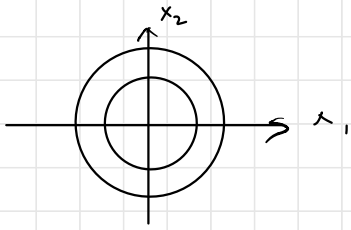


Lecture 5

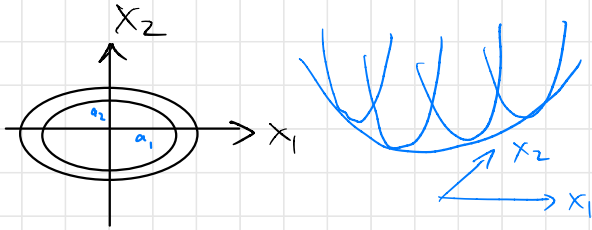


• Circle :



$$x_1^2 + x_2^2 = \text{const}$$

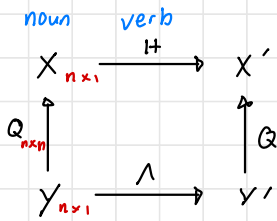
Ellipse :



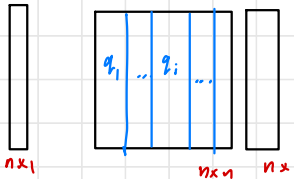
$$\frac{x_1^2}{a_1^2} + \frac{x_2^2}{a_2^2} = \text{const}$$

$$f(x) = \lambda_1 x_1^2 + \lambda_2 x_2^2$$

• Change of view :



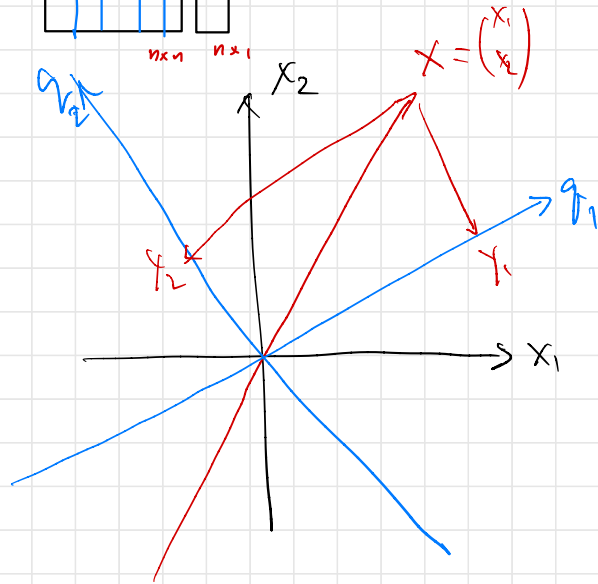
$$x = Q Y$$



$$x = (q_1, \dots, q_i, \dots, q_n) \begin{pmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{pmatrix}$$

$$x = q_1 y_1 + \dots + q_i y_i + \dots + q_n y_n$$

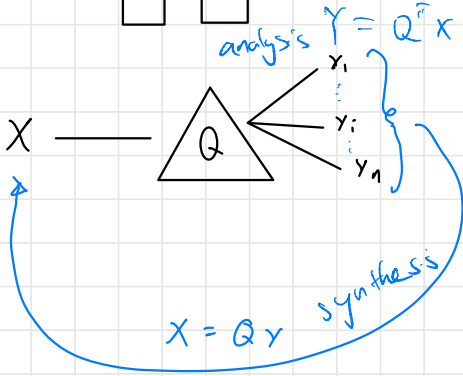
$$S_0 \langle q_i, q_j \rangle = \delta_{ij} = \begin{cases} 1, & i=j \\ 0, & i \neq j \end{cases}$$



y_i coordinate

$$y_i = \langle \begin{bmatrix} x_i \\ \vdots \end{bmatrix}, \begin{bmatrix} q_i \\ \vdots \end{bmatrix} \rangle = q_i^T x$$

$$Y = \begin{pmatrix} y_0 \\ \vdots \end{pmatrix} = \begin{pmatrix} q_0^T x \\ \vdots \end{pmatrix} = \begin{pmatrix} q_0^T \\ \vdots \end{pmatrix} x = Q^T x$$



• Note for any symmetric matrix H , we have

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

Consider $x' = Hx$

$$Qy' = HQY$$

$$\begin{aligned} Y' &= Q^{-1}HQY \\ &= Q^T HQY \\ &= Q^T Q \Lambda Q^T Q Y \\ &= \Lambda Y \end{aligned}$$

Note:

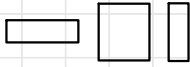
$$\begin{aligned} \langle x, x' \rangle &= |x| |x'| \cos \theta \\ &= |Y| |Y'| \cos \theta \end{aligned}$$

Rotation doesn't change length or angle

$$x^T x' = x^T H x = Y^T Y' = Y^T \Lambda Y$$

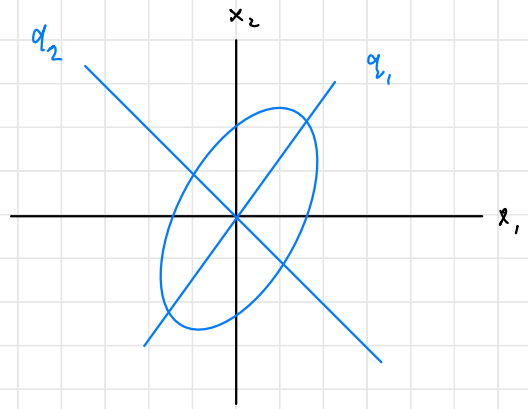
• Quadratic Form

$$f(x) = x^T H x$$



$$= y^T \Lambda y$$

$$= \sum_{i=1}^n \lambda_i y_i^2, \quad \lambda_i > 0$$



• Taylor Expansion:

$$F(x)$$

$$x = x_0 + \vec{u} t$$

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

$$= f(0) + f'(0)t + f''(0)t^2 + o(t^2)$$

small \circ
↑

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

↓
g

↓
H

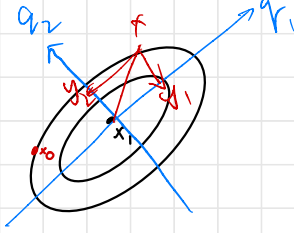
$$= F(x_0) + \langle g, x - x_0 \rangle + \frac{1}{2} (x - x_0)^T H (x - x_0) + o(|x - x_0|^2)$$

$$= \text{const} + \frac{1}{2} (x - x_0 + H^{-1}g)^T H (x - x_0 + H^{-1}g) \quad \text{complete the square}$$

$$= \text{const} + \frac{1}{2} (x - x_1)^T H (x - x_1)$$

$$\text{where } x_1 = x_0 - H^{-1}g = x_0 - F''(x_0)^{-1} F'(x_0)$$

• So the shape is given by

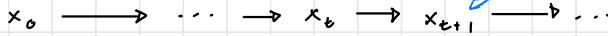


$y = Q^T(x - x_1)$, x_1 is a minimum since we have $\sum \lambda_i y_i^2$, $\lambda_i > 0$

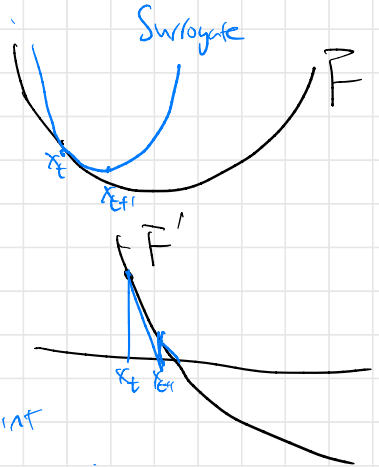
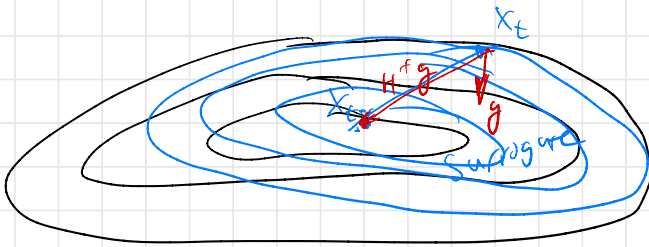
• Newton - Raphson

$\max_x F(x)$

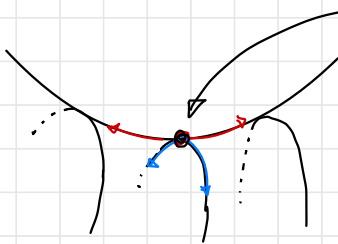
approx $F(x) \approx F(x_e) + \langle \underbrace{\tilde{F}'(x_e)}_{\text{Surrogate}}, \underbrace{(x-x_e)}_{\text{max}} \rangle + \frac{1}{2} (x-x_e)^T \tilde{F}''(x_e) (x-x_e)$



$x_{e+1} = x_e - \underbrace{F''(x_e)^{-1}}_{H^{-1}} \underbrace{F'(x_e)}_g$



saddle point can be problematic



• logistic regression:

	p	
i	x_i^T	y_i
	height weight	gender
n		

$$L(\beta) = \sum_{i=1}^n (y_i s_i - \log(1 + e^{s_i}))$$

$$s_i = x_i^T \beta$$

$$L'(\beta) = \sum_{i=1}^n (x_i y_i - x_i \frac{e^{s_i}}{1 + e^{s_i}})$$

$$= \sum_{i=1}^n x_i (y_i - p_i)$$

$$L''(\beta) = - \sum_{i=1}^n x_i x_i^T \underbrace{p_i (1 - p_i)}_{w_i}$$

$$= - \sum_{i=1}^n w_i x_i x_i^T$$

$$\beta_{t+1} = \beta_t + \left(\sum_{i=1}^n w_i x_i x_i^T \right)^{-1} \left(\sum_{i=1}^n x_i e_i \right) \quad \text{IRLS}$$

Newton-Raphson

• Gradient Descent / Ascent

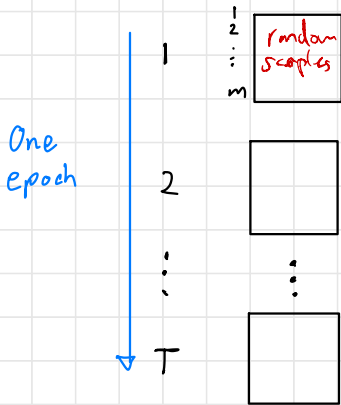
θ : unknown parameter

Loss(θ)

$$\theta_{t+1} = \theta_t - \eta_t \text{Loss}'(\theta_t)$$

↑
step-size learning rate

• Mini-batch:



$$\text{Loss}_b(\theta) = -\frac{1}{m} \sum_{i=1}^m (y_i s_i - \log(1 + e^{s_i}))$$

batch

$$g_t = \text{Loss}'_b(\theta_t)$$

$$= -\frac{1}{m} \sum_{i=1}^m x_i (y_i - p_i)$$

↑
 θ_t

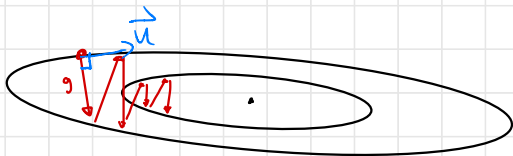
Stochastic Gradient Descent

$$\theta_{t+1} = \theta_t - \eta_t g_t$$

• Issues with SGD:

- consider the following loss function

$$f'(\theta) = \langle g, \vec{u} \rangle$$



- Think Newton-Raphson

$$\begin{array}{ccc} g & \xrightarrow{H^{-1}} & H^{-1}g \\ \uparrow \mathcal{R} & & \uparrow \mathcal{Q} \\ \tilde{g} & \xrightarrow{A^{-1}} & A^{-1}\tilde{g} = \begin{pmatrix} \tilde{g}_1 \\ \vdots \\ \tilde{g}_n \end{pmatrix} \end{array}$$

$$\text{GD: } \theta_{t+1} = \theta_t - \eta_t L'(\theta_t) \quad H = \frac{1}{\eta_t} I$$

$$\text{Newton: } \theta_{t+1} = \theta_t - H^{-1} L'(\theta_t)$$

GD: $\theta_{t+1} = \theta_t - \eta_t L'(\theta_t)$ $H = \frac{1}{\eta_t} I$

Objective: $L(\theta) \approx L(\theta_t) + \underbrace{\langle \theta - \theta_t, L'(\theta_t) \rangle}_{\text{Surrogate}} + \frac{1}{2} (\theta - \theta_t)^T \frac{1}{\eta_t} I (\theta - \theta_t)$

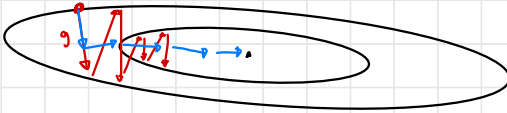
Poor approximation $\rightarrow \frac{1}{2\eta_t} \|\theta - \theta_t\|^2$
 no curvature info

Newton: $\theta_{t+1} = \theta_t - H^{-1} L'(\theta_t)$



• Momentum: $v_t = \gamma v_{t-1} + \eta_t g_t$ heavy ball

$\theta_{t+1} = \theta_t - v_t$ (cancelling oscillation / variance)

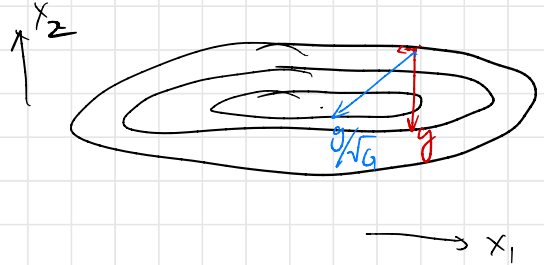


o Adaptive Gradient:

$$G_t = G_{t-1} + g_t^2$$

$$\theta_{t+1} = \theta_t - \eta_t \frac{g_t}{\sqrt{G_{t+2}}}$$

Element-wise operation



• Recall $g_t = -\frac{1}{m} \sum_{i=1}^m x_i (y_i - p_i)$

if $x_{ij} = 0$ for most i (sparse feature)
then $L'(A_j)$ small

→ then once we see $x_{ij} \neq 0$ we take a large step.