

# Mathematics for AI

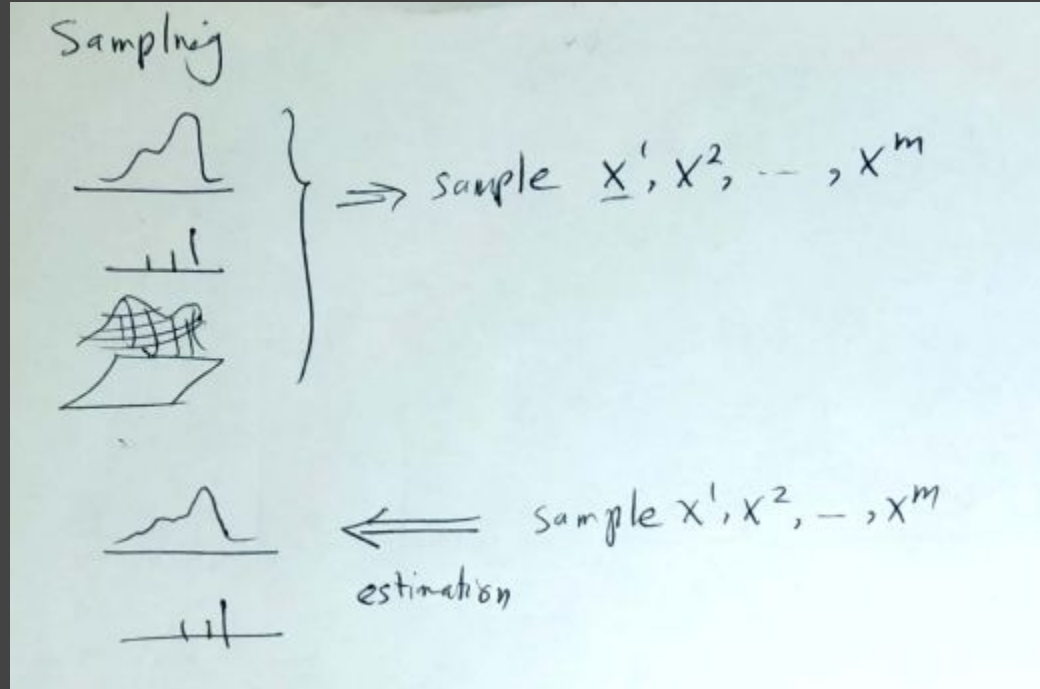
## Lecture 27

Statistical Estimation, Maximum Likelihood Solution  
Introduction to Optimization

# Statistical Estimation



K. N. Toosi  
University of Technology

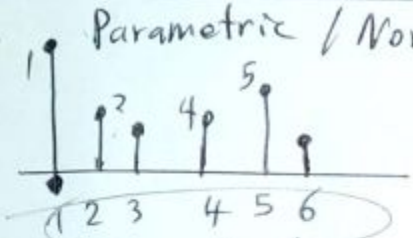



# Example: Rolling a dice



Estimation

Parametric / Nonparametric



$\Pr(X=1)=P_1$     $\Pr(X=2)=P_2$    ...    $P_6$

$P_1 + P_2 + \dots + P_6 = 1$

$X^1$	$X^2$	$X^3$	...	$X^m$
1	2	3		6

$\theta = \begin{bmatrix} P_1 \\ P_2 \\ \vdots \\ P_6 \end{bmatrix}$

$P_\theta(x) = \Pr(X=x) = P_x \quad x \in \{1, 2, \dots, 6\}$

$x^1, x^2, \dots, x^m$  i.i.d independent  
 $x^i \sim P$

# The likelihood function



Fix  $\theta = [p_1, p_2, \dots, p_6]$  find Prob.  $x^1, x^2, \dots, x^m$

$$\begin{aligned} \mathcal{L}(\theta) &= \text{Prob}(x^1, x^2, \dots, x^m | \theta) = \text{Prob}(x^1 | \theta) \text{Prob}(x^2 | \theta) \dots \text{Prob}(x^m | \theta) \\ &\stackrel{\text{likelihood}}{=} \prod_{i=1}^m \text{Prob}(x^i | \theta) = \prod_{i=1}^m p_{\theta}(x^i) = \prod_{i=1}^m p_{x^i} \\ &= p_1^{m_1} p_2^{m_2} \dots p_6^{m_6} = \mathcal{L}(\theta) = \mathcal{L}(p_1, p_2, \dots, p_6) \end{aligned}$$

$$m_1 = \#(x^i = 1)$$

$$m_2 = \#(x^i = 2)$$

$$m_6 = \#(x^i = 6)$$

# Maximum Likelihood (ML) Solution



K. N. Toosi  
University of Technology

$$L(\theta) = L(p_1, p_2, \dots, p_6) = p_1^{m_1} p_2^{m_2} \dots p_6^{m_6} = L(\theta) = L(\theta; D) \quad \text{MA 27 II}$$

$$\theta^* = \underset{\theta}{\operatorname{argmax}} L(\theta, D)$$

→ likelihood  
→ Maximum Likelihood solution

$$L(\theta; x_1, \dots, x^m) \\ = \Pr(D | \theta) \\ = \Pr(x_1, \dots, x^m | \theta)$$

# Example: Rolling a dice



$$\theta^* = \operatorname{argmax}_{\theta} L(\theta; D) \text{ subject to } \boxed{P_1, P_2, \dots, P_6 \geq 0}$$
$$P_1 + P_2 + \dots + P_6 = 1$$

$$\theta^* = \operatorname{argmax}_{\theta} P_1^{m_1} P_2^{m_2} \dots P_5^{m_5} (1 - P_1 - P_2 - \dots - P_5)^{m_6}$$

$$= \operatorname{argmax}_{\theta} \log (P_1^{m_1} P_2^{m_2} \dots P_5^{m_5} (1 - \sum_{i=1}^5 P_i)^{m_6})$$

$$= \operatorname{argmax}_{\theta} \log m_1 \log P_1 + m_2 \log P_2 + \dots + m_5 \log P_5$$

$$\theta = [P_1, P_2, \dots, P_5]$$

$$+ m_6 \log (1 - \sum_{i=1}^5 P_i)$$

log-likelihood

# Example: Rolling a dice



$$= \arg \max_{\theta} \log m_1 \log p_1 + m_2 \log p_2 + \dots + m_5 \log p_5 + m_6 \log \left( 1 - \sum_{i=1}^5 p_i \right)$$

$\theta = [p_1, p_2, \dots, p_5]$

log-likelihood  $LL(\theta)$

$$\frac{\partial}{\partial p_1} LL(\theta) = \frac{m_1}{p_1} + m_6 \frac{-1}{1 - \sum_{i=1}^5 p_i} = 0$$

$$\frac{\partial}{\partial p_i} LL(\theta) = \frac{m_i}{p_i} + m_6 \frac{-1}{1 - \sum p_i} = 0 \quad i=1, \dots, 5$$

$$\Rightarrow \frac{m_1}{p_1} = \frac{m_2}{p_2} = \frac{m_3}{p_3} = \dots = \frac{m_5}{p_5} = \frac{m_6}{1 - \sum_{i=1}^5 p_i} = \frac{m_6}{p_6} = \frac{1}{\alpha}$$

# Example: Rolling a dice



$$\frac{\partial}{\partial p_1} \mathcal{L}(\theta) = \frac{m_1}{p_1} + m_6 \frac{-1}{1 - \sum_{i=1}^5 p_i} = 0 \quad \mathcal{L}(\theta)$$

$$\frac{\partial}{\partial p_i} \mathcal{L}(\theta) = \frac{m_i}{p_i} + m_6 \frac{-1}{1 - \sum p_i} = 0 \quad i=1, \dots, 5$$

$$\Rightarrow \frac{m_1}{p_1} = \frac{m_2}{p_2} = \frac{m_3}{p_3} = \dots = \frac{m_5}{p_5} = \frac{m_6}{1 - \sum_{i=1}^5 p_i} = \frac{m_6}{p_6} = \frac{1}{\alpha}$$

$$\bullet \quad p_1 = \alpha m_1 \quad p_2 = \alpha m_2 \quad \dots \quad p_6 = \alpha m_6$$

$$\sum p_i = 1 \Rightarrow \alpha (m_1 + m_2 + \dots + m_6) = 1 = \alpha m = 1 \Rightarrow \alpha = \frac{1}{m}$$

$$p_1 = \alpha m_1 = \frac{m_1}{m} \quad p_2 = \frac{m_2}{m} \quad \dots \quad p_6 = \frac{m_6}{m}$$

maximum-likelihood solution



# Example: Normal Distribution

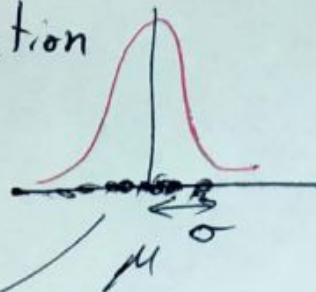


K. N. Toosi  
Technology

Example: Normal Distribution

MA27 (II)

$$p(x) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$$



$$\theta = [\mu, \sigma]$$

samples  $x^1, x^2, \dots, x^m$

$$\begin{aligned} \text{Like lihood } L(\theta) = L(\mu, \sigma) &= P(x^1 - x^m) = \prod_{i=1}^m p_{\theta}(x^i) \\ &= \prod_{i=1}^m \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2} \frac{(x^i - \mu)^2}{\sigma^2}} \end{aligned}$$

(see likelihood 10/1)

# Example: Normal Distribution Likelihood and log-likelihood



K. N. Toosi  
University of Technology

$$\begin{aligned} \text{Likelihood } L(\theta) &= L(\mu, \sigma) = P(n^1, \dots, n^m) = \prod_{i=1}^m P_{\theta}(n^i) \\ &= \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \frac{(n^i - \mu)^2}{\sigma^2}} \end{aligned}$$

$$\text{Log-likelihood } \ell(\theta) = \log_e L(\mu, \sigma) =$$

$$= \sum_{i=1}^m \left( -\log \sqrt{2\pi} - \log \sigma - \frac{1}{2} \frac{(n^i - \mu)^2}{\sigma^2} \right)$$

$$\Rightarrow \ell(\theta) = -m \log \sqrt{2\pi} - m \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^m (n^i - \mu)^2$$

# Example: Normal Distribution - ML Solution



$$\Rightarrow \ell(\theta) = -m \log \sqrt{2\pi} - m \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^m (x^i - \mu)^2$$

$$\frac{\partial \ell(\theta)}{\partial \mu} = -\frac{1}{2\sigma^2} \sum_{i=1}^m (-2)(x^i - \mu) = 0 \Rightarrow \sum_{i=1}^m (x^i - \mu) = 0$$

$$\Rightarrow \sum_{i=1}^m x^i - m\mu \Rightarrow \mu^* = \frac{1}{m} \sum_{i=1}^m x^i$$

$$\frac{\partial \ell(\theta)}{\partial \sigma} = \frac{-m}{\sigma} - \frac{1}{2\sigma^3} \sum (x^i - \mu)^2 = 0$$

$$\Rightarrow -m + \frac{1}{\sigma^2} \sum_{i=1}^m (x^i - \mu)^2 = 0$$

$$\sigma^{*2} = \frac{1}{m} \sum_{i=1}^m (x^i - \mu^*)^2$$

# Exercise



K. N. Toosi  
University of Technology

Exercise: find ML solution for  
uniform dist.  $U[a, b]$   $= p(x) = \begin{cases} \frac{1}{b-a} & x \in [a, b] \\ 0 & \text{otherwise} \end{cases}$   
exponential dist.  $p(x) = \lambda e^{-\lambda x}$   
 $x \geq 0$

# Exercise



Exercise:  $N(\mu, \Sigma)$  (IV)

multivariable Gaussian

$$P(x) = \frac{1}{(2\pi)^{d/2} \det(\Sigma)} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

$x \in \mathbb{R}^d$

$ll(\vec{\mu}, \Sigma)$

$x^1, x^2, \dots, x^m$

~~Maximum~~

ML Solution:

$$\mu^* = \frac{1}{m} \sum_{i=1}^m x^i$$

$$\Sigma^* = \frac{1}{m} \sum_{i=1}^m (x^i - \mu^*)(x^i - \mu^*)^T$$

# Maximum Likelihood General Approach



Find  $L(\theta)$  or  $\ln L(\theta)$

$$\nabla_L(\theta) = 0 \Rightarrow \theta = \checkmark$$

$$\theta = (\theta_1, \theta_2, \dots, \theta_n)$$

$$\nabla_L(\theta) = \begin{bmatrix} \frac{\partial L}{\partial \theta_1} \\ \frac{\partial L}{\partial \theta_2} \\ \vdots \\ \frac{\partial L}{\partial \theta_n} \end{bmatrix} = 0 = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \left. \begin{array}{l} \text{\{n-equations\}} \\ \text{\{n-unknowns\}} \end{array} \right\}$$

# Maximum Likelihood General Approach



$$\theta = (\theta_1, \theta_2, \dots, \theta_n)$$
$$\nabla_L(\theta) = \begin{bmatrix} \frac{\partial L}{\partial \theta_1} \\ \frac{\partial L}{\partial \theta_2} \\ \vdots \\ \frac{\partial L}{\partial \theta_n} \end{bmatrix} = \vec{0} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \left. \begin{array}{l} \text{\{ n-equation,} \\ \text{\} n-unknown} \end{array} \right\}$$

→ A (usually non-linear) system of equations.

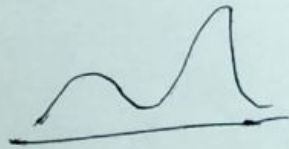
→ In many cases it is very hard to solve (impossible to solve!)

# Example: Mixture of Gaussians



solve comp

$$\text{Example: } p(x) = \frac{1}{3} \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{1}{2} \frac{(x-\mu_1)^2}{\sigma_1^2}} + \frac{2}{3} \frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{1}{2} \frac{(x-\mu_2)^2}{\sigma_2^2}}$$



mixture of Gaussians.

$$\theta = [\mu_1, \sigma_1, \mu_2, \sigma_2]$$

$$\Rightarrow x^* = \operatorname{argmax}_{\theta} L(\theta)$$

↳ Hard to find global optimum



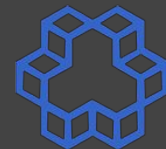
# Remember: Learning from data



K. N. Toosi  
University of Technology



- Training data  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$ 
  - choose  $\theta$  such that  $f(\theta, \mathbf{x}_i)$  is close to  $y_i$



# Remember: Learning from data: Cost function



- Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ 
  - choose  $\theta$  such that  $f(\theta, x_i)$  is close to  $y_i$

Cost function

# Remember: Learning from data: Cost function



K. N. Toosi  
University of Technology



- Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ 
  - choose  $\theta$  such that  $f(\theta, x_i)$  is close to  $y_i$
  - cost function:

$$C(\theta) = \sum_{i=1..N} d(f(\theta, x_i), y_i)$$

# Remember: Learning from data: Cost function



K. N. Toosi  
University of Technology



- Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ 
  - choose  $\theta$  such that  $f(\theta, x_i)$  is close to  $y_i$
  - cost function:

$$C(\theta) = \sum_{i=1..n} d(f(\theta, x_i), y_i)$$

↓  
data output

# Remember: Learning from data: Cost function



K. N. Toosi  
University of Technology



- Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ 
  - choose  $\theta$  such that  $f(\theta, x_i)$  is close to  $y_i$
  - cost function:

$$C(\theta) = \sum_{i=1..n} d(f(\theta, x_i), y_i)$$



model output given  $x_i$

# Remember: Learning from data: Cost function



K. N. Toosi  
University of Technology



- Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ 
  - choose  $\theta$  such that  $f(\theta, x_i)$  is close to  $y_i$
  - cost function:

$$C(\theta) = \sum_{i=1..n} d(f(\theta, x_i), y_i)$$

↓  
distance

# Remember: Learning from data: Cost function



K. N. Toosi  
University of Technology



- Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ 
  - choose  $\theta$  such that  $f(\theta, x_i)$  is close to  $y_i$
  - cost function:

$$C(\theta) = \sum_{i=1..n} \| f(\theta, x_i) - y_i \|^2$$

distance

# Remember: Learning from data: Cost function



K. N. Toosi  
University of Technology



- Training data  $(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_N, \mathbf{y}_N)$ 
  - choose  $\theta$  such that  $f(\theta, \mathbf{x}_i)$  is close to  $\mathbf{y}_i$
  - cost function:

$$C(\theta) = \sum_{i=1..n} d(f(\theta, \mathbf{x}_i), \mathbf{y}_i)$$



# Remember: Learning from data: Cost function



K. N. Toosi  
University of Technology



- Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ 
  - choose  $\theta$  such that  $f(\theta, x_i)$  is close to  $y_i$
  - cost function:

$$C(\theta) = \sum_{i=1..n} d(f(\theta, x_i), y_i)$$

choose  $\theta$  such that  $C(\theta)$  is small

# Remember: Learning from data: Cost function



K. N. Toosi  
University of Technology



- Training data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ 
  - choose  $\theta$  such that  $f(\theta, x_i)$  is close to  $y_i$
  - cost function:

$$C(\theta) = \sum_{i=1..N} d(f(\theta, x_i), y_i)$$

$$\theta^* = \operatorname{argmin}_{\theta} C(\theta)$$

# Optimization: Continuous vs Discrete



Optimization

$f(x)$

$f: \mathbb{R}^d \rightarrow \mathbb{R}$   
 $f: S \rightarrow \mathbb{R}$   
 $S \subseteq \mathbb{R}^d$

$x^* = \operatorname{argmin}_x f(x)$   
 $x^* = \operatorname{argmax}_x f(x)$

~~Discrete Optimization  
(Combinatorial)~~

~~$f(\theta, x) = Ax + b$~~

# Remember: Linear Regression



K. N. Toosi  
University of Technology

$$C(\theta) = C(A, b) = \sum_{i=1}^{m'} \|Ax^i - b - y^i\|^2$$

(VI)

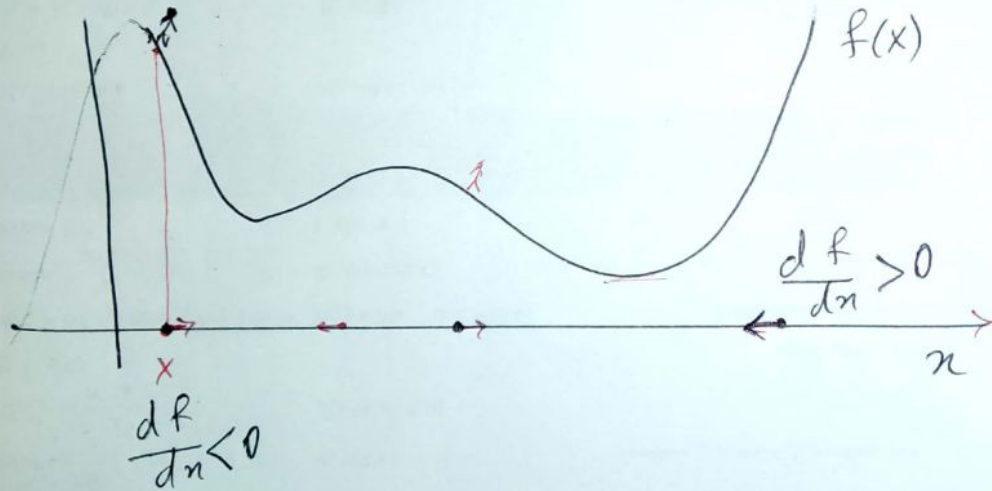
$$\left. \begin{array}{l} \frac{\partial C}{\partial A} = 0_{m \times n} \\ \frac{\partial C}{\partial b} = \vec{0} \end{array} \right\} \Rightarrow \text{closed form solution}$$

# Iterative Optimization Algorithms

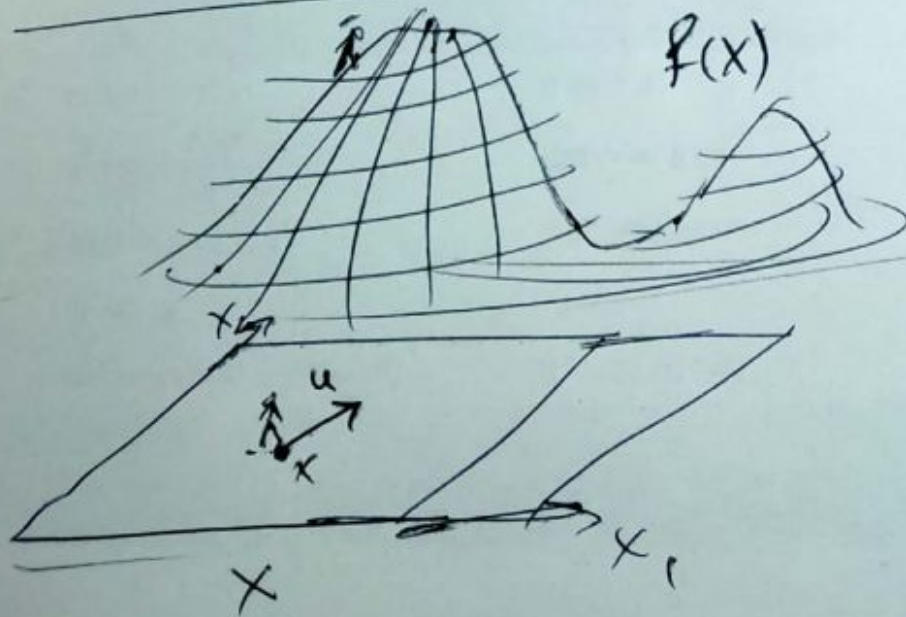


K. N. Toosi  
University of Technology

In most cases we cannot  
solve  $\frac{\partial C}{\partial \theta} = 0$  for  $\theta$ .

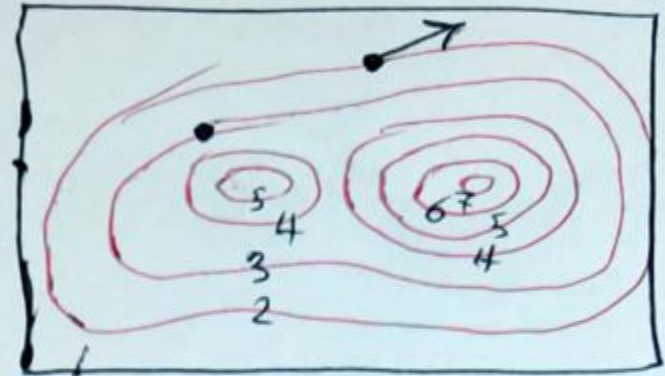
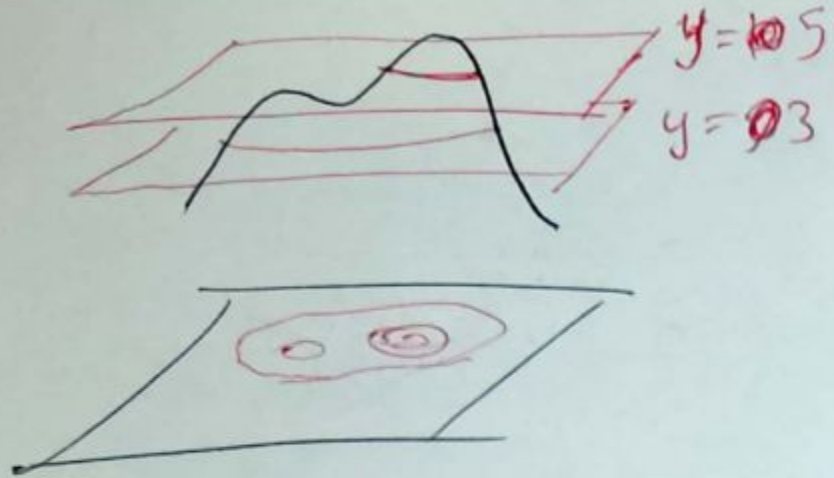


# Steepest Descent



$$\|u\|=1 \quad \left\{ \begin{array}{l} D[u]f(x) = \nabla_f^T u \\ u = \frac{\nabla_f}{\|\nabla_f\|} = \operatorname{argmax} \nabla_f^T u \\ \text{steepest ascend} \quad \text{s.t. } \|u\|=1 \\ u = -\frac{\nabla_f}{\|\nabla_f\|} = \operatorname{argmin} \nabla_f^T u \\ \text{steepest descent} \quad \|u\|=1 \end{array} \right.$$

# Level Curves

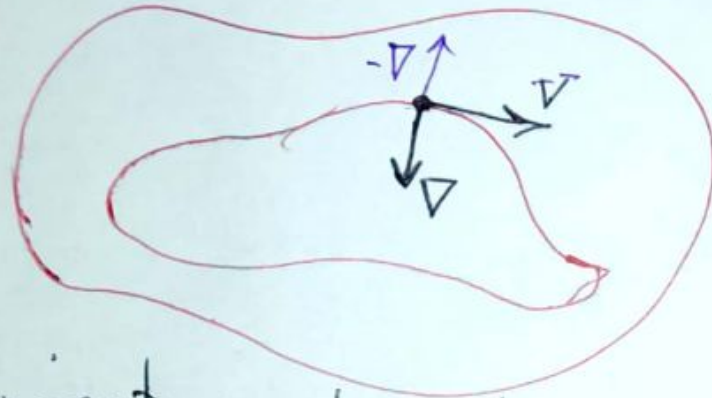


level curves  
~~contour~~ contour lines

# Level Curves and the Gradient Vector



K. N. Toosi  
University of Technology



$$v \perp \nabla$$



$$D[\mathbf{v}]f = 0$$

$$\Rightarrow v^T \nabla = 0$$

any tangent vector to contour curves  
is perpendicular to the gradient  
of  $f$  at that point.