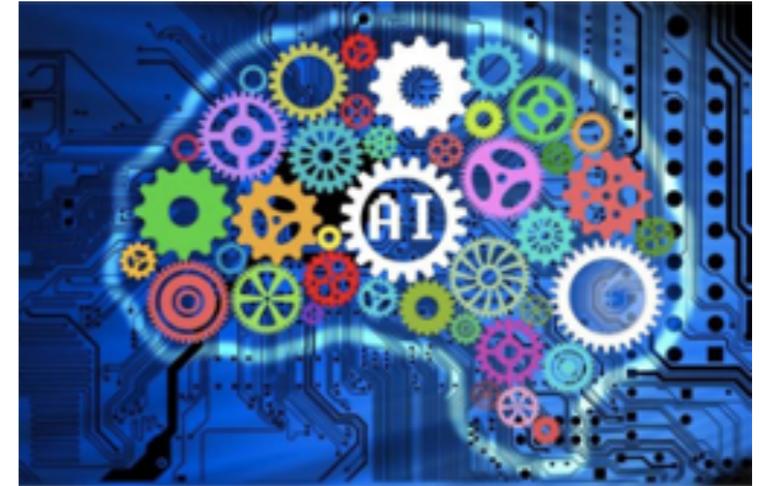


AI & Machine Learning in Physics



PHYS3151 (6 credits)

Time & Place : Tue 13:30-14:20, 14:30-15:20 MB 122
Fri 14:30-15:20 MB 142

Teachers: Zi Yang Meng (zymeng@hku.hk), HOC 231

<https://quantummc.xyz/hku-phys3151-machine-learning-in-physics-2023/>

Tutor: Ting-Tung Wang (leowdd@connect.hku.hk), HOC 217

AI & Machine Learning in Physics

Teaching Materials:

<https://quantummc.xyz/hku-phys3151-machine-learning-in-physics-2023/>

Slides / Reading materials

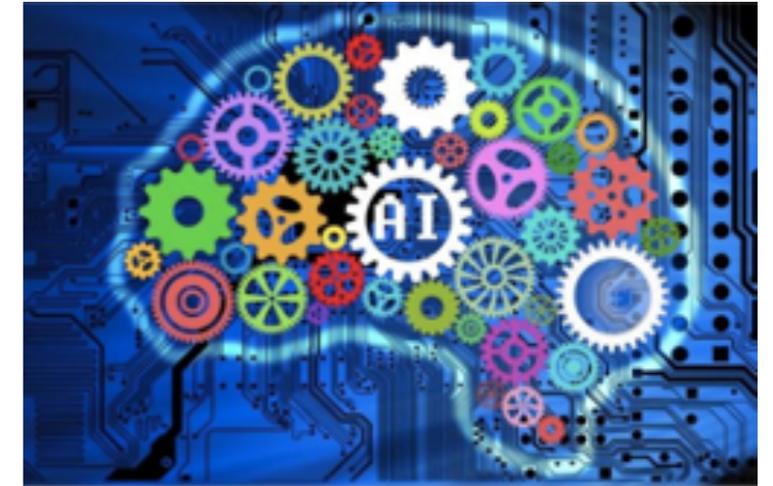
Python notebooks

Assignments

Assessment Methods and Weighting

- Assignments 30%
- Presentation 20%
- Project report 20%
- Exam. 30%

Content



0. Introduction

1. Regression

1.1 Multivariate Linear Regression (curve fitting)

1.2 Regularization (Lagrange multiplier)

1.3 Logistic Regression (Fermi-Dirac distribution)

1.4 Support Vector Machine (high-school geometry)

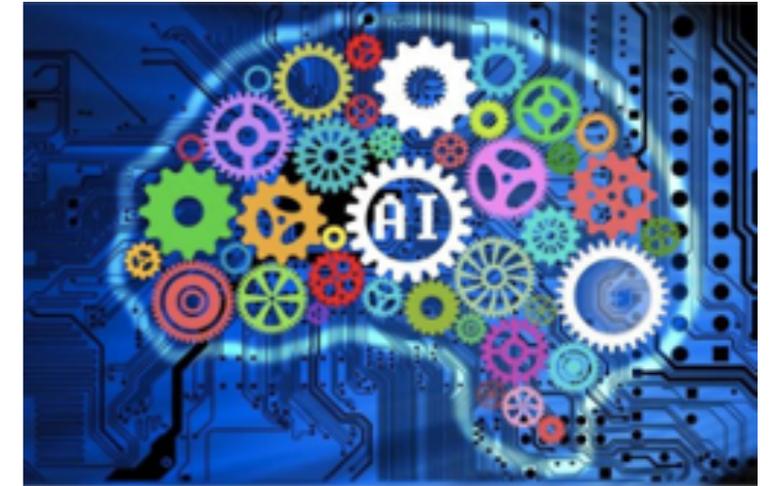
2. Dimensionality Reduction/feature extraction

2.1 Principal Component Analysis (order parameters)

2.2 Recommender Systems

2.3 Clustering (phase transition)

Content



3. Neural Networks

3.1 Biological neural networks

3.2 Mathematical representation

3.3 Factoring biological ingredient

3.4 Feed-forward neural networks

3.5 Learning algorithm

3.6 Universal Approximation Theorem

Multivariate Linear Regression

- Ethem Alpaydin, Introduction to Machine Learning, Third Edition, MIT Press 2014

Chap. 4. Parametric Methods

4.6 Regression

4.7 Bias/Variance Dilemma

Chap. 5. Multivariate Methods

5.1 Multivariate data

5.8 Multivariate Regression

- Linear algebra books:

Gilbert Strang, Linear Algebra and its Applications, 1988

David Harville, Matrix Algebra from a Statistician's Perspective, 1997, Springer

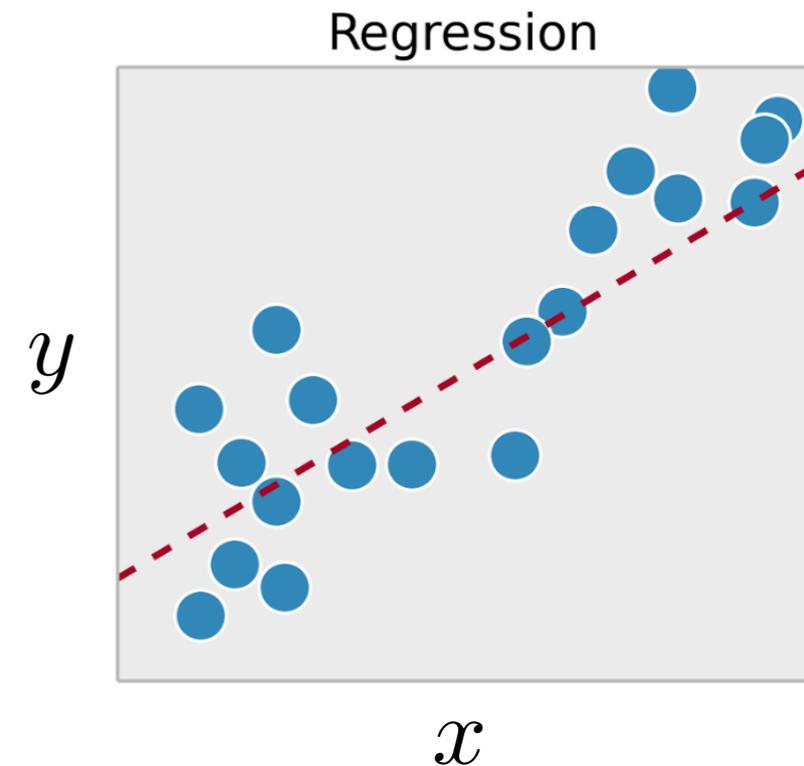
- Most important

<https://quantummc.xyz/hku-phys3151-machine-learning-in-physics-2023/>

Multivariate Linear Regression

Regression: $y = h_{\Theta}(x) = \Theta \cdot x$

Statistics	Machine Learning	Notation	Remarks
independent variable	feature	$x_j^{(i)}$	$j = 1, \dots, N$
dependent variable	outcome	$y^{(i)}$	
sample	example	$(x^{(i)}, y^{(i)})$	$i = 1, \dots, M$
model	hypothesis	$h_{\theta}(x)$	
parameter	parameter	θ_j	
intercept	bias	θ_0	



$$\{(x_j^{(i)}, y^{(i)}), \theta_j\}; j = 1, 2, \dots, N; i = 1, 2, \dots, M; N < M$$

$$y^{(i)} = \theta_0 + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \dots + \theta_N x_N^{(i)}$$

$$\mathbb{R}^{M \times (N+1)} \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_N^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_N^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(M)} & x_2^{(M)} & \dots & x_N^{(M)} \end{bmatrix} \cdot \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_N \end{bmatrix} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(M)} \end{bmatrix}$$

$$\underline{\underline{X}} \cdot \underline{\underline{\Theta}} = \underline{\underline{Y}}$$

Multivariate Linear Regression

Basic assumption: examples are independent and identically distributed (i.i.d.).

Likelihood function $\mathcal{L}(\Theta|X) = \sum_{i=1}^M \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y^{(i)} - \theta^T \cdot x^{(i)})^2}{2\sigma^2}\right)$

$$\Theta^* = \arg \min_{\Theta} \left\{ \sum_{i=1}^M \frac{(y^{(i)} - \theta^T \cdot x^{(i)})^2}{2M} \right\}$$

Cost / Loss function $J(\Theta) = \frac{1}{2M} \sum_{i=1}^M (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Convex optimisation $\Delta_{\theta} J(\theta) = \nabla_{\theta} \cdot \nabla_{\theta} J(\theta)$

$$= \nabla_{\theta} \cdot \nabla_{\theta} \left[\frac{1}{2M} \sum_{i=1}^M \left(\theta^T x^{(i)} - y^{(i)} \right)^2 \right]$$
$$= \nabla_{\theta} \cdot \left[\frac{1}{M} \sum_{i=1}^M \left(\theta^T x^{(i)} - y^{(i)} \right) x^{(i)} \right]$$
$$= \frac{1}{M} \sum_{i=1}^M \left| x^{(i)} \right|^2 \geq 0$$

There exist global minimum, our goal is to find it.

Multivariate Linear Regression

Feature Scaling:
$$x_j^{(i)} := \frac{x_j^{(i)} - \mu_j^x}{\sigma_j^x} \quad \mu_j^x = \frac{1}{M} \sum_{i=1}^M x_j^{(i)} \quad \sigma_j^x = \sqrt{\frac{\sum_{i=1}^M (x_j^{(i)} - \mu_j^x)^2}{M}}$$

Ensure all the features θ_j lie in the same range

$$J(\theta) = \frac{1}{2M} \sum_{i=1}^M \left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right)^2 = \frac{1}{2M} \|X\Theta - Y\|^2$$

$$= \frac{1}{2M} (\Theta^T X^T - Y^T)(X\Theta - Y)$$

$$= \frac{1}{2} \left(\Theta^T \frac{1}{M} X^T X \Theta - 2Y^T X \Theta + Y^T Y \right)$$

$$= \frac{1}{2} \Theta^T Q \Theta - b^T \Theta + c$$

$$Q = \frac{1}{M} X^T X \geq 0$$

$$b = \frac{1}{M} X^T Y$$

$$c = \frac{1}{2M} Y^T Y$$

Paraboloid

Multivariate Linear Regression

Structure of $Q = \frac{1}{M} X^T X$

$$\frac{1}{M} \underbrace{\begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1^{(1)} & x_1^{(2)} & \cdots & x_1^{(M)} \\ x_2^{(1)} & x_2^{(2)} & \cdots & x_2^{(M)} \\ \vdots & \vdots & \ddots & \vdots \\ x_N^{(1)} & x_N^{(2)} & \cdots & x_N^{(M)} \end{bmatrix}}_{\mathbb{R}^{(N+1) \times M}} \cdot \underbrace{\begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \cdots & x_N^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \cdots & x_N^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(M)} & x_2^{(M)} & \cdots & x_N^{(M)} \end{bmatrix}}_{\mathbb{R}^{M \times (N+1)}}$$

$$= \underbrace{\begin{bmatrix} \frac{\sum_{i=1}^M 1}{M} & \frac{\sum_{i=1}^M x_1^{(i)}}{M} & \frac{\sum_{i=1}^M x_2^{(i)}}{M} & \cdots & \frac{\sum_{i=1}^M x_N^{(i)}}{M} \\ \frac{\sum_{i=1}^M x_1^{(i)}}{M} & \frac{\sum_{i=1}^M (x_1^{(i)})^2}{M} & \frac{\sum_{i=1}^M x_1^{(i)} x_2^{(i)}}{M} & \cdots & \frac{\sum_{i=1}^M x_1^{(i)} x_N^{(i)}}{M} \\ \frac{\sum_{i=1}^M x_2^{(i)}}{M} & \frac{\sum_{i=1}^M x_2^{(i)} x_1^{(i)}}{M} & \frac{\sum_{i=1}^M (x_2^{(i)})^2}{M} & \cdots & \frac{\sum_{i=1}^M x_2^{(i)} x_N^{(i)}}{M} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{\sum_{i=1}^M x_N^{(i)}}{M} & \frac{\sum_{i=1}^M x_N^{(i)} x_1^{(i)}}{M} & \frac{\sum_{i=1}^M x_N^{(i)} x_2^{(i)}}{M} & \cdots & \frac{\sum_{i=1}^M (x_N^{(i)})^2}{M} \end{bmatrix}}_{\mathbb{R}^{(N+1) \times (N+1)}} = \underbrace{\begin{bmatrix} \langle 1 \rangle & \langle x_1 \rangle & \langle x_2 \rangle & \cdots & \langle x_N \rangle \\ \langle x_1 \rangle & \langle (x_1)^2 \rangle & \langle x_1 x_2 \rangle & \cdots & \langle x_1 x_N \rangle \\ \langle x_2 \rangle & \langle x_2 x_1 \rangle & \langle (x_2)^2 \rangle & \cdots & \langle x_2 x_N \rangle \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \langle x_N \rangle & \langle x_N x_1 \rangle & \langle x_N x_2 \rangle & \cdots & \langle (x_N)^2 \rangle \end{bmatrix}}_{\mathbb{R}^{(N+1) \times (N+1)}}$$

Symmetric and positive semi-definite

Multivariate Linear Regression

Positive semi-definite $v^T Q v = v^T X^T X v = (Xv)^T (Xv) = u^T u \geq 0$

$$J(\theta) = \frac{1}{2M} \sum_{i=1}^M \left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right)^2 = \frac{1}{2M} \|X\Theta - Y\|^2$$

$$\nabla_{\theta} J(\theta) = \frac{1}{M} \left((X^T X) \theta - (X^T Y) \right)$$

$$Q = \frac{1}{M} X^T X \quad b = X^T Y$$

Normal Equation

$$Q\theta \equiv \frac{1}{M} (X^T X) \theta = \frac{1}{M} X^T Y$$

$$\theta = (X^T X)^{-1} X^T Y = Q^{-1} b$$

$$O((N + 1)^3)$$

Multivariate Linear Regression

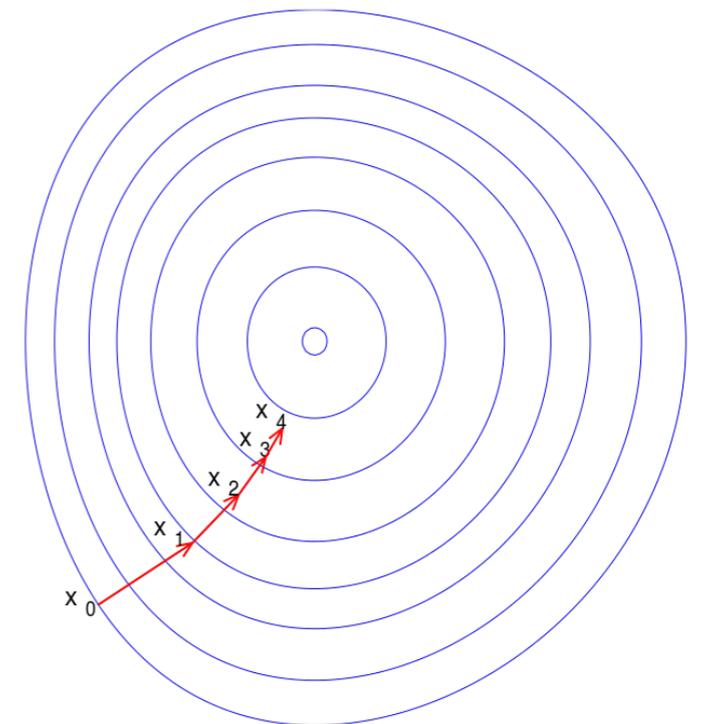
Gradient descent method

Algorithm:

- 1 Set a precision ϵ , learning rate α , and set the initial guess θ_0
- 2 Let $\theta_{j+1} := \theta_j - \alpha \nabla_{\theta} J(\theta_j)$
- 3 Calculate $J(\theta_{j+1})$
 - If $|\theta_{j+1} - \theta_j| < \epsilon$, stop and return θ_{j+1}
 - Else, return to step 2

The algorithm is based on the Taylor expansion

$$J(\theta_{j+1}) = J(\theta_j - \alpha \nabla_{\theta} J(\theta_j)) = J(\theta_j) - \alpha (\nabla_{\theta} J(\theta_j))^2 + O(\alpha^2)$$



Multivariate Linear Regression

- Deterministic method
 - Normal equation, computational complexity and instability issues
- Stochastic method
 - Gradient Descent
 - Steepest Descent
 - Conjugate Gradient

Good reference:

Jonathan Richard Shewchuk

An Introduction to Conjugate Gradient Method Without the Agonizing Pain

<http://www.cs.cmu.edu/~quake-papers/painless-conjugate-gradient.pdf>



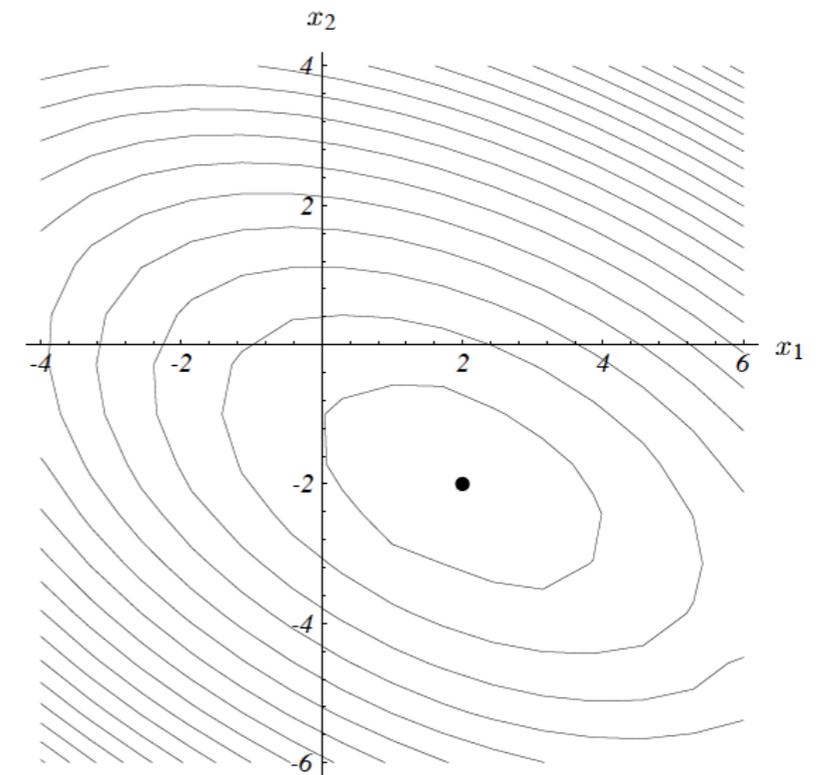
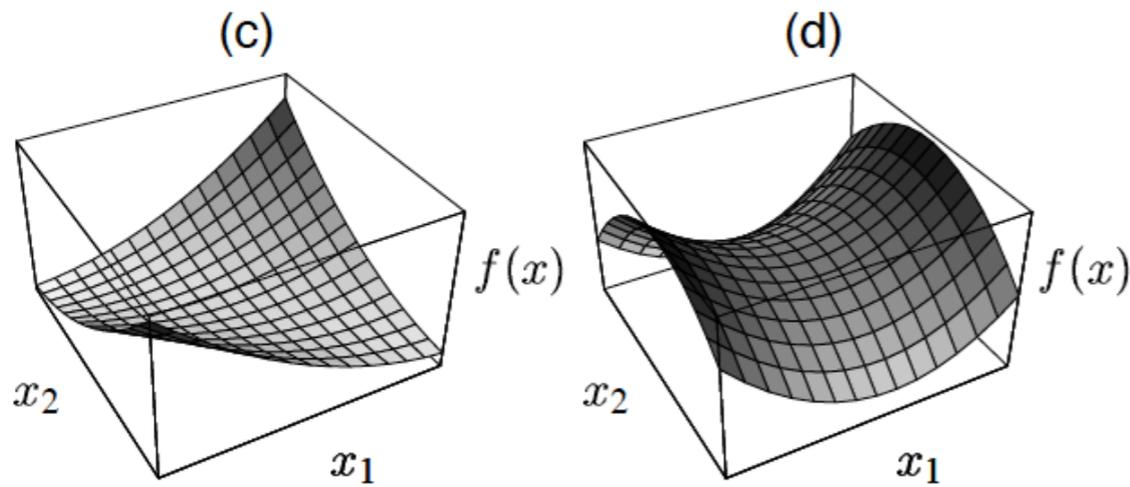
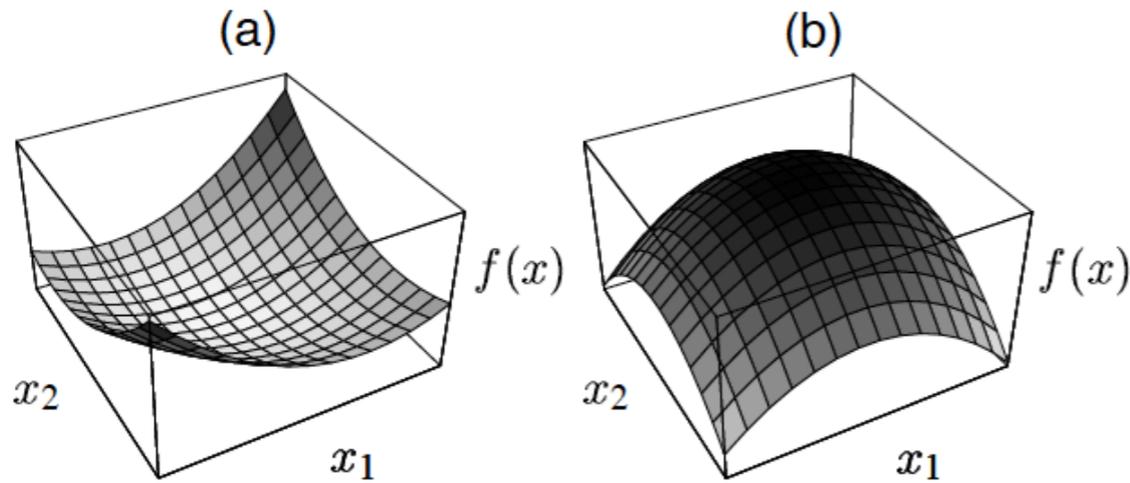
“I am trained to only sleep during national holidays.”

Paraboloid and positive-definite

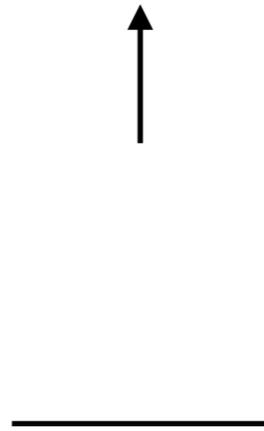
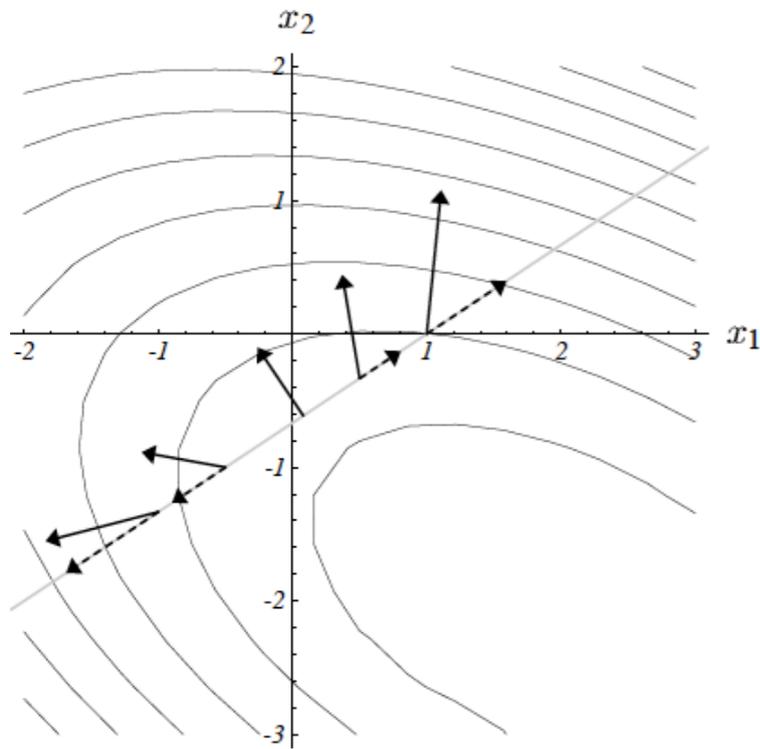
$$f(x) = \frac{1}{2}x^T Ax - b^T x + c \quad \nabla_x f(x) = Ax - b = 0 \quad Ax = b \quad A = \frac{1}{M}X^T X \quad \mathbb{R}^{(N+1) \times (N+1)}$$

$$b = X^T Y \quad \mathbb{R}^{(N+1) \times 1}$$

$$A = \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix} \quad b = \begin{bmatrix} 2 \\ -8 \end{bmatrix} \quad c = 0$$

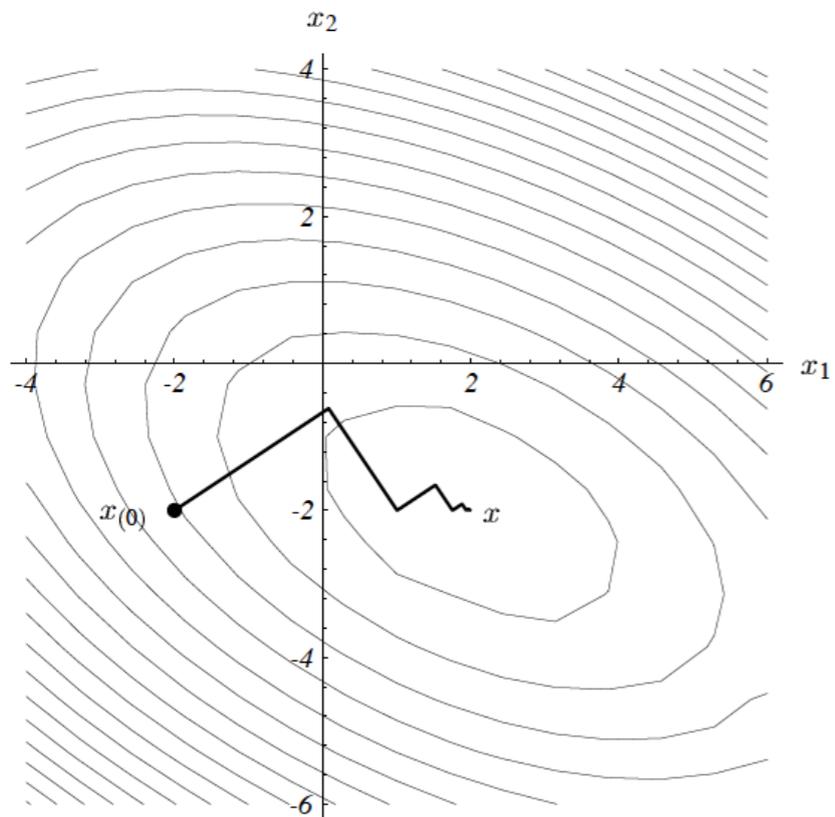


Steepest Descent



gradient

search line,
 $f(x)$ is minimised where the gradient is
orthogonal to the search line



Steepest descent

The search direction is orthogonal, but this is not
sufficient,
The searching direction needs to be A-orthogonal

Steepest Descent

$$A = \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix} \quad b = \begin{bmatrix} 2 \\ -8 \end{bmatrix} \quad c = 0$$

$$x_{(0)} = [-2, -2]^T$$

Residual $r_{(i)} = b - Ax_{(i)} \quad r_i = -\nabla_{x_{(i)}} f(x)$

Learning rate $\alpha_{(i)} = \frac{r_{(i)}^T r_{(i)}}{r_{(i)}^T A r_{(i)}}$

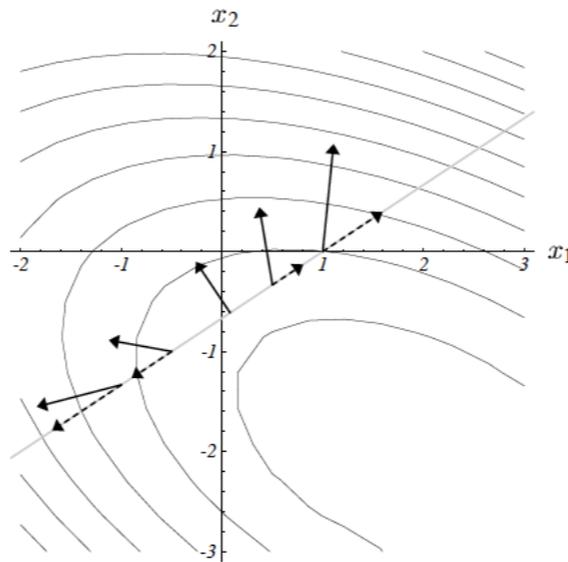
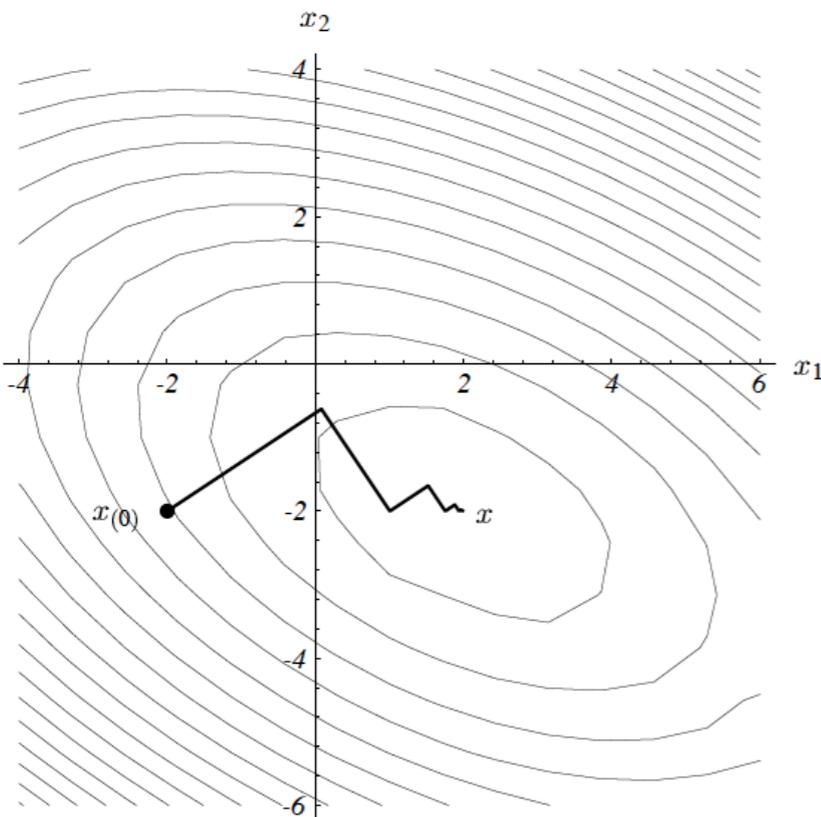
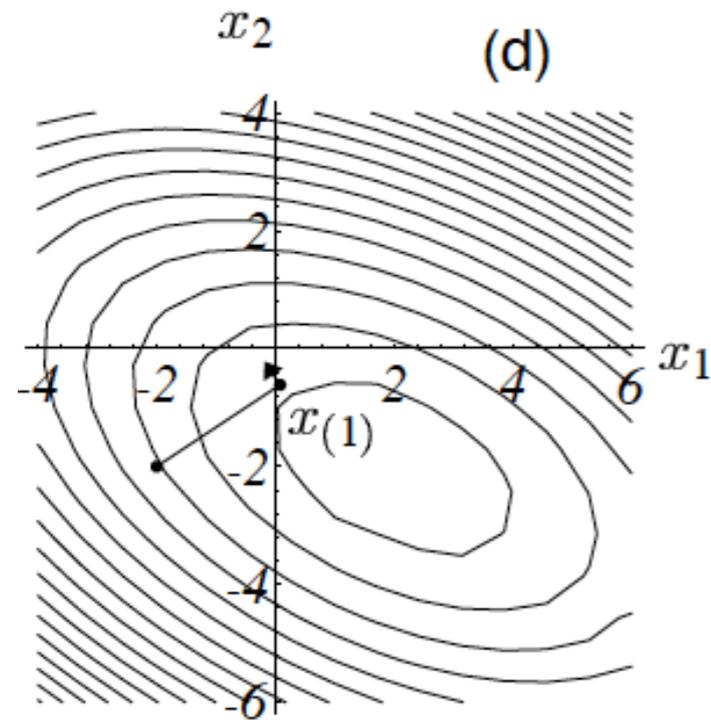
$$r_{(0)} = [12, 8]^T$$

$$r_{(0)}^T r_{(0)} = 208$$

$$r_{(0)}^T A r_{(0)} = 1200$$

Iteration $x_{(i+1)} = x_{(i)} + \alpha_{(i)} r_{(i)}$

$$x_{(1)} = [-2, -2]^T + \frac{208}{1200} [12, 8]^T = [0.08, -0.62]^T$$



$$r_{(i+1)}^T r_{(i)} = 0$$

$$(b - Ax_{(i+1)})^T r_{(i)} = 0$$

$$(b - A(x_{(i)} + \alpha r_{(i)}))^T r_{(i)} = 0$$

$$(b - Ax_{(i)})^T r_{(i)} = \alpha (Ar_{(i)})^T r_{(i)}$$

$$r_{(i)}^T r_{(i)} = \alpha r_{(i)}^T A r_{(i)}$$

$$\alpha = \frac{r_{(i)}^T r_{(i)}}{r_{(i)}^T A r_{(i)}}$$

Each gradient is orthogonal to the previous gradient

Let's eigen do it

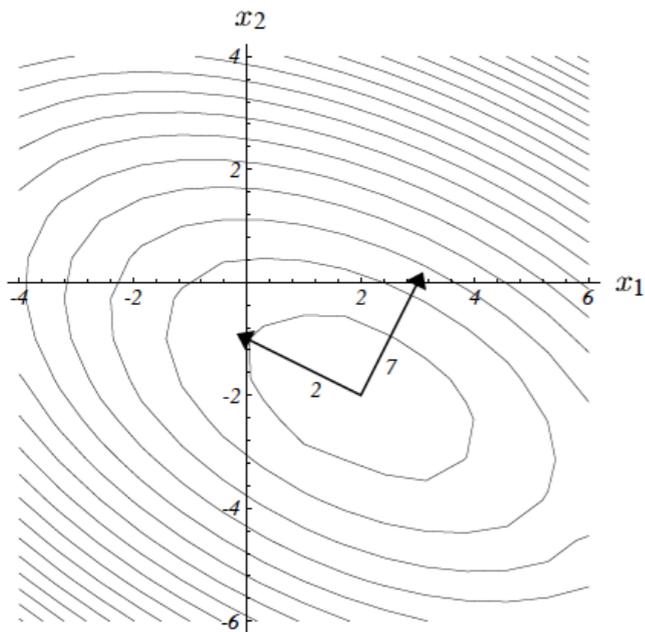
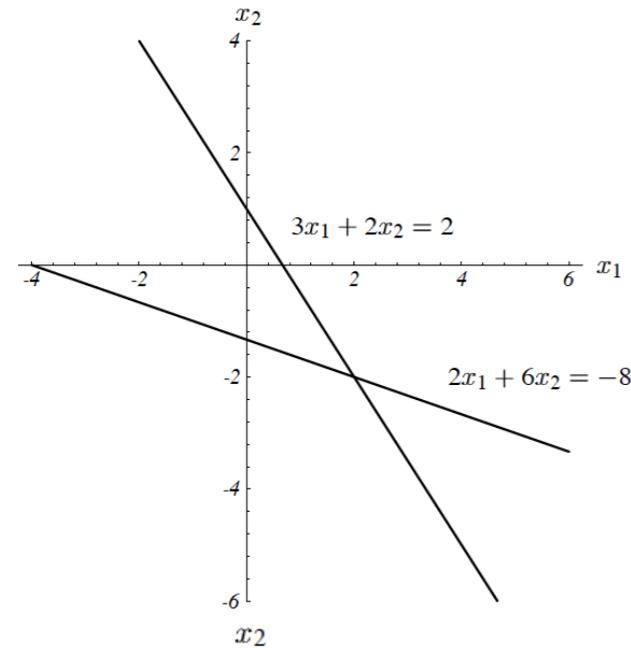
$$f(x) = \frac{1}{2} x^T A x - b^T x + c$$

$$A = \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix} \quad b = \begin{bmatrix} 2 \\ -8 \end{bmatrix} \quad c = 0$$

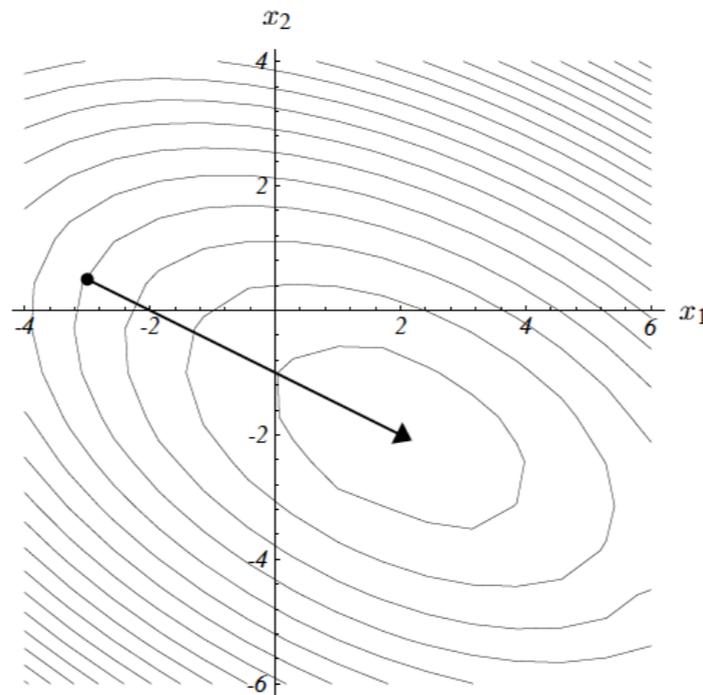
Eigenvalues and eigenvectors:

$$\lambda_1 = 7, \quad e_{(1)} = [1, 2]^T$$

$$\lambda_2 = 2, \quad e_{(2)} = [-2, 1]^T$$



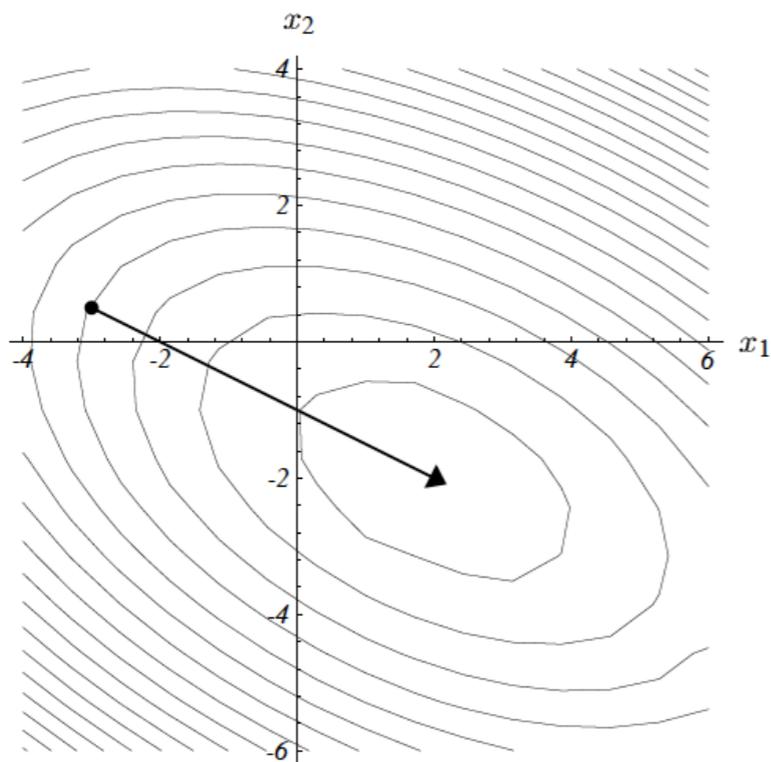
If $r_{(i)}$ is the eigenvector, only one step to converge to the exact solution



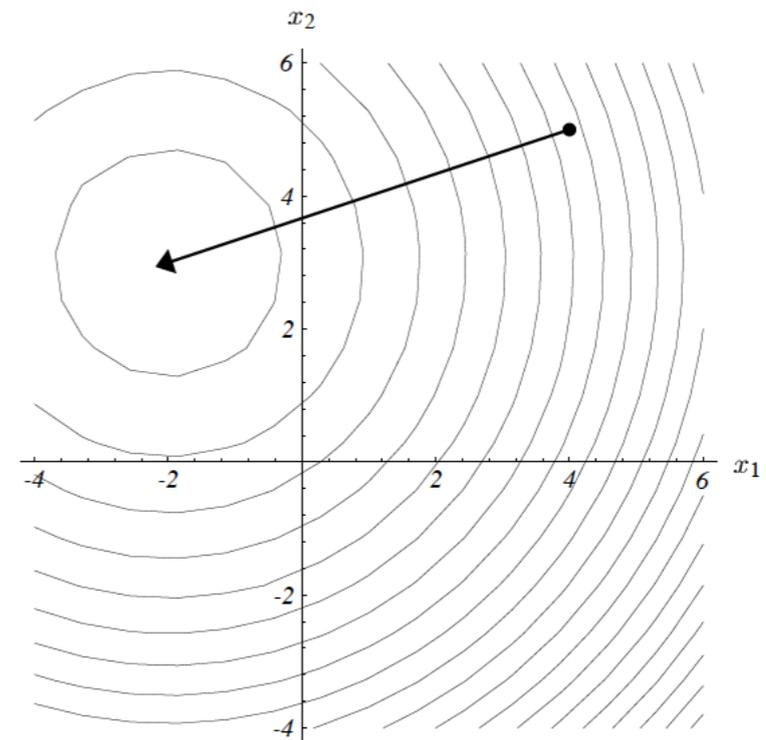
$$A r_{(i)} = \lambda r_{(i)} \quad \alpha = 1/\lambda$$

$$\begin{aligned} r_{(i+1)} &= b - A(x_{(i)} + \alpha r_{(i)}) \\ &= b - A x_{(i)} - \alpha A r_{(i)} \\ &= 0 \end{aligned}$$

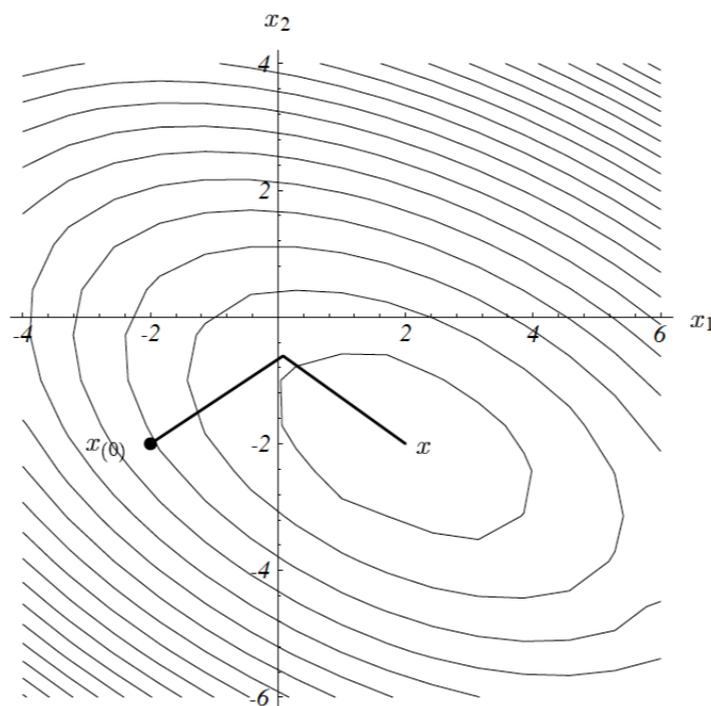
“Eigenvectors are useful tools, and not just bizarre torture devices inflicted on you by your professors for the pleasure of watching you suffer (although the latter is a nice fringe benefit)”



Paraboloid is ellipsoidal, only if r are eigenvectors, can one find the minimal

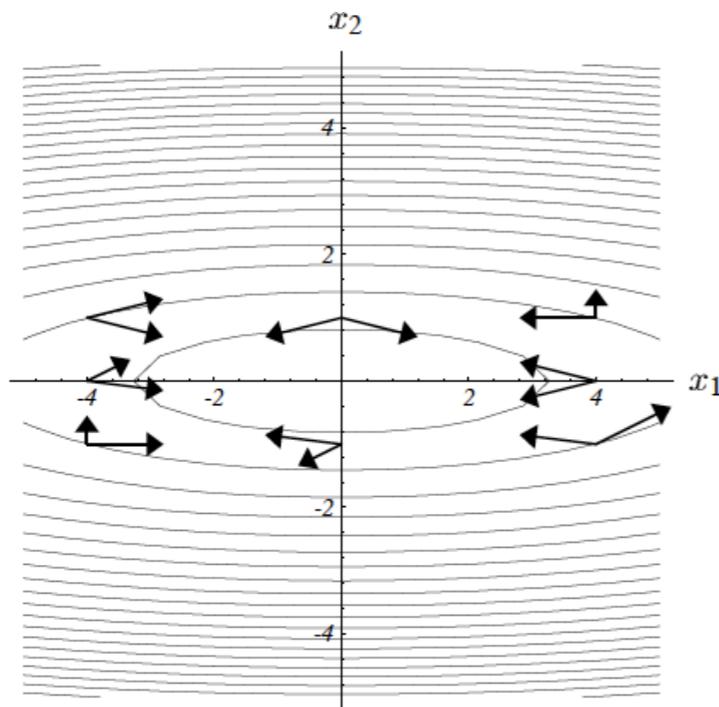
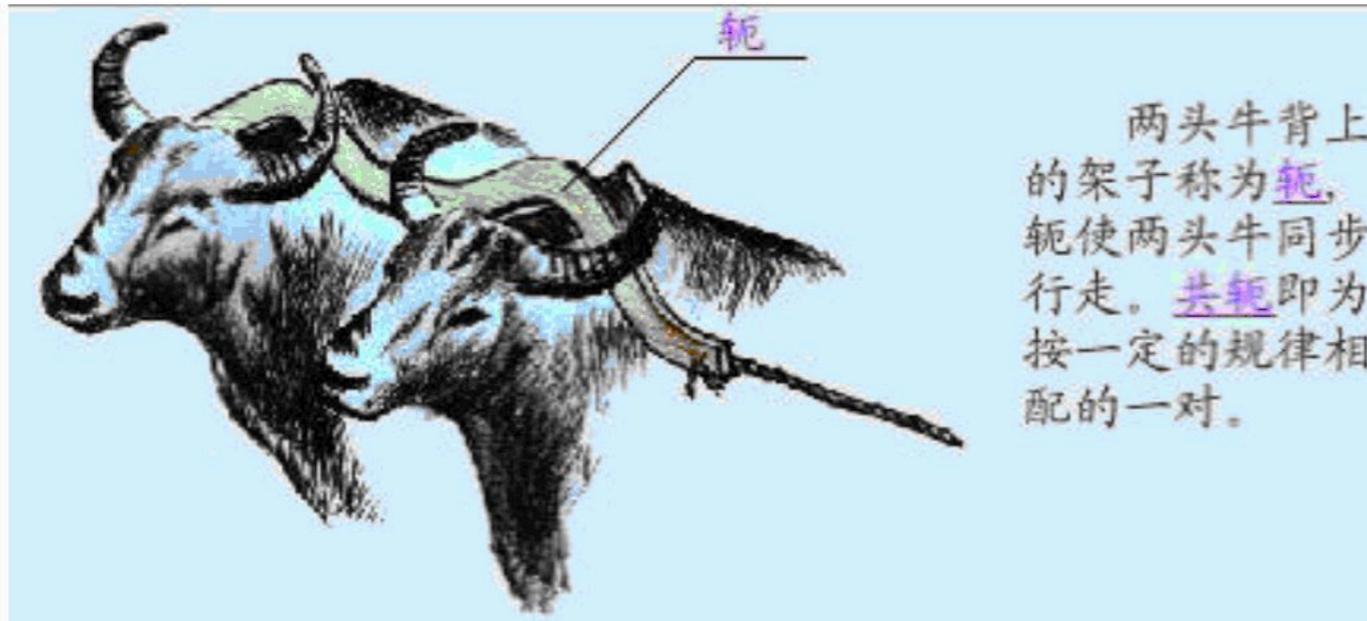


Paraboloid is spherical, no matter what point we start, always find the minimal minimal

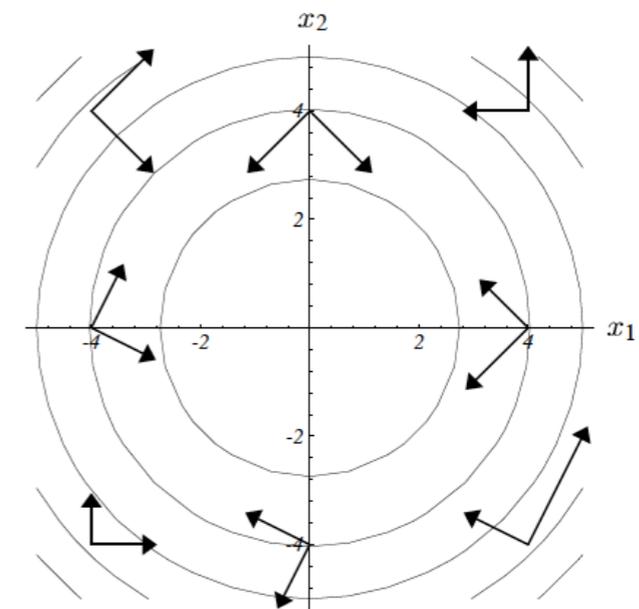


CG is to find the orthogonal directions in a stretched (scaled) space.

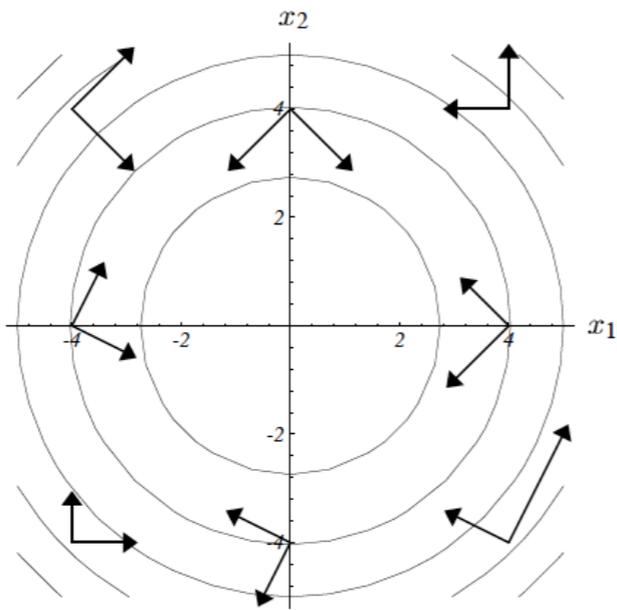
Conjugate Gradient Method



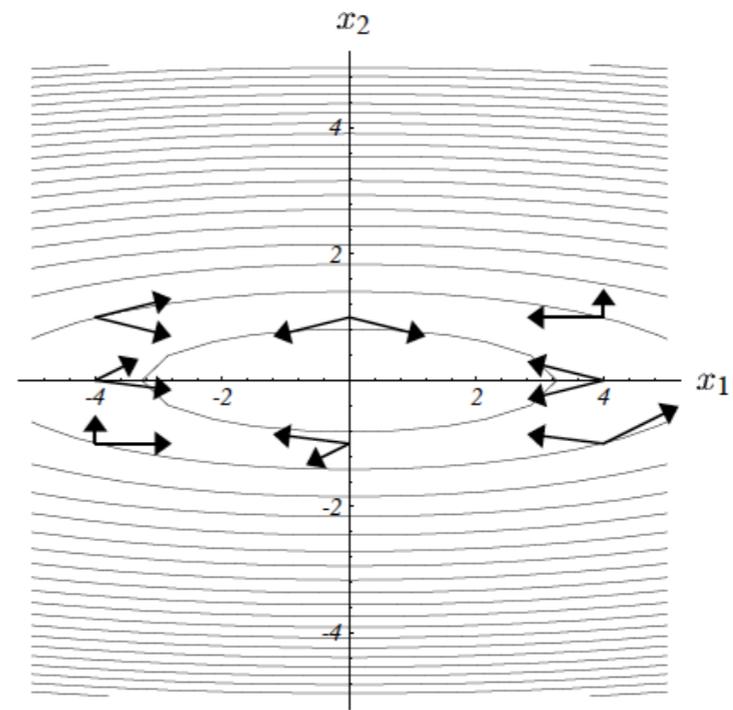
$$d_{(i)}^T A d_{(j)} = 0$$



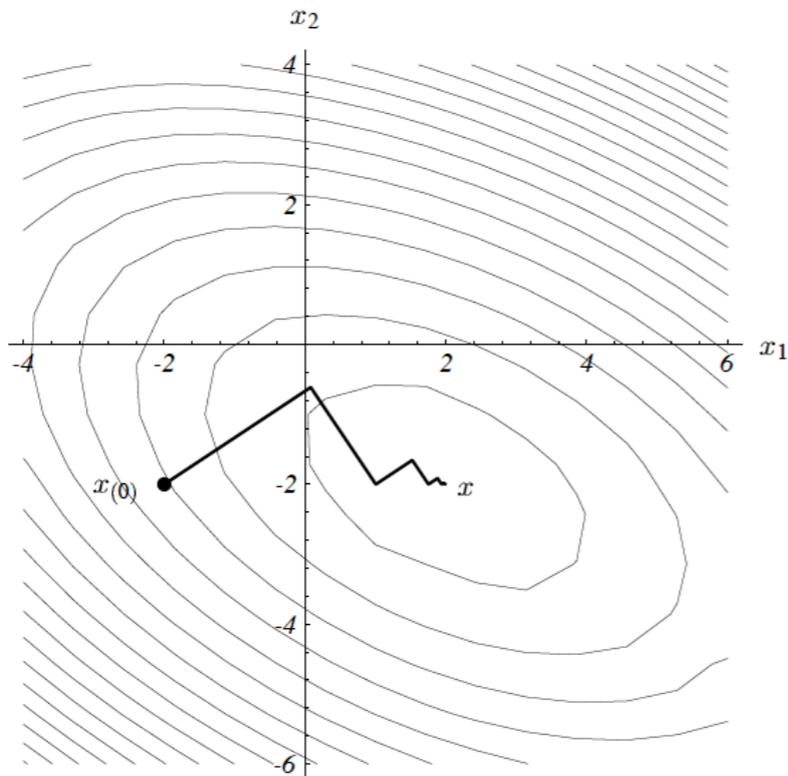
Pairs of vectors that are A-orthogonal, conjugate



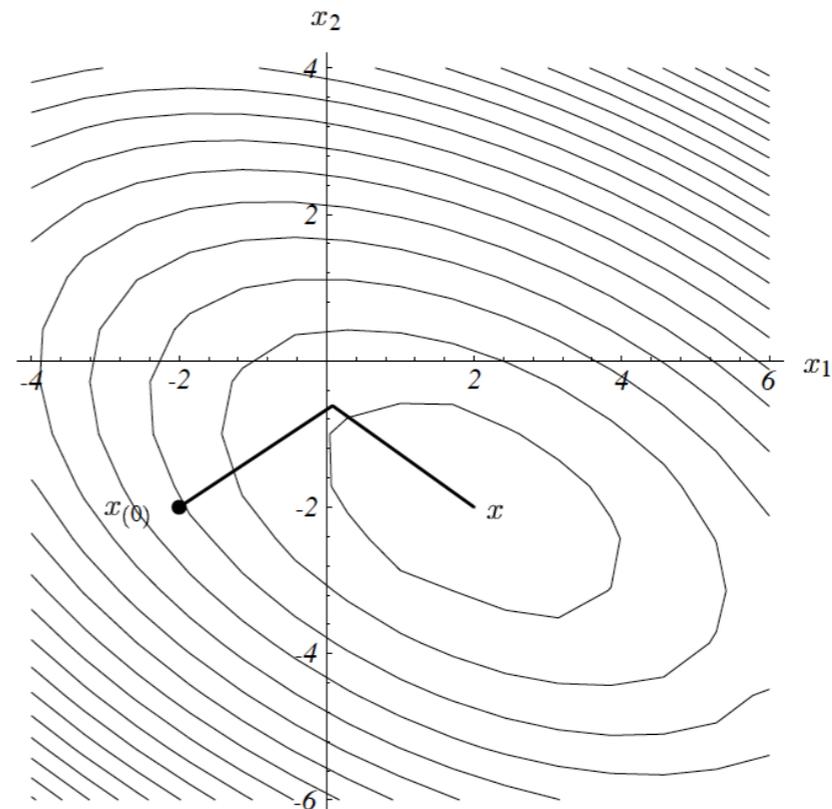
Pairs of vectors that are orthogonal



Pairs of vectors that are A-orthogonal, conjugate



Many steps to find the solution



Only takes N-steps to find the solution

$$x^* = \sum_{i=1}^n \alpha_i p_i$$

a set of n mutually conjugate vectors (with respect to A)

$$P = \{P_1, P_2, \dots, P_n\}$$

One can express the solution

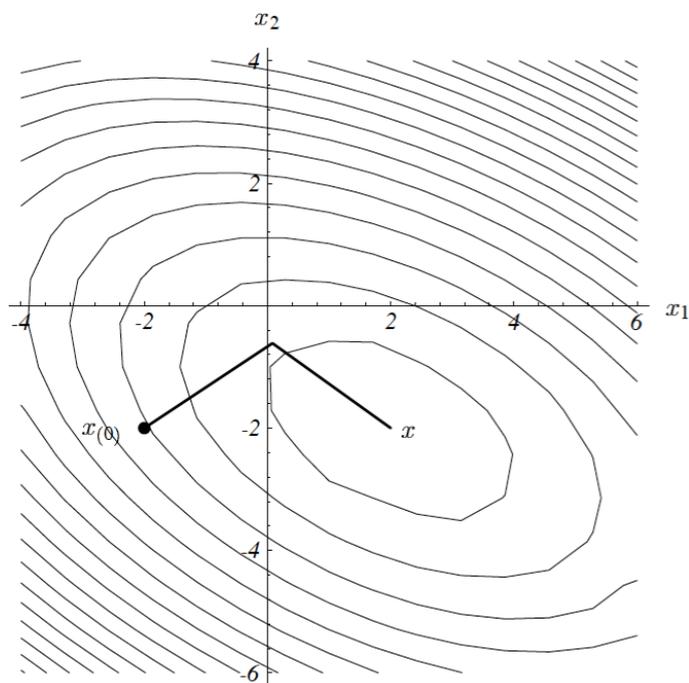
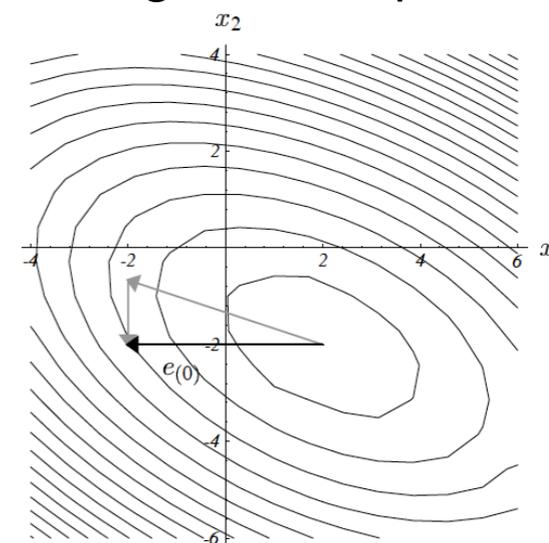
$$x^* = \sum_{i=1}^n \alpha_i P_i$$

Decompose a vector x^* to a sum of A -orthogonal components

$$Ax^* = \sum_{i=1}^n \alpha_i AP_i$$

$$P_k^T Ax^* = \sum_{i=1}^n \alpha_i P_k^T AP_i \quad \forall i \neq k \quad P_k^T AP_i = 0 \quad Ax^* = b$$

$$\alpha_k = \frac{P_k^T b}{P_k^T AP_k}$$



$$r_{(0)} = b - Ax_{(0)} \quad \text{if } r_{(0)} < \epsilon, \text{ return } x_{(0)}$$

$$P_{(0)} = r_{(0)}$$

Repeat from $k=0$

$$\alpha_{(k)} = \frac{r_{(k)}^T r_{(k)}}{P_{(k)}^T AP_{(k)}}$$

$$x_{(k+1)} = x_{(k)} + \alpha_{(k)} P_{(k)}$$

$$r_{(k+1)} = b - Ax_{(k+1)} = b - A(x_{(k)} + \alpha_{(k)} P_{(k)}) = r_{(k)} - \alpha_{(k)} AP_{(k)}$$

Gram-Schmidt conjugation

$$\beta_{(k)} = \frac{r_{(k+1)}^T r_{(k+1)}}{r_{(k)}^T r_{(k)}}$$

if $r_{(k+1)} < \epsilon$, return $x_{(k+1)}$

$$P_{(k+1)} = r_{(k+1)} + \beta_{(k)} P_{(k)}$$

$$k = k + 1$$

One example on the fly

$$r_{(0)} = b - Ax_{(0)}$$

$$P_{(0)} = r_{(0)}$$

Repeat from $k=0$

$$\alpha_{(k)} = \frac{r_{(k)}^T r_{(k)}}{P_{(k)}^T A P_{(k)}}$$

$$x_{(k+1)} = x_{(k)} + \alpha_{(k)} P_{(k)}$$

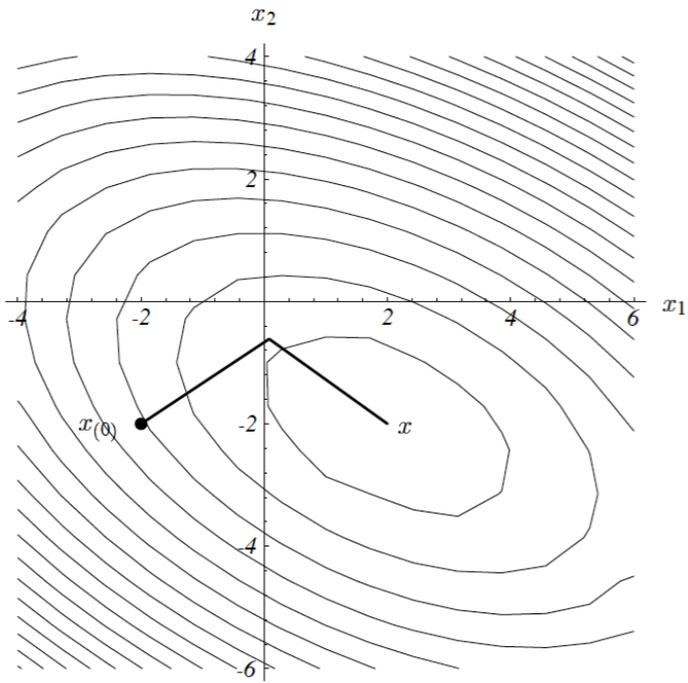
$$r_{(k+1)} = r_{(k)} - \alpha_{(k)} A P_{(k)}$$

if $r_{(k+1)} < \epsilon$, return $x_{(k+1)}$

$$\beta_{(k)} = \frac{r_{(k+1)}^T r_{(k+1)}}{r_{(k)}^T r_{(k)}}$$

$$P_{(k+1)} = r_{(k+1)} + \beta_{(k)} P_{(k)}$$

$$k = k + 1$$



$r_{(k)}$ and $r_{(k+1)}$ are orthogonal

$P_{(k)}$ and $P_{(k+1)}$ are A conjugate

N step converge

$$A = \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix} \quad b = \begin{bmatrix} 2 \\ -8 \end{bmatrix}$$

$$\eta_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad \eta_2 = 2 \begin{bmatrix} -2 \\ 1 \end{bmatrix}$$

$$r_{(0)}^T \cdot r_{(0)} = 0$$

$$x_{(0)} = \begin{bmatrix} -2 \\ -2 \end{bmatrix} \quad r_{(0)} = b - Ax_{(0)} = \begin{bmatrix} 2 \\ -8 \end{bmatrix} - \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix} \begin{bmatrix} -2 \\ -2 \end{bmatrix} = \begin{bmatrix} 12 \\ 8 \end{bmatrix}$$

$$\frac{104}{600} \quad \frac{52}{300} \quad \frac{26}{150}$$

$$P_{(0)} = r_{(0)}$$

$$\begin{matrix} & 2 & 10 \\ -10 & & -8 + 16 \\ -4 - 12 & & \end{matrix}$$

$$\frac{13}{75}$$

$$\alpha_{(0)} = \frac{r_{(0)}^T r_{(0)}}{P_{(0)}^T A P_{(0)}} = \frac{208}{1200} = \frac{13}{75}$$

$$[12 \quad 24] \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix} \begin{bmatrix} 12 \\ 24 \end{bmatrix}$$

$$\frac{156}{75} \quad \frac{15}{75}$$

$$x_{(1)} = x_{(0)} + \alpha_{(0)} P_{(0)} = \begin{bmatrix} -2 \\ -2 \end{bmatrix} + \frac{13}{75} \begin{bmatrix} 12 \\ 8 \end{bmatrix} = \begin{bmatrix} \frac{2}{25} \\ -\frac{46}{75} \end{bmatrix}$$

$$\begin{matrix} [84 & 168] \\ \begin{bmatrix} 12 \\ 24 \end{bmatrix} \end{matrix}$$

$$r_{(1)} = b - Ax_{(1)} = \begin{bmatrix} 2 \\ -8 \end{bmatrix} - \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix} \begin{bmatrix} \frac{2}{25} \\ -\frac{46}{75} \end{bmatrix} = \begin{bmatrix} 2 \\ -8 \end{bmatrix} - \begin{bmatrix} \frac{74}{75} \\ -\frac{264}{75} \end{bmatrix}$$

$$\frac{6}{25} - \frac{92}{75}$$

$$\beta_{(0)} = \frac{r_{(1)}^T r_{(1)}}{r_{(0)}^T r_{(0)}} = 0.1394$$

$$= \begin{bmatrix} \frac{224}{75} \\ \frac{336}{75} \end{bmatrix}$$

$$\frac{18-92}{75}$$

$$P_{(1)} = r_{(1)} + \beta_{(0)} P_{(0)} = \begin{bmatrix} \frac{224}{75} \\ -\frac{336}{75} \end{bmatrix} + 0.1394 \begin{bmatrix} 12 \\ 8 \end{bmatrix} = \begin{bmatrix} 4.6547 \\ -3.3648 \end{bmatrix}$$

$$\frac{4}{75} - \frac{276}{75} = \frac{12-276}{75}$$

$$\alpha_{(1)} = \frac{r_{(1)}^T r_{(1)}}{P_{(1)}^T A P_{(1)}} = \frac{163072}{5625} = 0.412$$

$$\begin{matrix} (4.6547 & -3.3648) \\ \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix} \\ \begin{bmatrix} 4.6547 \\ -3.3648 \end{bmatrix} \end{matrix}$$

$$70.3505$$

$$x_{(2)} = x_{(1)} + \alpha_{(1)} P_{(1)} = \begin{bmatrix} \frac{2}{25} \\ -\frac{46}{75} \end{bmatrix} + 0.412 \begin{bmatrix} 4.6547 \\ -3.3648 \end{bmatrix} = \begin{bmatrix} 2 \\ -2 \end{bmatrix}$$

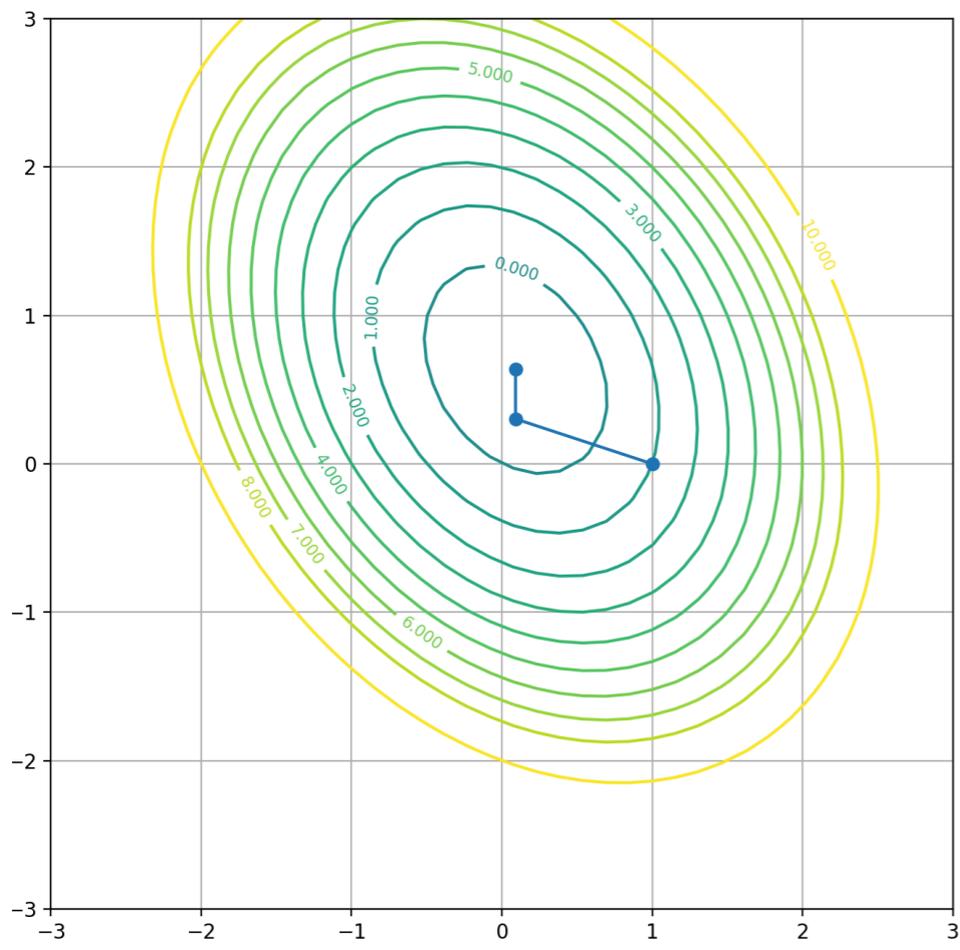
$$(7.24001 - 10.86006)$$

check

$$r_{(1)}^T \cdot r_{(0)} = \begin{bmatrix} \frac{224}{75} & -\frac{336}{75} \end{bmatrix} \begin{bmatrix} 12 \\ 8 \end{bmatrix} = \frac{224 \cdot 12 - 336 \cdot 8}{75} = \frac{2688 - 2688}{75} = 0$$

$$P_{(1)}^T A P_{(0)} = \begin{bmatrix} 4.6547 & -3.3648 \end{bmatrix} \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix} \begin{bmatrix} 12 \\ 8 \end{bmatrix} = \begin{bmatrix} 7.24001 & -10.86006 \end{bmatrix} \begin{bmatrix} 12 \\ 8 \end{bmatrix} = 0.02684$$

round off errors



$$\underline{A} \underline{x} = \underline{b}, \quad \begin{bmatrix} 4 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \quad \text{solution } \underline{x} = \begin{bmatrix} 0.0909 \\ 0.6364 \end{bmatrix} = \begin{bmatrix} \frac{1}{11} \\ \frac{7}{11} \end{bmatrix}$$

$$\underline{x}_0 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\underline{p}_0 := \underline{r}_0 = \underline{b} - \underline{A} \cdot \underline{x}_0 = \begin{bmatrix} 1 \\ 2 \end{bmatrix} - \begin{bmatrix} 4 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} -3 \\ 1 \end{bmatrix}$$

$$\alpha_0 = \frac{\underline{r}_0^T \cdot \underline{r}_0}{\underline{p}_0^T \cdot \underline{A} \cdot \underline{p}_0} = \frac{[-3 \ 1] \cdot \begin{bmatrix} -3 \\ 1 \end{bmatrix}}{[-3 \ 1] \begin{bmatrix} 4 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} -3 \\ 1 \end{bmatrix}} = \frac{10}{33}$$

$$\underline{x}_1 = \underline{x}_0 + \alpha_0 \underline{p}_0 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \frac{10}{33} \begin{bmatrix} -3 \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{11} \\ \frac{10}{33} \end{bmatrix}$$

$$\underline{r}_1 = \underline{r}_0 - \alpha_0 \underline{A} \cdot \underline{p}_0 = \begin{bmatrix} -3 \\ 1 \end{bmatrix} - \frac{10}{33} \begin{bmatrix} 4 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} -3 \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{3} \\ 1 \end{bmatrix}$$

$$\beta_0 = \frac{\underline{r}_1^T \cdot \underline{r}_1}{\underline{r}_0^T \cdot \underline{r}_0} = \frac{[\frac{1}{3} \ 1] \begin{bmatrix} \frac{1}{3} \\ 1 \end{bmatrix}}{[-3 \ 1] \begin{bmatrix} -3 \\ 1 \end{bmatrix}} = \frac{1}{9}$$

$$\underline{p}_1 = \underline{r}_1 + \beta_0 \underline{p}_0 = \begin{bmatrix} \frac{1}{3} \\ 1 \end{bmatrix} + \frac{1}{9} \begin{bmatrix} -3 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ \frac{10}{9} \end{bmatrix} \quad \frac{\frac{1}{9}}{\frac{10}{9}} = \frac{1}{10}$$

$$\alpha_1 = \frac{\underline{r}_1^T \cdot \underline{r}_1}{\underline{p}_1^T \cdot \underline{A} \cdot \underline{p}_1} = \frac{[\frac{1}{3} \ 1] \begin{bmatrix} \frac{1}{3} \\ 1 \end{bmatrix}}{[0 \ \frac{10}{9}] \begin{bmatrix} 4 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} 0 \\ \frac{10}{9} \end{bmatrix}} = \frac{30}{10} = \frac{3}{1}$$

$$\underline{x}_2 = \underline{x}_1 + \alpha_1 \underline{p}_1 = \begin{bmatrix} \frac{1}{11} \\ \frac{10}{33} \end{bmatrix} + \frac{3}{1} \begin{bmatrix} 0 \\ \frac{10}{9} \end{bmatrix} = \begin{bmatrix} \frac{1}{11} \\ \frac{7}{11} \end{bmatrix} \checkmark$$

$$\underline{r}_2 = \underline{r}_1 - \alpha_1 \underline{A} \cdot \underline{p}_1 = \begin{bmatrix} \frac{1}{3} \\ 1 \end{bmatrix} - \frac{3}{1} \begin{bmatrix} 4 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} 0 \\ \frac{10}{9} \end{bmatrix} = \begin{bmatrix} \frac{1}{3} \\ 1 \end{bmatrix} - \frac{3}{1} \begin{bmatrix} \frac{10}{9} \\ \frac{10}{3} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \checkmark$$

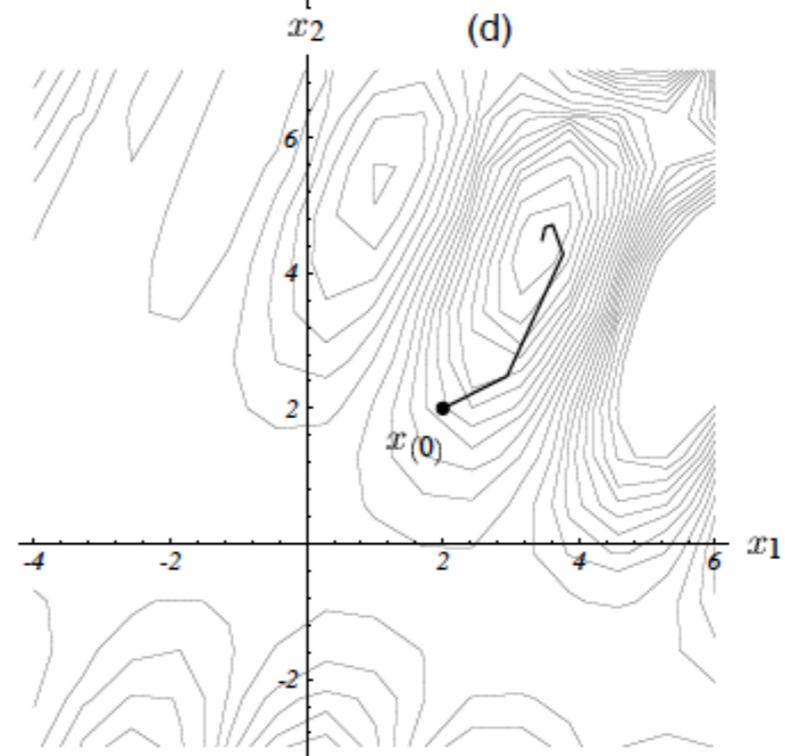
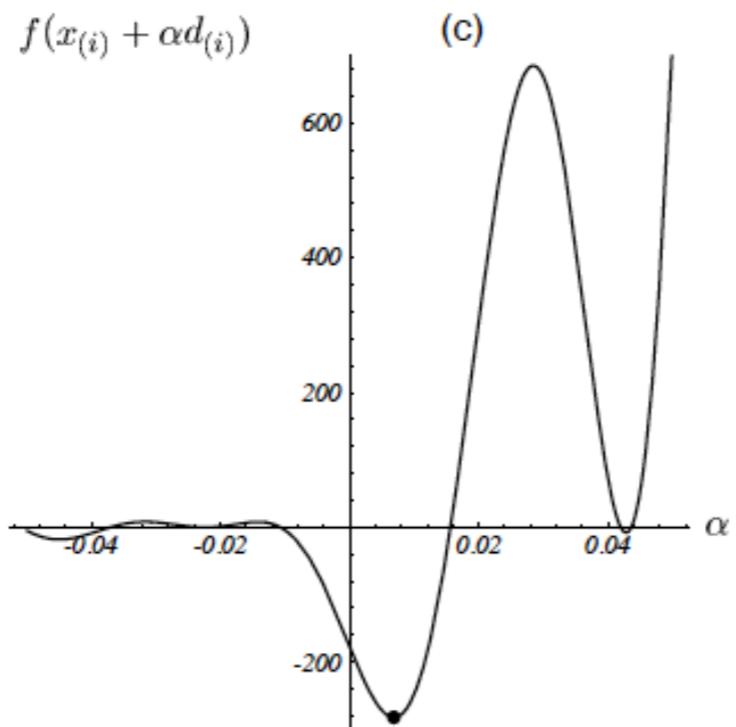
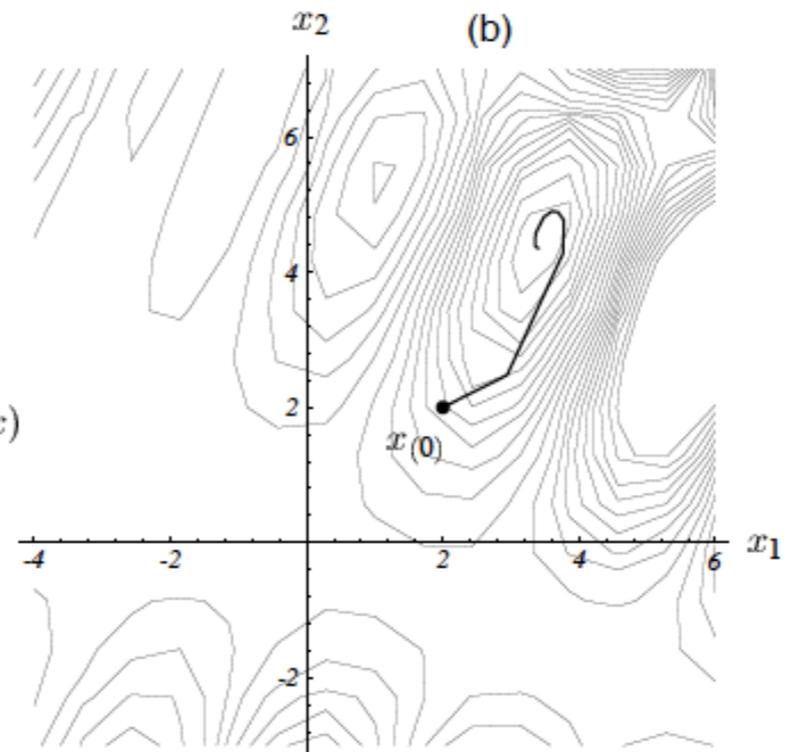
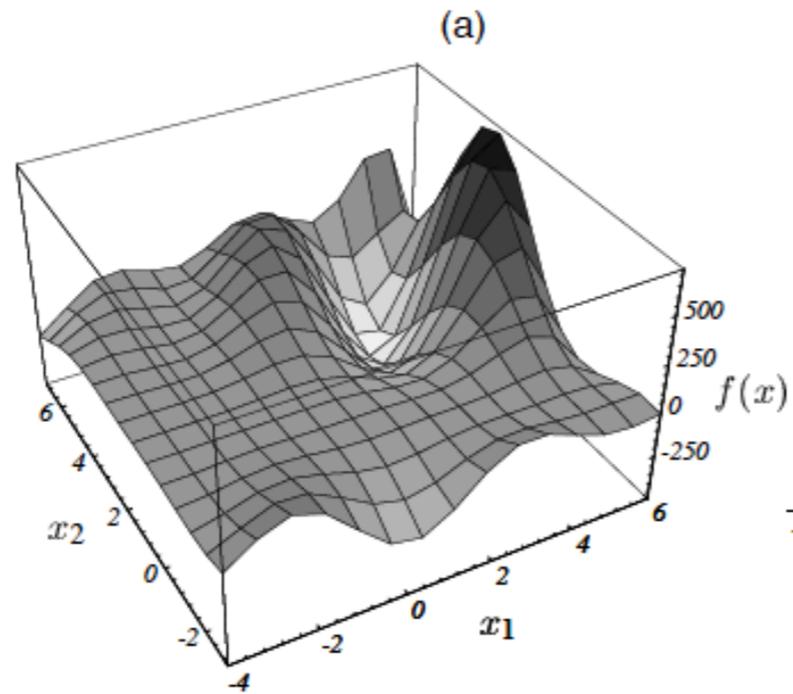
check: $\underline{r}_k, \underline{r}_{k+1}$ orthogonal \checkmark

$\underline{p}_k \cdot \underline{A} \underline{p}_{k+1}$ orthogonal \checkmark

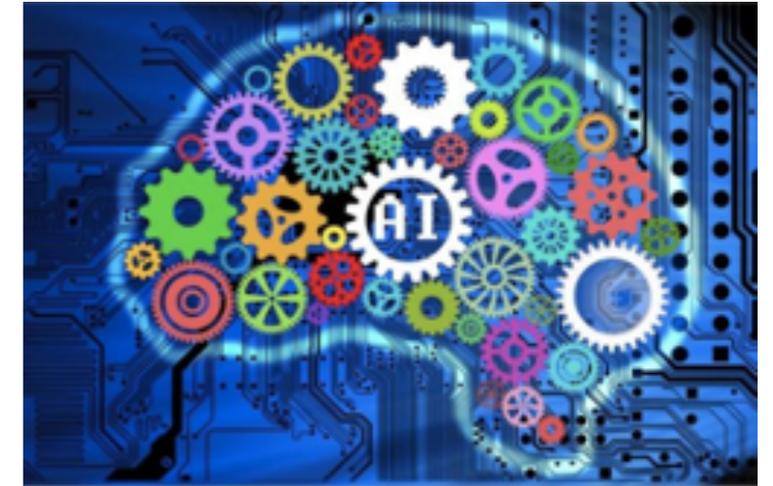
n-stop converge \checkmark

$$\underline{r}_0^T \cdot \underline{r}_1 = [-3 \ 1] \begin{bmatrix} \frac{1}{3} \\ 1 \end{bmatrix} = 0$$

$$\underline{p}_0^T \cdot \underline{A} \underline{p}_1 = [-3 \ 1] \begin{bmatrix} 4 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} 0 \\ \frac{10}{9} \end{bmatrix} = [-11 \ 0] \begin{bmatrix} 0 \\ \frac{10}{9} \end{bmatrix} = 0$$



Content



0. Introduction

1. Regression

1.1 Multivariate Linear Regression (curve fitting)

1.2 Regularization (Lagrange multiplier)

1.3 Logistic Regression (Fermi-Dirac distribution)

1.4 Support Vector Machine (high-school geometry)

2. Dimensionality Reduction/feature extraction

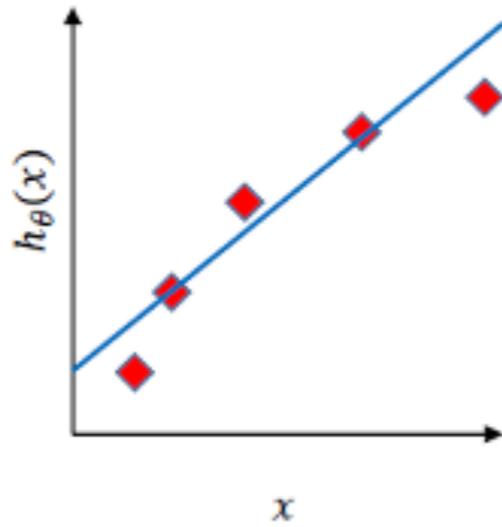
2.1 Principal Component Analysis (order parameters)

2.2 Recommender Systems

2.3 Clustering (phase transition)

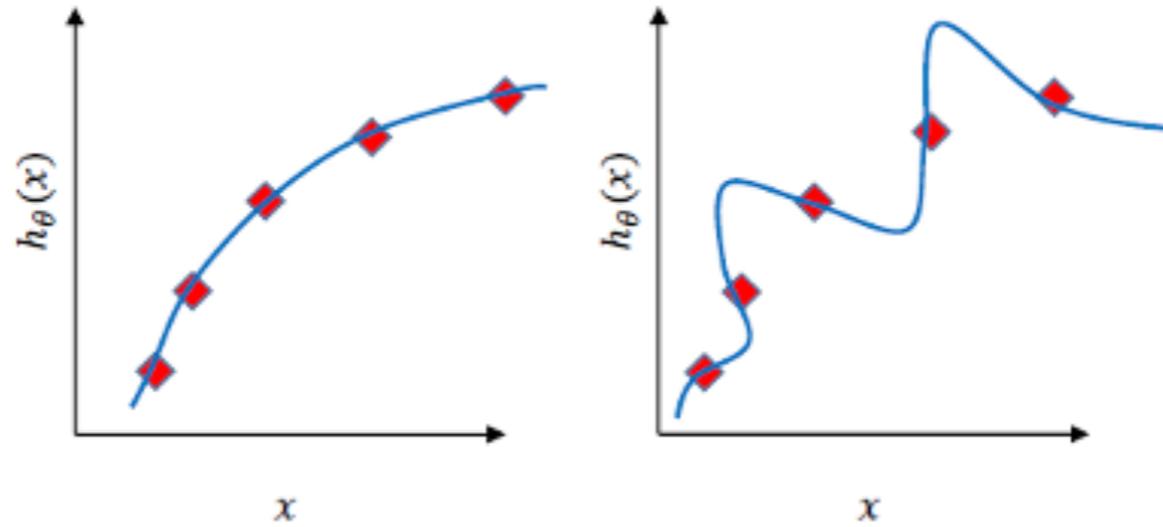
The quality of fitting

Underfitting (high bias)



High bias

Overfitting (high variance)

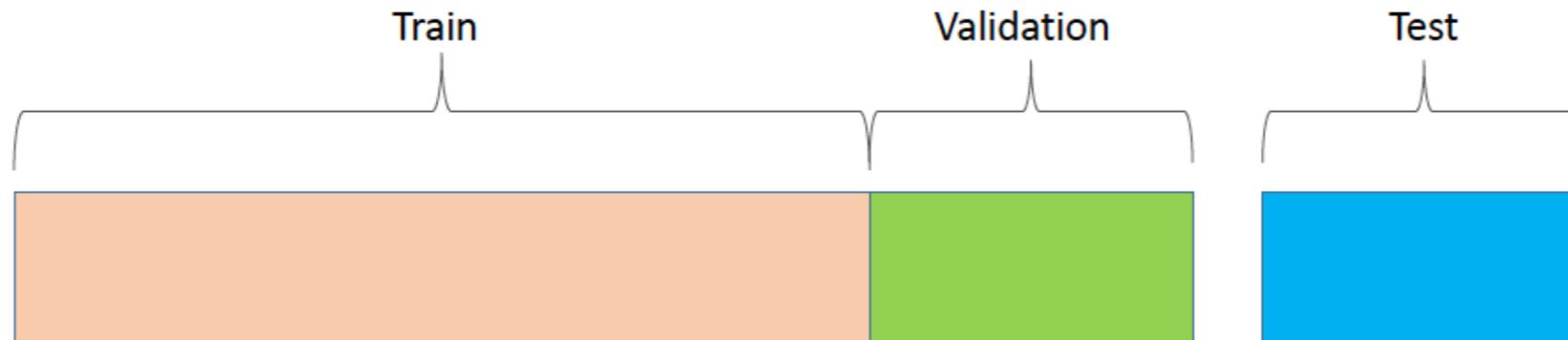


High variance

Polynomial regression

Forecastability

$$h_{\theta}^{[d]}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_d x^d = \theta^T x$$

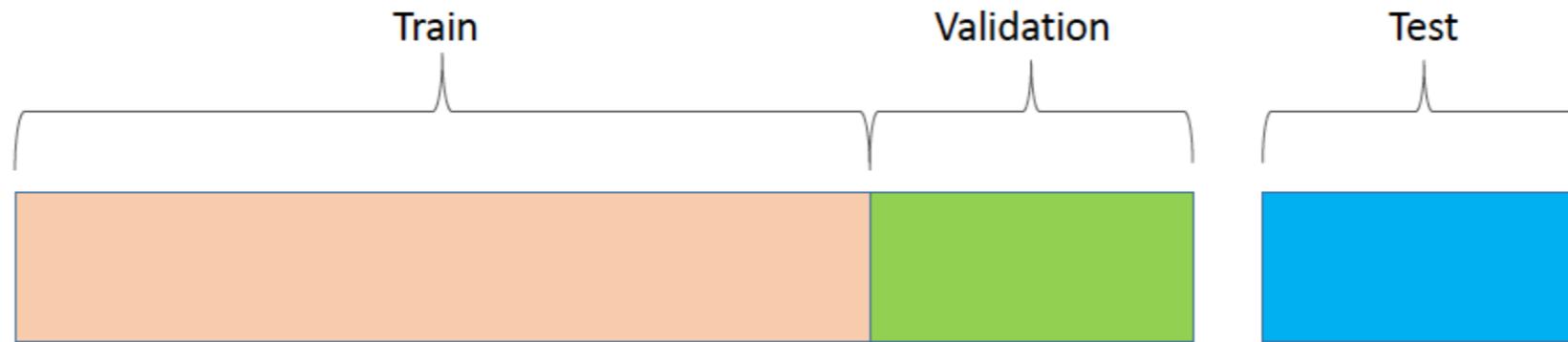


$$J_{\text{train}}(\theta) = \frac{1}{2M_{\text{train}}} \sum_{i_{\text{train}}=1}^{M_{\text{train}}} (y^{(i_{\text{train}})} - h_{\theta}(x))^2$$

$$J_{\text{cv}}(\theta) = \frac{1}{2M_{\text{cv}}} \sum_{i_{\text{cv}}=1}^{M_{\text{cv}}} (y^{(i_{\text{cv}})} - h_{\theta}(x))^2$$

$$J_{\text{test}}(\theta) = \frac{1}{2M_{\text{test}}} \sum_{i_{\text{test}}=1}^{M_{\text{test}}} (y^{(i_{\text{test}})} - h_{\theta}(x))^2$$

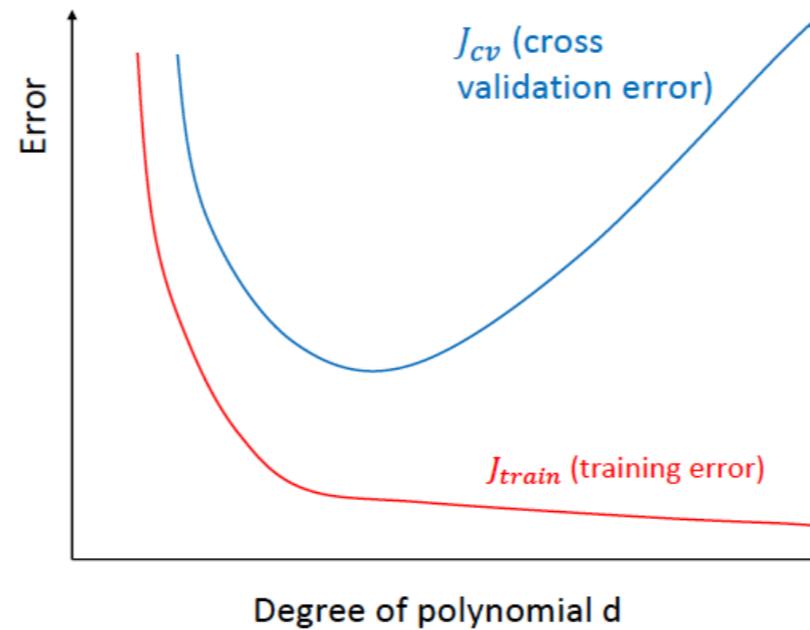
Learning Curves and Regularization



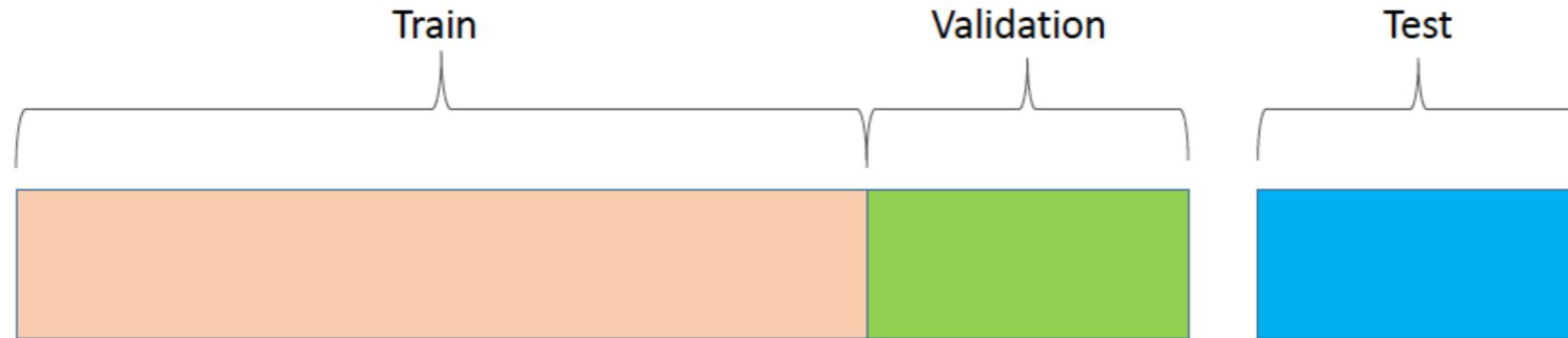
Forecastability

$$h_{\theta}^{[d]}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_d x^d = \theta^T x \quad \text{Polynomial regression}$$

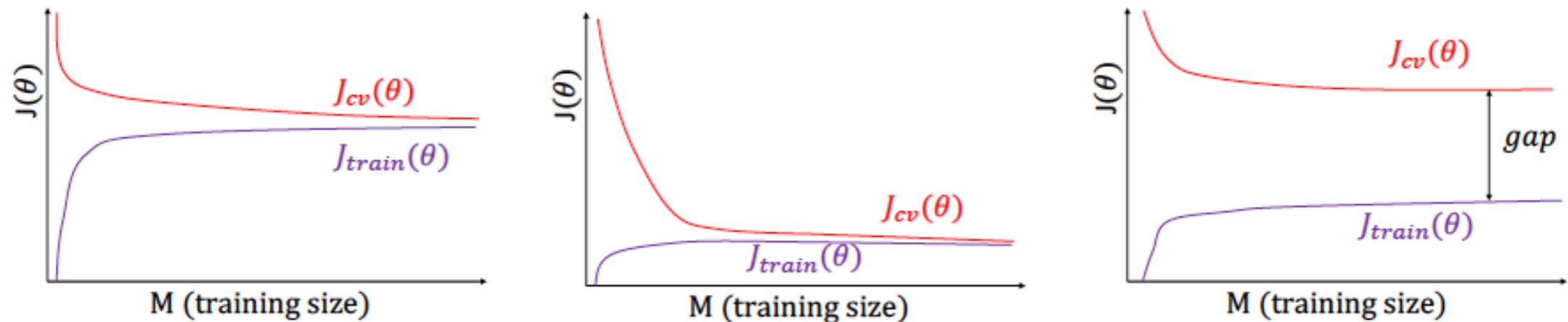
Model Selection



Learning Curves and Regularization



Learning curves



Regularization

$$J(\theta) = \frac{1}{2M} \left[\sum_{i=1}^M (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^N \theta_j^2 \right] \quad \lambda \geq 0$$