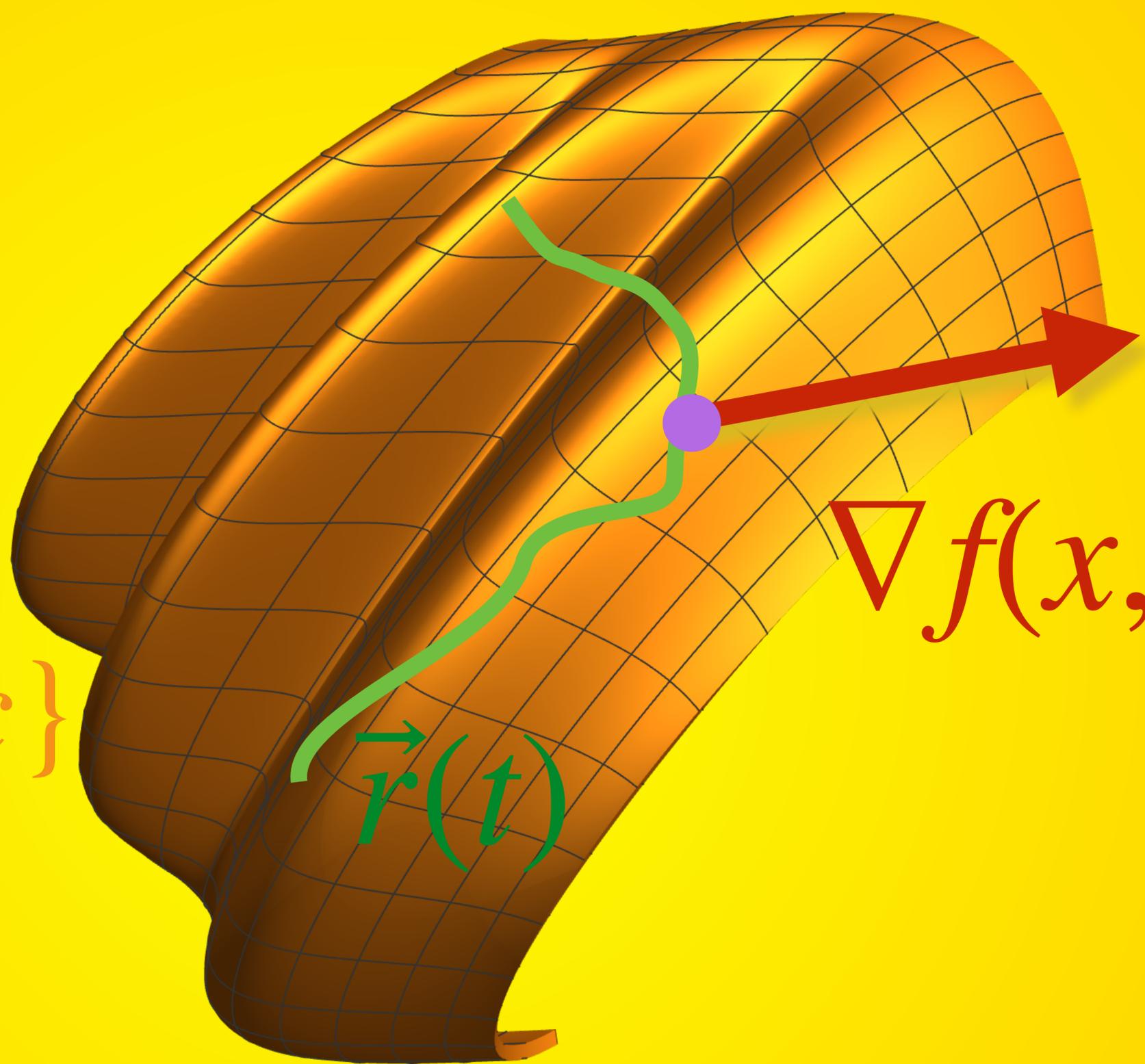# Lecture 24

# Gradient Rule

# Quadratic Approx

# Gradient Theorem

# Gradient Theorem

$$\nabla f \quad \text{is perpendicular to the level surface } \{f = c\}$$

*Proof*

$\{f(x, y, z) = c\}$

$\nabla f(x, y, z)$

$\vec{r}(t)$

$$0 = \frac{d}{dt} f(\vec{r}(t)) = \nabla f(\vec{r}(t)) \cdot \vec{r}'(t)$$

Steapest ascent

# Steapest Ascent

$\nabla f$ is the direction in which f increases most

# Proof

$$\frac{d}{dt} f(\vec{r}(t)) = \nabla f(\vec{r}(t)) \cdot \vec{r}\,'(t) = |\nabla f(\vec{r}(t))| \cdot |\vec{r}\,'(t)| \cos(\theta)$$

is maximal for $\cos(\theta) = 1$

---

**Algorithm 1:** Gradient ascent with backtracking line search

1  initialize $\lambda = \lambda_0$, $\theta_{old} = \theta_0$, define $\epsilon$, $c$, $\tilde{f}(\theta)$
2  **while** $|\tilde{f}'(\theta_{old})| > \epsilon$ **do**

3  $\quad$ compute $\begin{cases} \tilde{f}(\theta_{old}) & \text{(Function value)} \\ \tilde{f}'(\theta_{old}) & \text{(1}^{st}\text{ derivative)} \end{cases}$

4  $\quad$ compute $\theta_{new} = \theta_{old} + \underbrace{\lambda \tilde{f}'(\theta_{old})}_{d}$

5  $\quad$ **if** $\tilde{f}(\theta_{new}) < \tilde{f}(\theta_{old}) + c \cdot d\tilde{f}'(\theta_{old})$ **then**
6  $\quad\quad$ assign $\lambda = \lambda/2$ $\quad$ (Backtracking)
7  $\quad\quad$ Goto 4

8  $\quad$ assign $\lambda = \lambda_0$, $\theta_{old} = \theta_{new}$

9  $\theta^* = \theta_{old}$

---



loop #1
$\theta_{old} = 3.50$, $\lambda_0 = 3$
$\tilde{f}(\theta_{old}) = 0.21$
$\tilde{f}'(\theta_{old}) = 0.32$
$\theta_{new} = \theta_{old} + \lambda_0 \tilde{f}'(\theta_{old})$
$= 4.47$
$\tilde{f}(\theta_{new}) = 0.82$ △

loop #2
$\theta_{old} = 4.47$, $\lambda_0 = 3$
$\tilde{f}(\theta_{old}) = 0.82$
$\tilde{f}'(\theta_{old}) = 0.72$
$\theta_{new} = \theta_{old} + \lambda_0 \tilde{f}'(\theta_{old})$
$= 6.64$
$\tilde{f}(\theta_{new}) = -0.03$ △
$\theta_{new} = \theta_{old} + \frac{\lambda_0}{2} \tilde{f}'(\theta_{old})$
$= 5.55$
$\tilde{f}(\theta_{new}) = 0.70$ △
$\theta_{new} = \theta_{old} + \frac{\lambda_0}{4} \tilde{f}'(\theta_{old})$
$= 5.01$
$\tilde{f}(\theta_{new}) = 1.01$ △

Figure 5.4: Example of application of gradient ascent with backtracking for finding the maximum of a function.

Figure 5.4 presents the first two steps of the application of algorithm 1 to a non-convex/non-concave function with an initial value $\theta_0 = 3.5$ and a scaling factor $\lambda_0 = 3$. For the second step, the scaling factor $\lambda$ has to be reduced twice in order to satisfy the Armijo rule. One of the difficulties with gradient ascent is that the convergence speed depends on the choice of $\lambda_0$. If $\lambda_0$ is too small, several steps will be wasted and convergence will be slow. If $\lambda_0$ is too large, the algorithm may not converge.

Figure 5.5 presents a limitation common to all convex optimization methods when applied to functions involving local maxima; if the starting location $\theta_0$ is not located on the slope segment leading to the global maximum, the algorithm will most likely miss it and converge to a local maximum. The task of selecting a proper value $\theta_0$ is nontrivial because in most cases, it is not possible to visualize $\tilde{f}(\theta)$. This issue can be tackled by attempting multiple starting locations $\theta_0$ and by using domain knowledge to identify proper starting locations.

Gradient ascent can be applied to search for the maximum of a multivariate function by replacing the univariate derivative by the gradient so that

$$\theta_{new} = \theta_{old} + \lambda \cdot \nabla_\theta \tilde{f}(\theta_{old}).$$

As illustrated in figure 5.6, because gradient ascent follows the direction where the gradient is maximal, it often displays an oscillatory pattern. This issue can be mitigated by introducing a *momentum* term in the calculation of $\theta_{new}$,[3]

$$\begin{aligned} \mathbf{v}_{new} &= \gamma \cdot \mathbf{v}_{old} + \lambda \cdot \nabla_\theta \tilde{f}(\theta_{old}), \\ \theta_{new} &= \theta_{old} + \mathbf{v}_{new} \end{aligned}$$

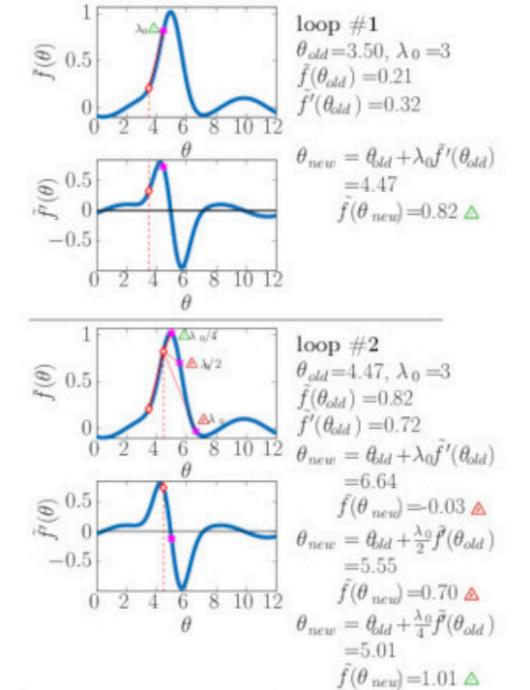where $\mathbf{v}$ can be interpreted as a velocity that carries the momentum from the previous iterations.
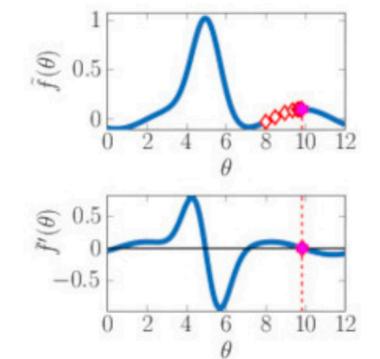


Figure 5.5: Example of application of gradient ascent converging to a local maximum for a function.
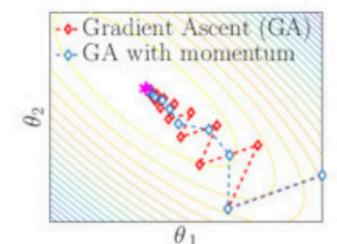


Figure 5.6: Comparison of gradient ascent with and without momentum.

[3] Rumelhart, D. E., G. E. Hinton, and R. J. Williams (1986). Learning representations by back-propagating errors. *Nature 323*, 533–536.

## 5.1 Gradient Ascent

A *gradient* is a vector containing the partial derivatives of a function with respect to its variables. For a continuous function, the maximum is located at the point where its gradient equals zero. *Gradient ascent* is based on the principle that as long as we move in the direction of the gradient, we are moving toward a maximum. For the unidimensional case, we choose to move to a new position $\theta_{\text{new}}$ defined as the old value $\theta_{\text{old}}$ plus a search direction $d$ defined by a scaling factor $\lambda$ times the derivative estimated at $\theta_{\text{old}}$,

$$\theta_{\text{new}} = \theta_{\text{old}} + \underbrace{\lambda \cdot \tilde{f}'(\theta_{\text{old}})}_{d}.$$

A common practice for setting $\lambda$ is to employ *backtracking line search* where a new position is accepted if the *Armijo rule*[2] is satisfied so that

$$\tilde{f}(\theta_{\text{new}}) \geq \tilde{f}(\theta_{\text{old}}) + c \cdot d\tilde{f}'(\theta_{\text{old}}), \quad \text{with } c \in (0,1). \tag{5.2}$$

Figure 5.3 presents a comparison of the application of equation 5.2 with the two extreme cases, $c = 0$ and $c = 1$. For $c = 1$, $\theta_{\text{new}}$ is

**Derivative**

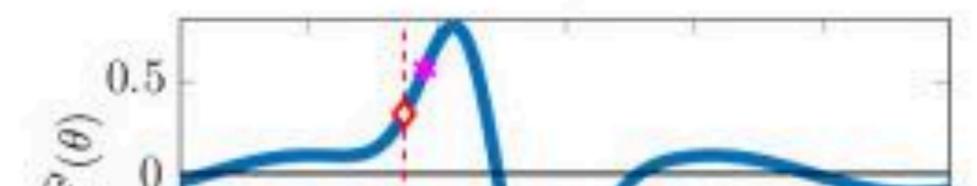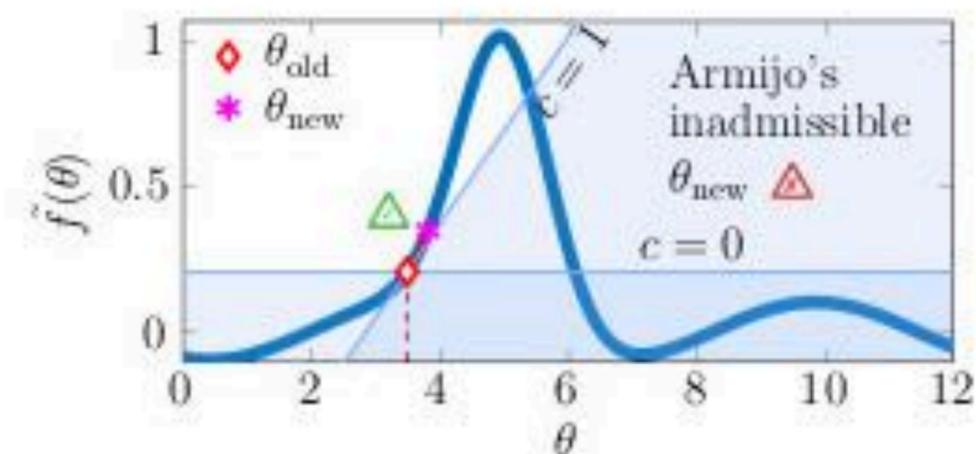$$\tilde{f}'(\theta) \equiv \frac{d\tilde{f}(\theta)}{d\theta}$$

**Gradient**

$$\nabla \tilde{f}(\theta) \equiv \nabla_{\boldsymbol{\theta}} \tilde{f}(\boldsymbol{\theta})$$
$$= \left[ \frac{\partial \tilde{f}(\boldsymbol{\theta})}{\partial \theta_1} \quad \frac{\partial \tilde{f}(\boldsymbol{\theta})}{\partial \theta_2} \quad \cdots \quad \frac{\partial \tilde{f}(\boldsymbol{\theta})}{\partial \theta_n} \right]^{\mathsf{T}}$$

**Maximum of a concave function**

$$\theta^* = \arg \max_{\theta} \tilde{f}(\theta) : \frac{d\tilde{f}(\theta^*)}{d\theta} = 0$$

$\lambda$ is also known as the *learning rate* or *step length*.

[2] Armijo, L. (1966). Minimization of functions having Lipschitz continuous first partial derivatives. *Pacific Journal of Mathematics 16*(1), 1–3

Mt.
Ephraim
Reservoir
Turkey
Hill
ARLINGTON
ARLINGTON
Arlington
Heights
Arlington
MED
Med
Brook

Turkey
Hill

ARLINGTON

RLINGTON

Arlington
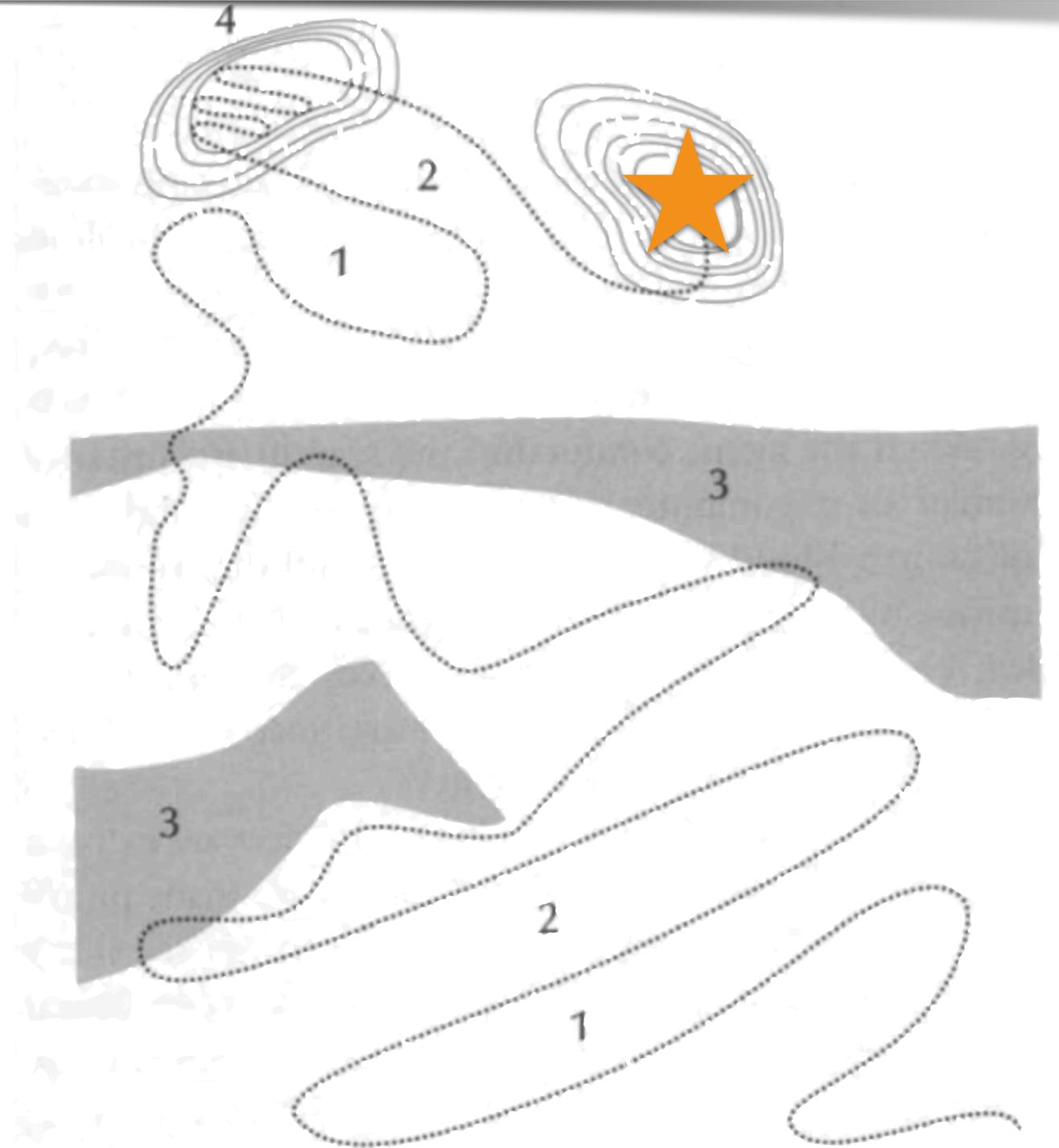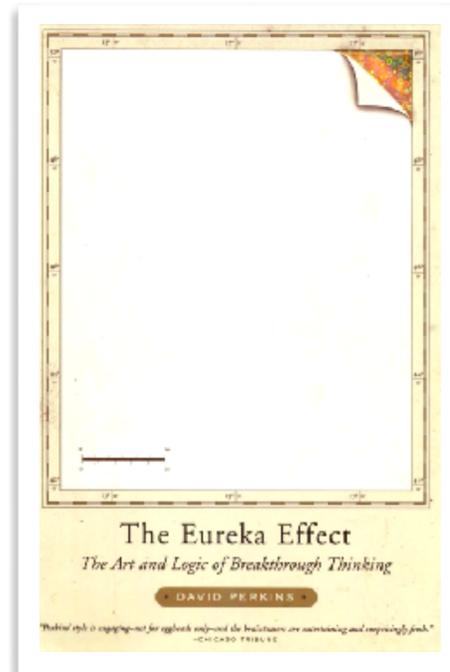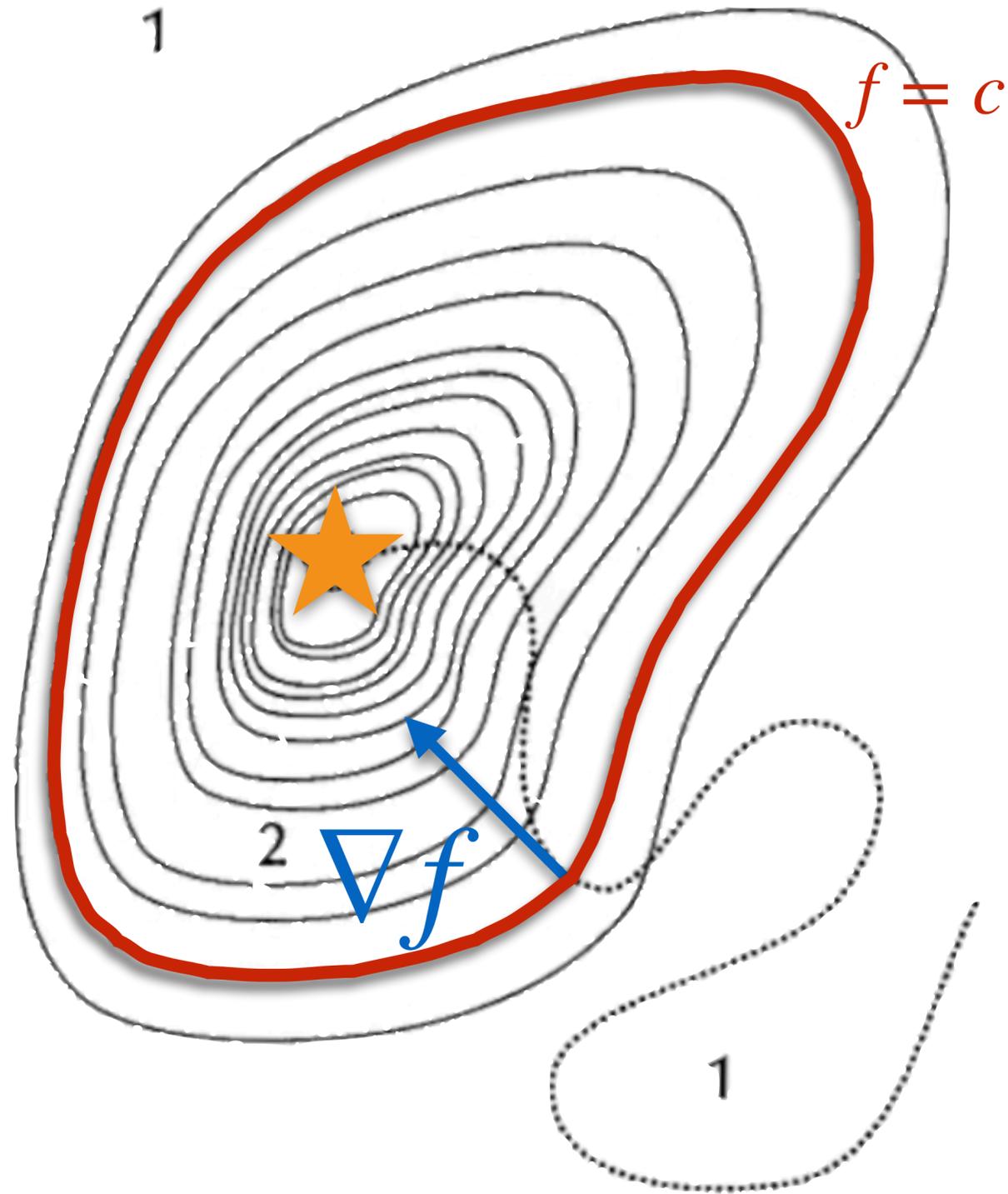Heights

Arlingt

Grindelwald

Eiger

Challifirn

nördliches

ALCOHOL & CALCULUS DON'T MIX. NEVER DRINK & DERIVE.

$$r(t) = x + t\,\nabla f(x)$$

$$\frac{d}{dt}f(x) = \nabla f(x) \cdot \nabla f(x)$$

1

$f = c$

$\nabla f$

2

1

The Eureka Effect
The Art and Logic of Breakthrough Thinking
DAVID PERKINS

Search in a Homing Space: 1. Clueless regions.
2. Large clued regions leading to the target.

4

2

1

3

3

2

1

Search in a Klondike Space: 1. A large space wtih few solutions (a wilderness trap). 2. Regions with no clues pointing direction (plateau traps). 3. A barrier isolates the solution (creating a canyon trap). 4. An area of high promise but no solution (an oasis trap).
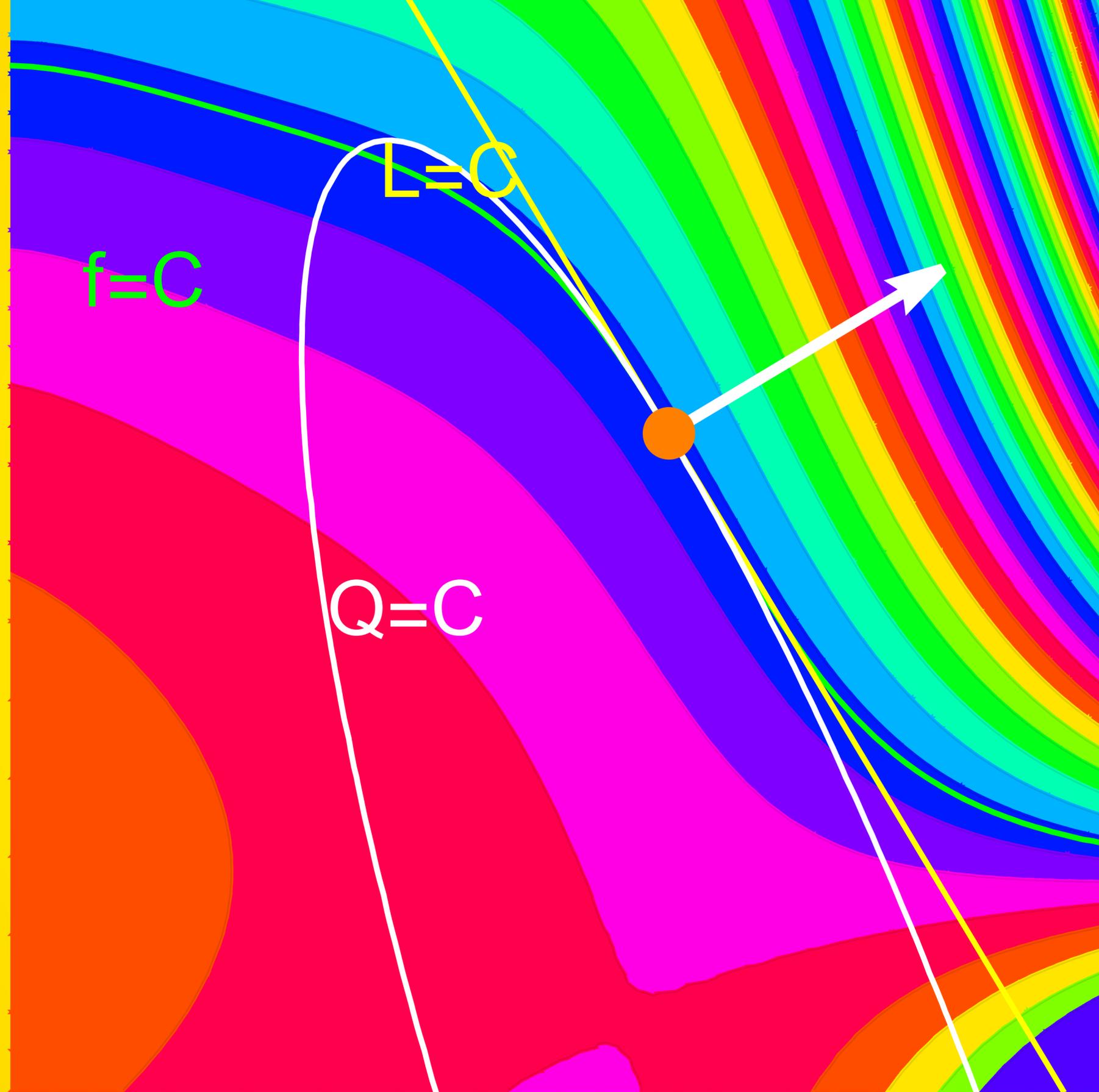
# Quadratic Approximation

# Quadratic Approximation

$$Q(x,y) = f(x_0, y_0) + f_x(x_0, y_0)(x - x_0) + f_y(x_0, y_0)(y - y_0)$$

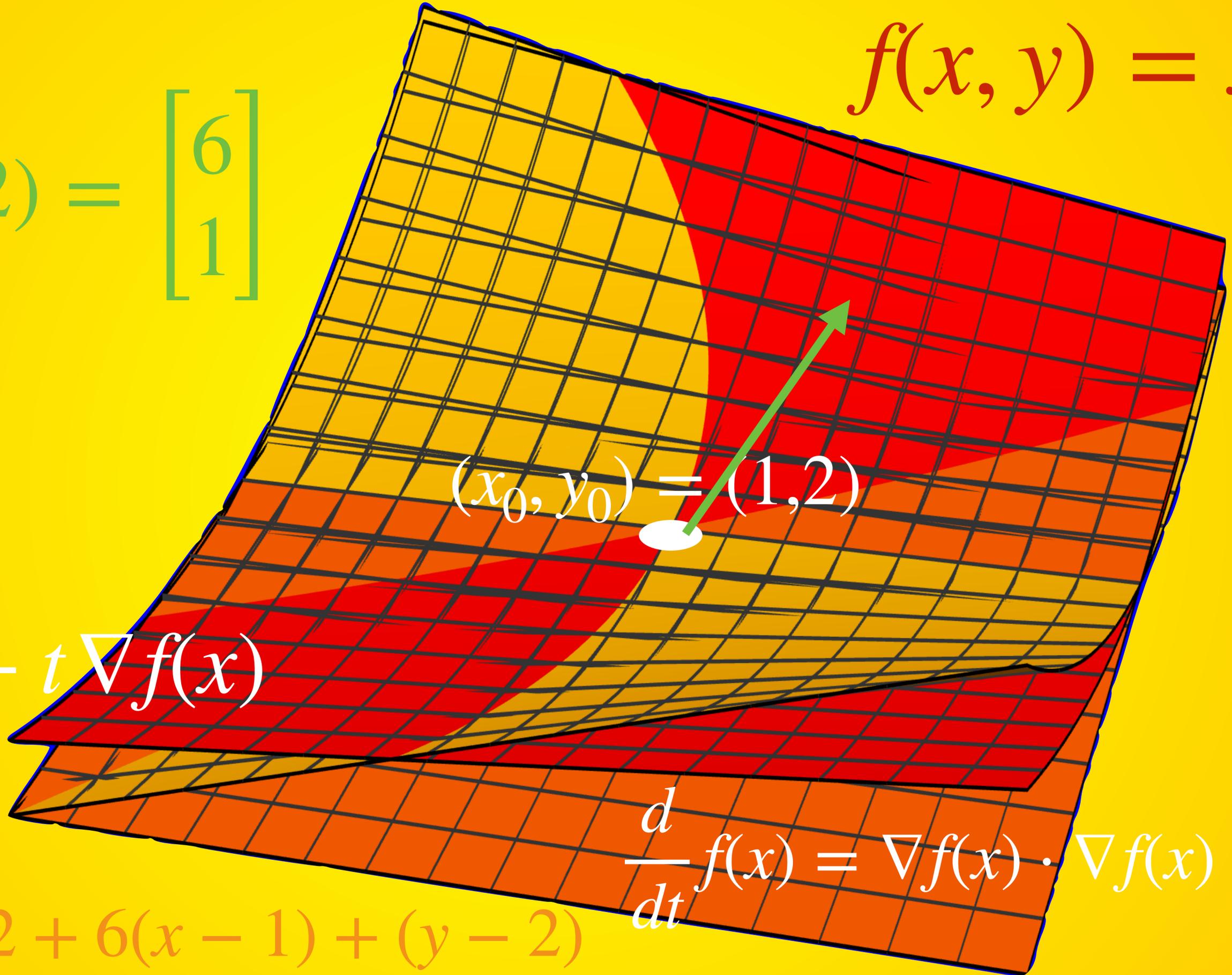$$\frac{f_{xx}(x_0, y_0)(x - x_0)^2 + f_{yy}(x_0, y_0)(y - y_0^2) + 2f_{xy}(x_0, y_0)(x - x_0)(y - y_0)}{2}$$

The best quadratic function near the point $(x_0, y_0)$

$$f(x, y) = x^3 y$$
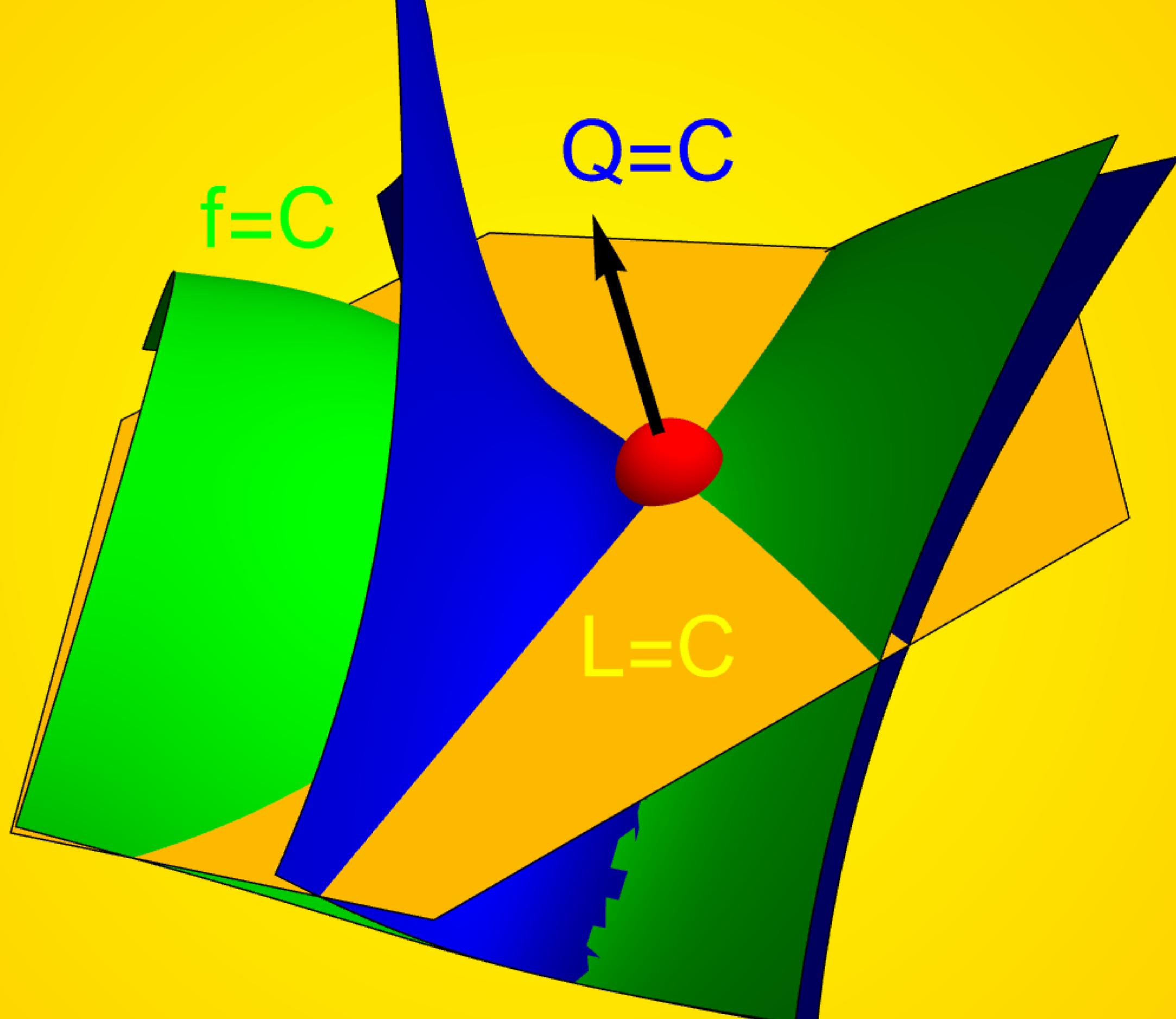
$$\nabla f(1,2) = \begin{bmatrix} 6 \\ 1 \end{bmatrix}$$

$$(x_0, y_0) = (1,2)$$

$$r(t) = x + t\nabla f(x)$$

$$\frac{d}{dt}f(x) = \nabla f(x) \cdot \nabla f(x)$$

$$L(x, y) = 2 + 6(x - 1) + (y - 2)$$

THE END