = 10)$$

Saddle point

$$\lambda_1 > 0 \qquad \lambda_2 < 0$$
$$(\lambda_1 = 1) \qquad (\lambda_2 = -5)$$
$$\text{curve up} \qquad \text{curve down}$$