

# **COMP0130: Robotic Vision and Navigation**

## Lecture 06D: Inference and Optimization

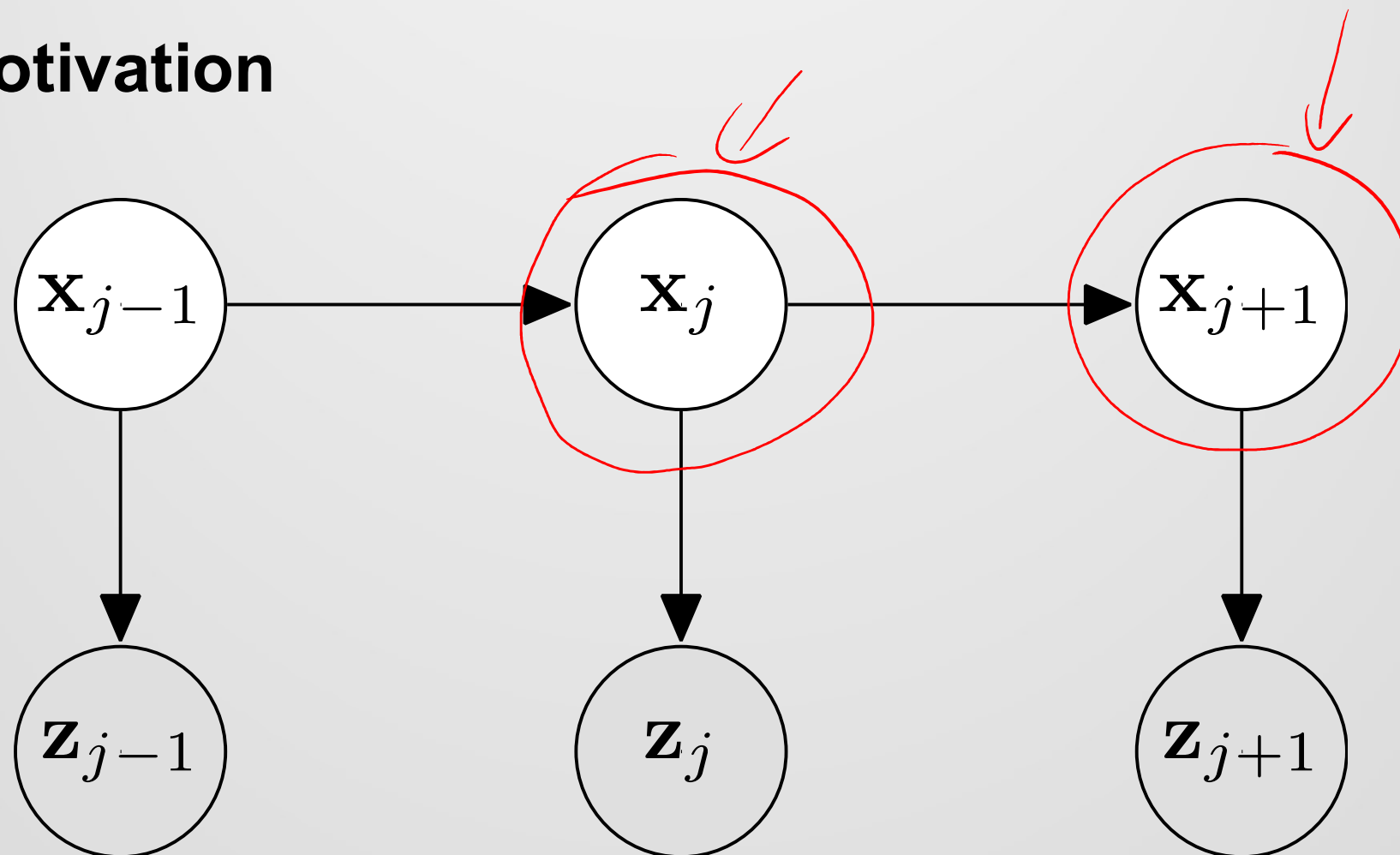
# Structure

- Motivation
- The Problem with Means
- Maximum A posteriori Estimation
- Optimization Algorithms
- Approximate Covariance Calculations

## Motivation

- The factor graph stores the entire probability density in a single equation
- Theoretically, we can ask any question we like of it and compute the answer
- In robotics, we are often interested in computing a point estimate which we can use in control

# Motivation



# Expected Values

- *Motivation*
- Expected Values
- *Maximum A posteriori Estimation*
- *Optimization Algorithms*
- *Approximate Covariance Calculations*

## Computing the Mean

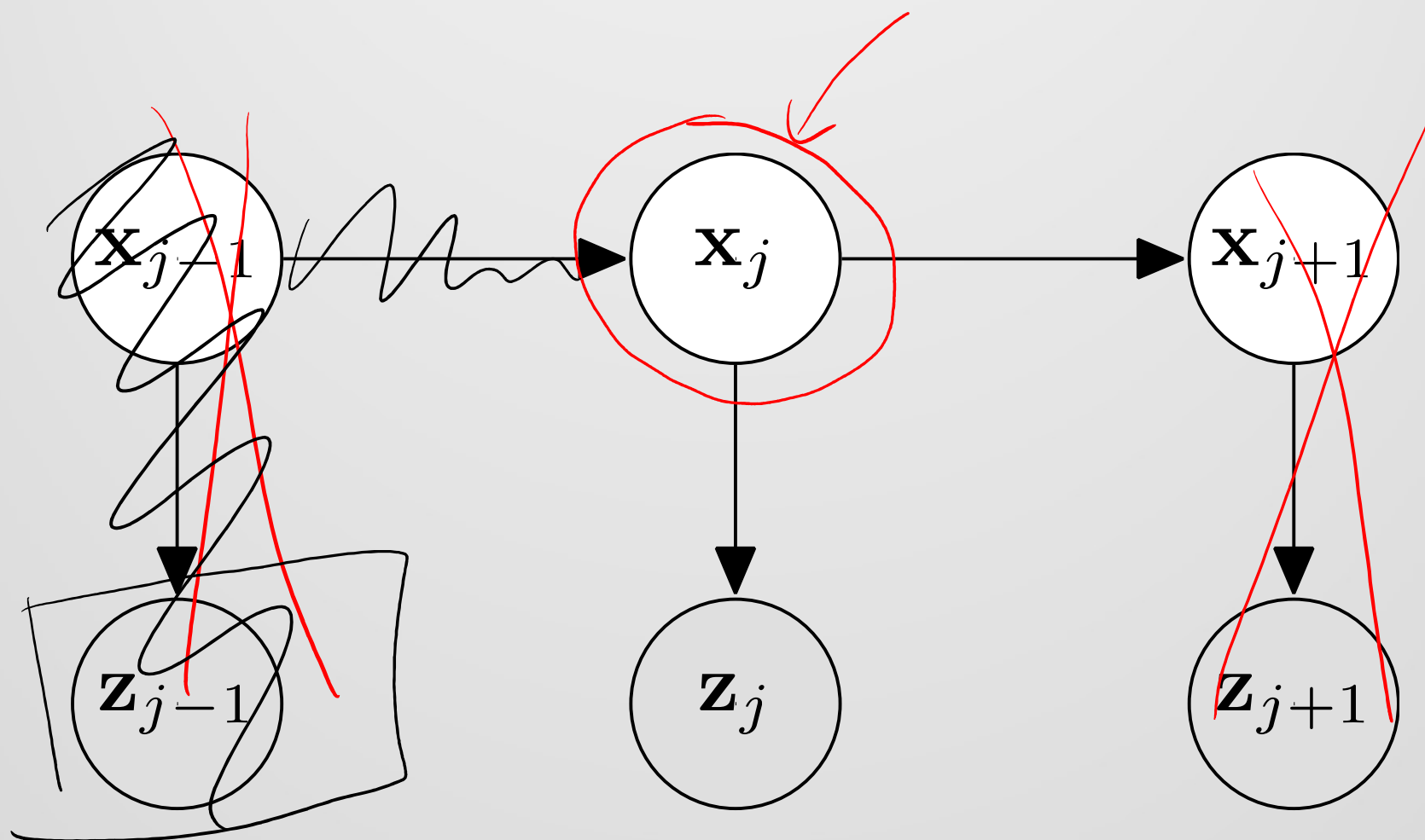
- Suppose we want to compute the mean of the  $j$ th state
- The mean is given by

$$\mathbb{E}[\mathbf{x}_j] = \int \mathbf{x}_j f(\mathbf{x}_j | \mathbb{I}_k) d\mathbf{x}_j$$

*Handwritten annotations:* An arrow points from  $\mathbf{x}_j$  to  $\mathbb{E}[\mathbf{x}_j]$ . Another arrow points from  $\mathbf{x}_j$  to  $\mathbf{x}_j$  in the integrand. A third arrow points from  $\mathbf{x}_j$  to  $d\mathbf{x}_j$ .

- To do this, we have to marginalize

# Computing the Mean



# Marginalization

- Marginalization is the “proper” way to remove variables from a graph
- We do this by integrating all the other variables out

$$f(\mathbf{x}_j | \mathbb{I}_k) = \int \underbrace{f(\mathbf{x}_{1:k} | \mathbb{I}_k)}_{\text{joint distribution}} d\mathbf{x}_{\underbrace{1:k \setminus j}_{\text{variables to integrate out}}}$$

$$\underbrace{1:k \setminus j}_{\text{variables to integrate out}} = [1, 2, \dots, j-1, j+1, \dots, k]$$

$\uparrow \quad \uparrow \quad \quad \quad \uparrow \quad \quad \uparrow \quad \quad \quad \uparrow$



## Computed Expected Values

- More generally, we can specify any kind of function,

$$\mathbb{E} [\underbrace{\mathbf{b}}^{\swarrow} [\underbrace{\mathbf{x}_{1:k}}^{\swarrow}]] = \int \underbrace{\mathbf{b}}^{\swarrow} [\mathbf{x}_{1:k}] f(\mathbf{x}_{1:k} | \mathbb{I}_k) d\mathbf{x}_{1:k}$$

- This can be used to compute things like covariances and cross correlations as well

## Computing the Mean

- However, performing the marginalization is awful
- We have to do all the integration work we were desperately trying to avoid
- Therefore, we need to look at other approaches

# Structure

- *Motivation*
- *The Problem with Means*
- Maximum A posteriori Estimation      MAP
- *Optimization Algorithms*
- *Approximate Covariance Calculations*

# Maximum A Posteriori Estimation

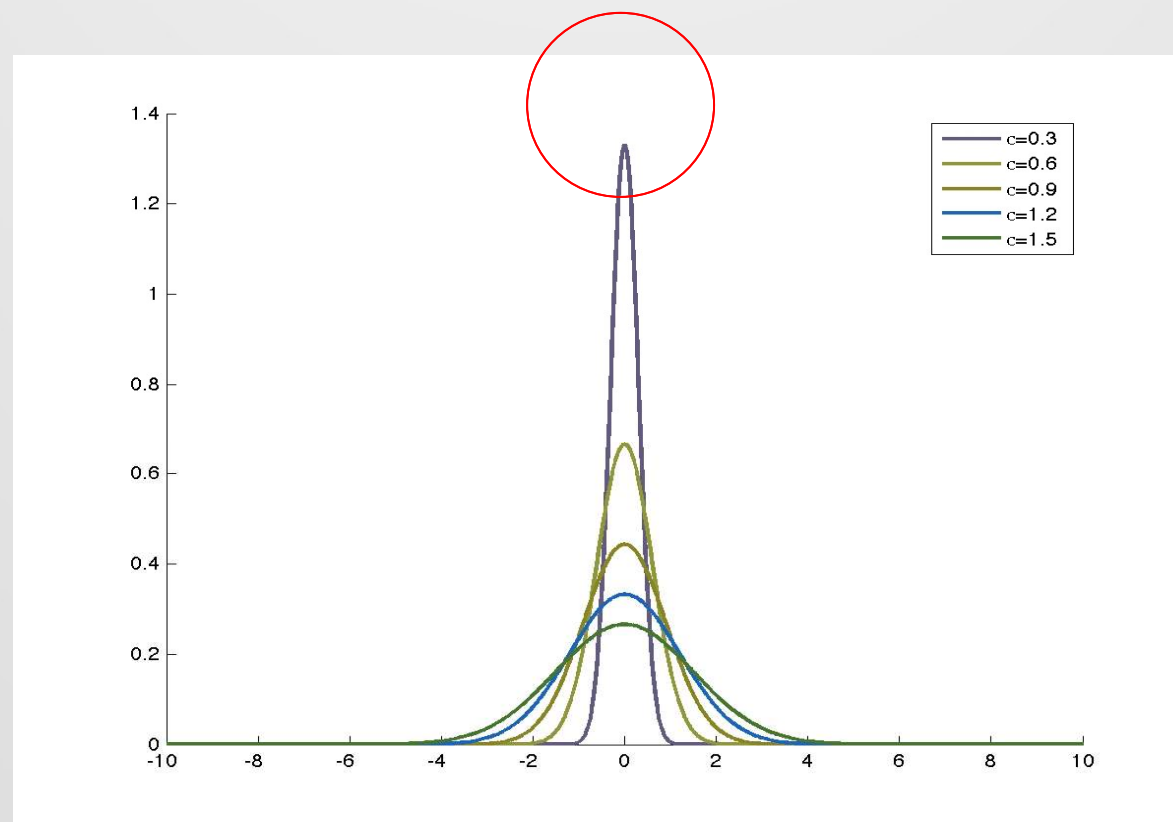
- The problem with estimating the mean is that we have to integrate over the entire distribution, which we can't do in practice
- An alternative idea is to just extract information directly from the distribution without integration
- One particularly easy approach is to do peak finding

# Maximum A Posteriori Estimation

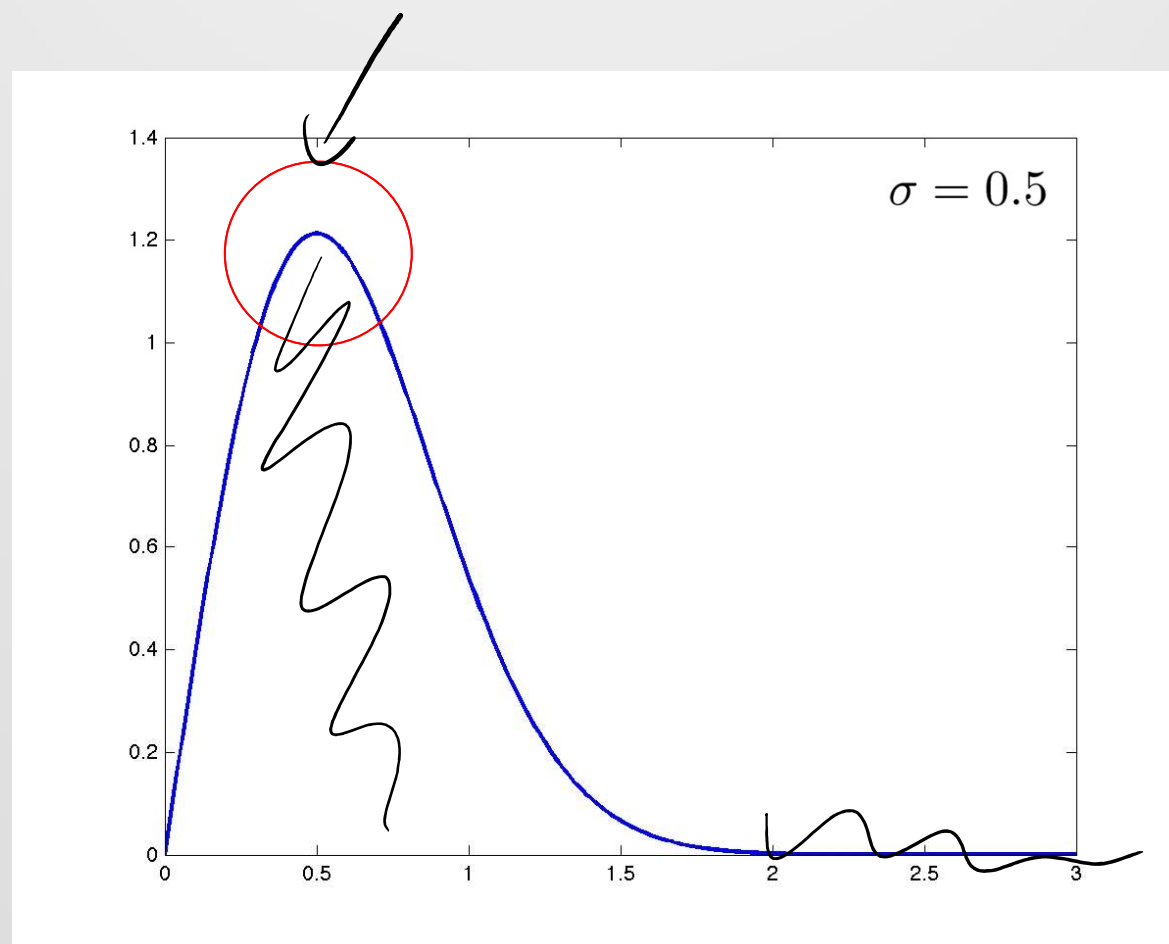
- Define the output estimate to be

$$\underset{\hat{\mathbf{x}}}{\mathbf{x}}^* = \arg \max_{\hat{\mathbf{x}}} f(\mathbf{x})$$

# MAP of a Gaussian PDF



# MAP of a Rayleigh PDF



# Maximum A Posteriori Estimation

- Define the output estimate to be

$$\mathbf{x}^* = \arg \max_{\mathbf{x}} f(\mathbf{x})$$

- Advantages:
  - We do not need to know the normalization constant (no integration!)
  - The estimate itself is often pretty sensible



# Normalization Constant not Required

- Suppose we have the solution

$$\mathbf{x}^* = \arg \max_{\mathbf{x}} \underline{f(\mathbf{x})}$$

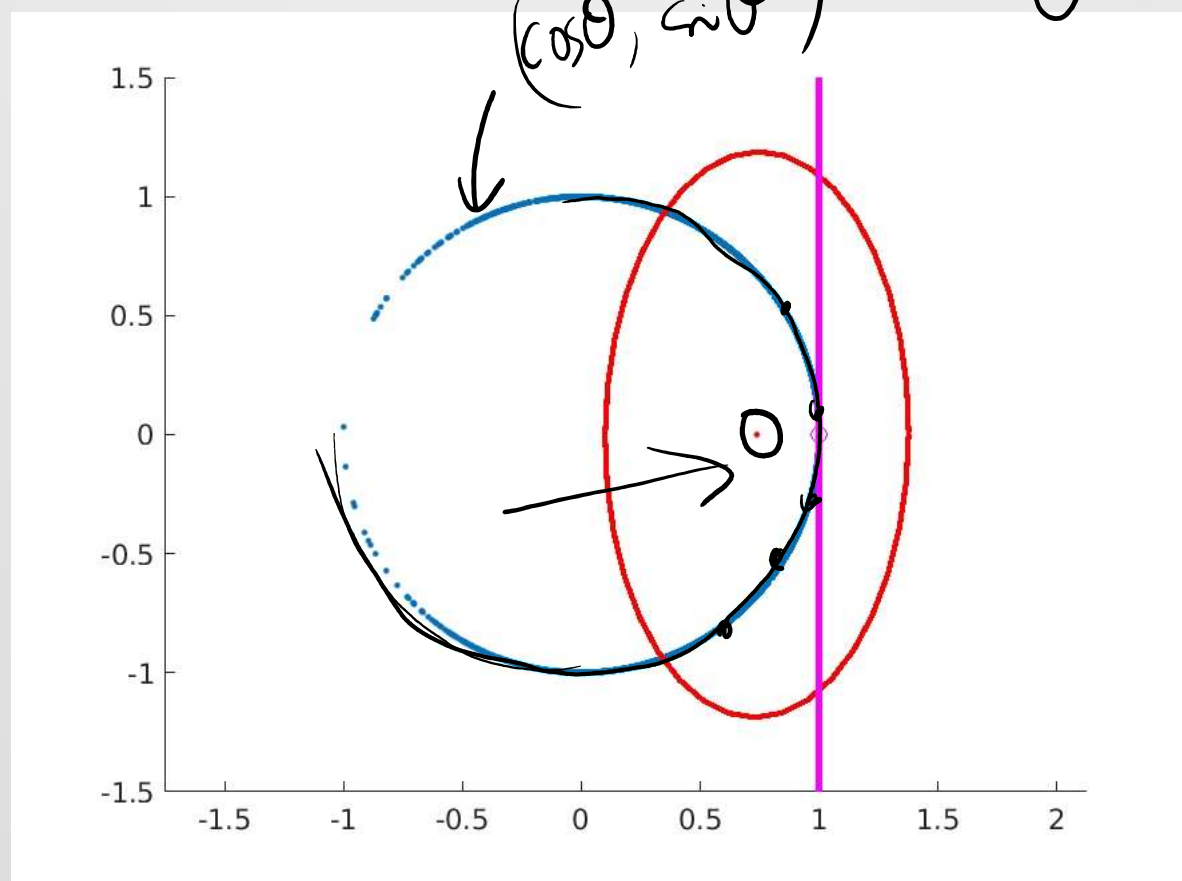
- The value is the same as

$$\mathbf{x}^* = \arg \max_{\mathbf{x}} c f(\mathbf{x})$$

$$c > 0$$

$$c = \int f(\mathbf{z}) d\mathbf{x} \geq 0$$

# MAP Estimates Tend to be Sensible



# Computing the MAP Estimate

- We want to compute the maximum of the function

$$f'(\mathbf{x}_{1:k} | \mathbb{I}_k) = f(\mathbf{x}_0) \prod_{i=1}^k f(\mathbf{x}_i | \mathbf{x}_{i-1}, \mathbf{u}_i) \times \prod_{i=1}^k L(\mathbf{x}_i; \mathbf{z}_i)$$

Handwritten annotations:

- A downward arrow points from  $\mathbf{x}_{1:k}$  to  $\mathbf{x}_1$ .
- A handwritten note "argmax  $\mathbf{x}_{1:n}$ " is written next to the first term.
- A handwritten note " $\geq 0$ " is written above the product term.
- Two arrows point from the product terms to the right, indicating a relationship or comparison.

## Computing the MAP Estimate

- Just directly optimizing this can lead to numerical over or underflow issues
- However, the form of this function is the product of terms, each of which is non-negative
- Therefore, we can take the log of the distribution

# MAP and Logs

$$\begin{array}{l} x_1 > x_2 \\ \log x_1 > \log x_2 \end{array}$$



- The log of a function has the property that it is asymptotically increasing
- Therefore,

$$\begin{aligned} \mathbf{x}^* &= \arg \max_{\mathbf{x}} \underline{f'(\mathbf{x})} \\ &= \arg \max_{\mathbf{x}} \ln f'(\mathbf{x}) \end{aligned}$$

# Computing the MAP Estimate

- Substituting,

$$\ln f'(\mathbf{x}_{0:k} | \mathbb{I}_k) = \ln f(\mathbf{x}_0) + \sum_{i=1}^k \ln f(\mathbf{x}_i | \mathbf{x}_{i-1}, \mathbf{u}_i) + \sum_{i=1}^k \ln L(\mathbf{x}_i; \mathbf{z}_i)$$

Handwritten arrows indicate dependencies: a downward arrow from  $\mathbb{I}_k$  to the first sum, a downward arrow from  $k$  to the first sum, an upward arrow from  $\mathbf{u}_i$  to the first sum, a downward arrow from  $k$  to the second sum, and an upward arrow from  $\mathbf{z}_i$  to the second sum.

## State Transition Probabilities and Likelihoods

- We now need to substitute for the expressions for the transition density and likelihood functions

$$f(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{u}_k) = f_{\mathbf{v}}(\mathbf{v}_k \stackrel{\swarrow}{=} \mathbf{e}[\mathbf{x}_k, \mathbf{x}_{k-1}, \mathbf{u}_k])$$

$$f(\mathbf{z}_k | \mathbf{x}_k) = f_{\mathbf{w}}(\mathbf{w}_k \stackrel{\nwarrow}{=} \mathbf{l}[\mathbf{x}_k, \mathbf{z}_k])$$



## Linear Case

- For the sensor likelihood equations, these are

$$L(\mathbf{x}_k; \mathbf{z}_k) \propto \exp \left\{ -\frac{1}{2} \underbrace{\mathbf{l}[\mathbf{x}_k, \mathbf{z}_k]}_{\substack{\omega \uparrow \\ \mathbf{z} - \mathbf{h}_x}}^\top \mathbf{R}_k^{-1} \mathbf{l}[\mathbf{x}_k, \mathbf{z}_k] \right\}$$

- Taking logs,

$$\ln L(\mathbf{x}_k; \mathbf{z}_k) = -\frac{1}{2} \mathbf{l}[\mathbf{x}_k, \mathbf{z}_k]^\top \mathbf{R}_k^{-1} \mathbf{l}[\mathbf{x}_k, \mathbf{z}_k] + \underbrace{c}_{\downarrow}$$



## Linear Case

- For the state transition equations, we can substitute for the inverse process model to get

$$f(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{u}_k) \propto \exp \left\{ -\frac{1}{2} \underbrace{\mathbf{e}[\mathbf{x}_k, \mathbf{x}_{k-1}, \mathbf{u}_k]}_{\mathbf{e} = \mathbf{x}_k - \mathbf{F}\mathbf{x}_{k-1}}^\top \mathbf{Q}_k^{-1} \mathbf{e}[\mathbf{x}_k, \mathbf{x}_{k-1}, \mathbf{u}_k] \right\}$$

- Taking logs,

$$\ln f(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{u}_k) = -\frac{1}{2} \mathbf{e}[\mathbf{x}_k, \mathbf{x}_{k-1}, \mathbf{u}_k]^\top \mathbf{Q}_k^{-1} \mathbf{e}[\mathbf{x}_k, \mathbf{x}_{k-1}, \mathbf{u}_k] + d$$

## Putting Everything Together

- Substituting Everything together, we get:

$$\ln f'(\mathbf{x}_{1:k} | \mathbb{I}_k) = \underbrace{-\frac{1}{2} \sum_{i=1}^k \mathbf{e}[\mathbf{x}_k, \mathbf{x}_{k-1}, \mathbf{u}_k]^\top \mathbf{Q}_k^{-1} \mathbf{e}[\mathbf{x}_k, \mathbf{x}_{k-1}, \mathbf{u}_k]}_{\text{first term}} + \underbrace{\left( -\frac{1}{2} \sum_{i=1}^k \mathbf{l}[\mathbf{x}_k, \mathbf{z}_k]^\top \mathbf{R}_k^{-1} \mathbf{l}[\mathbf{x}_k, \mathbf{z}_k] + e \right)}_{\text{second term}}$$

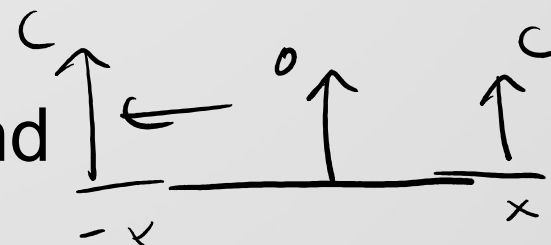
## Need for Optimization

- We want to compute the solution

$$\mathbf{x}^* = \arg \max_{\mathbf{x}} \ln f'(\mathbf{x})$$

- However, the solution includes lots of –ve signs which is a bit inconvenient
- We can flip signs and scale to find

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} -2 \ln f'(\mathbf{x})$$



# Putting Everything Together

- Therefore, the goal is to compute

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} \underbrace{c(\mathbf{x}_{1:k} | \mathbb{I}_k)}$$

where

$$c(\mathbf{x}_{1:k} | \mathbb{I}_k) = \sum_{i=1}^k \overbrace{\mathbf{e}[\mathbf{x}_k, \mathbf{x}_{k-1}, \mathbf{u}_k]^\top \mathbf{Q}_k^{-1} \mathbf{e}[\mathbf{x}_k, \mathbf{x}_{k-1}, \mathbf{u}_k]} + \sum_{i=1}^k \overbrace{\mathbf{l}[\mathbf{x}_k, \mathbf{z}_k]^\top \mathbf{R}_k^{-1} \mathbf{l}[\mathbf{x}_k, \mathbf{z}_k]} \quad \text{✗}$$

# Optimization Algorithms

- *Motivation*
- *Interpreting the Graph*
- *Maximum Likelihood Estimation*
- Optimization
- *Approximate Covariance Calculations*

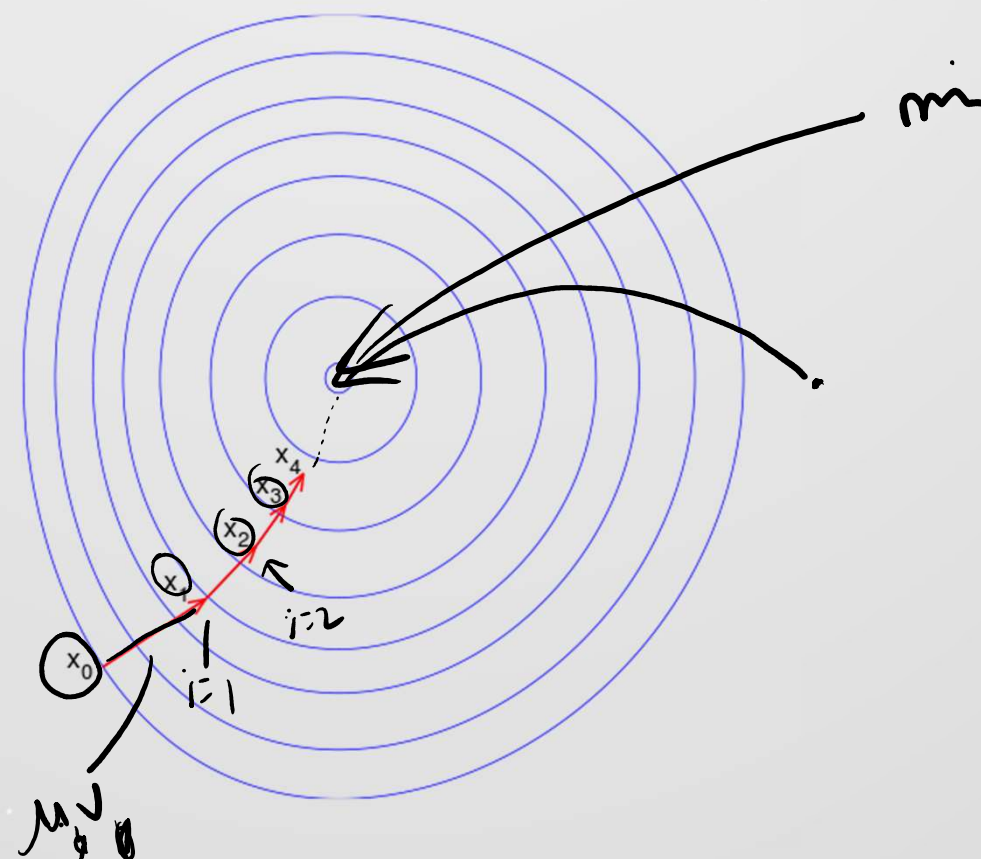
## Need for Optimization

- The MAP estimate is computed from an optimization process
- If the system is completely linear, you can show that the MAP estimate can be solved directly using matrix inverses
- However, in general the system is not linear
- Therefore, we are solving a nonlinear least squares estimation

# Nonlinear Least Squares Optimization

- Nonlinear systems rarely have a closed form solution
- Rather, we use an iterative approach which converges to the minimum

# Iterative Nonlinear Optimization





# Nonlinear Least Squares Optimization

- Nonlinear systems rarely have a closed form solution
- Rather, we use an iterative approach which converges to the minimum
- Since we clearly need more notation, we will write a solution in the optimizer as



# Iterative Nonlinear Optimization

$\mathbf{X}_k^0 \leftarrow$  Initial value  $\leftarrow$

$c^0 \leftarrow c(\mathbf{X}_k^0) \leftarrow$

$i \leftarrow 1$

**repeat**

$\mathbf{v}_k^i \leftarrow$  Step direction  $\leftarrow$

$\mu_k^i \leftarrow$  Step length  $\leftarrow$

$\mathbf{X}_k^i = \mathbf{X}_k^{i-1} + \mu_k^i \mathbf{v}_k^i$

$\rightarrow c^i \leftarrow c(\mathbf{X}_k^i)$

$i \leftarrow i + 1$

**until**  $|c^i - c^{i-1}| \leq \epsilon_c$  or  $|\mathbf{X}_k^i - \mathbf{X}_k^{i-1}| \leq \epsilon_X$

# Gradient Descent Optimization

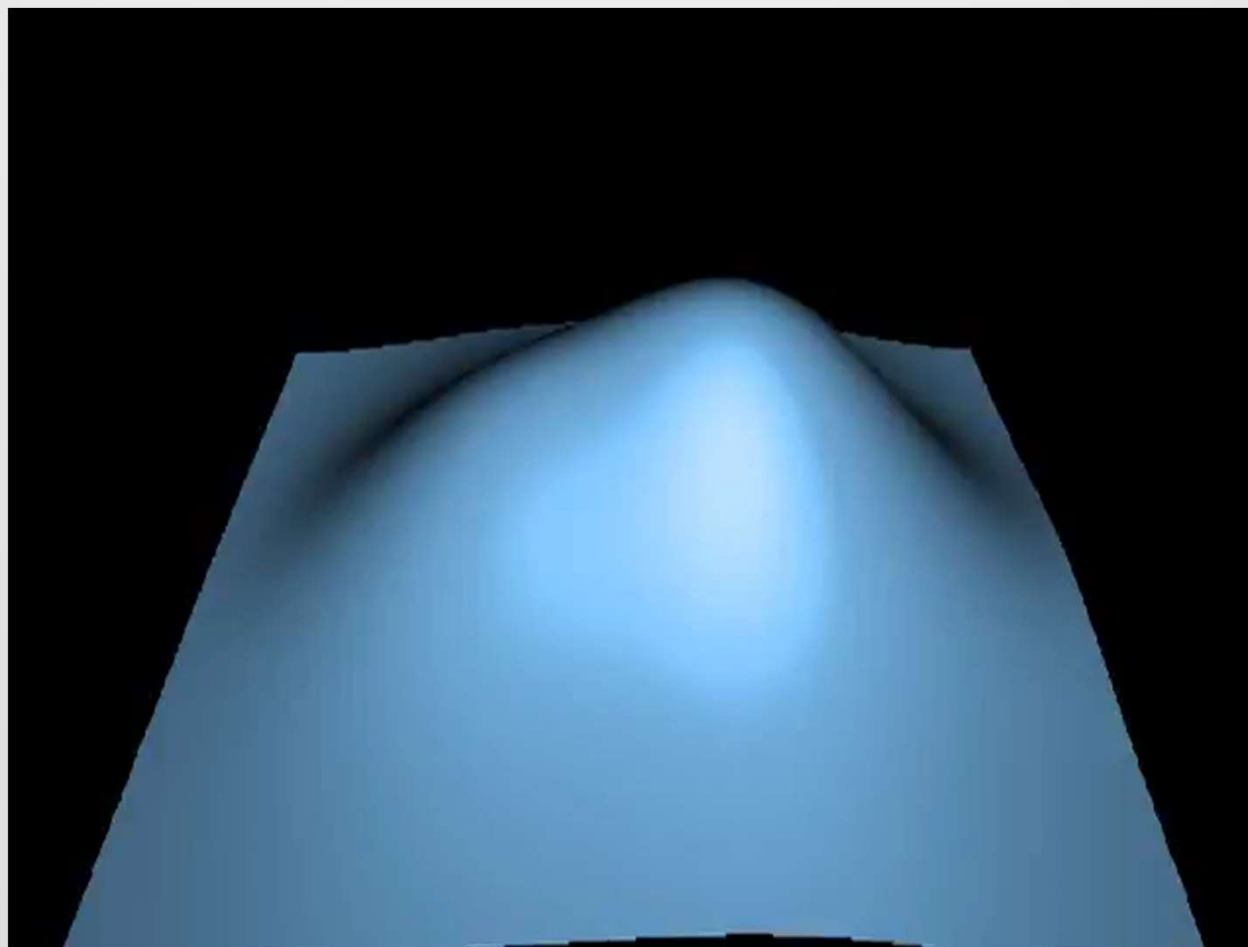
—  $\nabla$  —  $\frac{d}{dx}$

- Chooses the direction in which the cost function will fall the fastest
- The direction of greatest descent is the *negative* of the gradient of the cost function

$$\mathbf{X}_k^{i+1} = \mathbf{X}_k^i - \mathbf{J}(\mathbf{X}_k^i)$$

$$\mathbf{J} = \nabla c(\mathbf{X}_k^i) = \left[ \frac{\partial c(\mathbf{X}_k^i)}{\partial x_1} \quad \dots \quad \frac{\partial c(\mathbf{X}_k^i)}{\partial x_n} \right]^T$$

# Gradient Descent (Ascent!) at Work

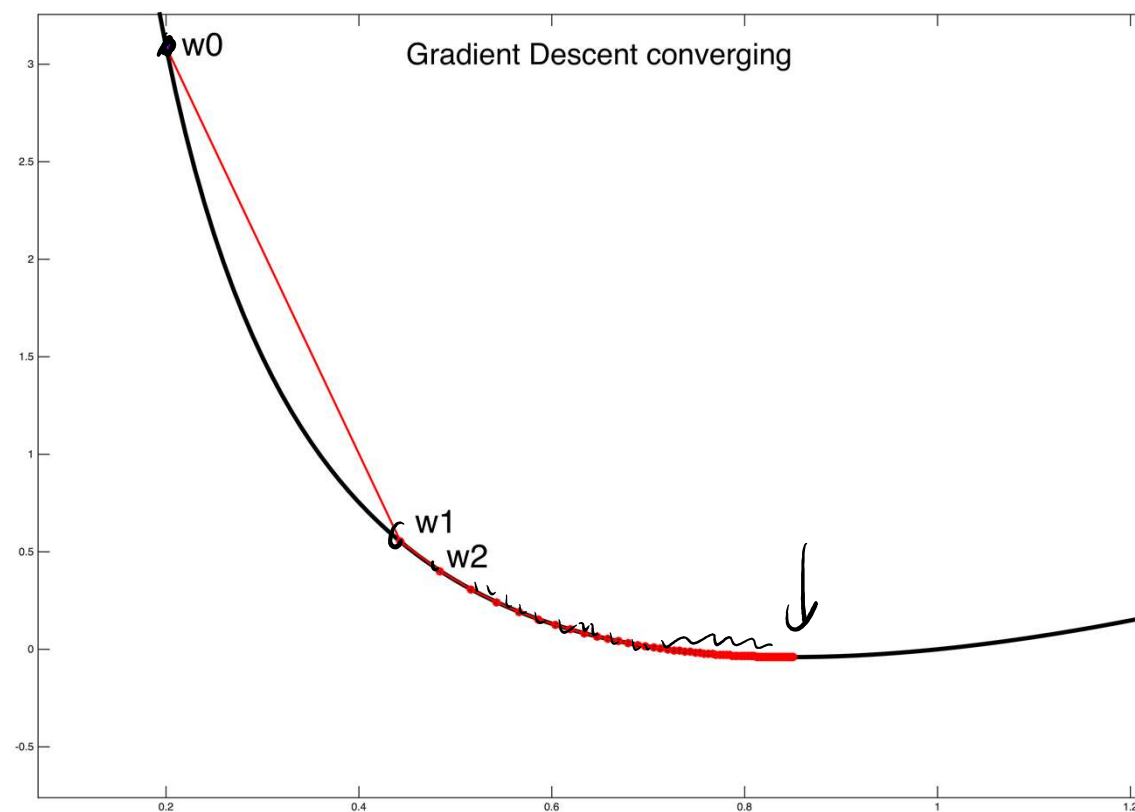


From <https://www.youtube.com/watch?v=L2W5SfGu09M>

## Effect of the Update

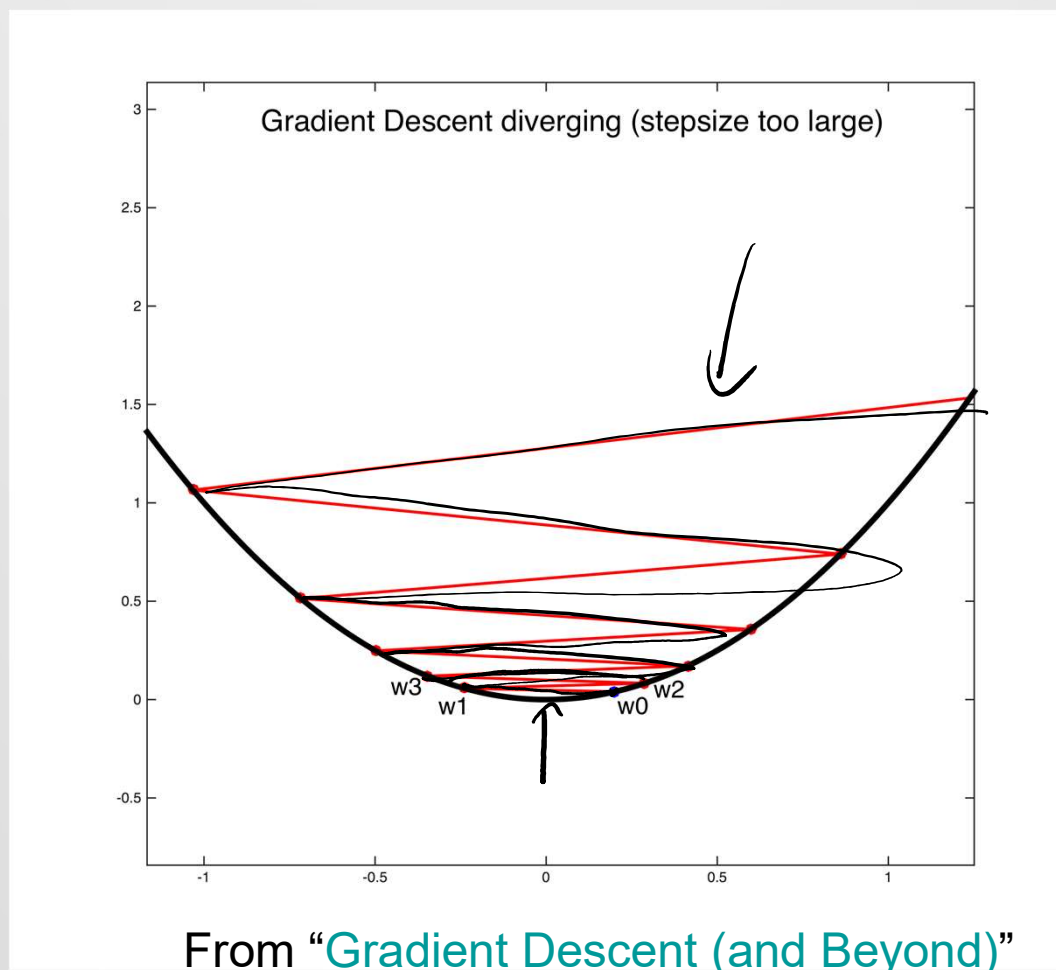


# Good Step Size



From "[Gradient Descent \(and Beyond\)](#)"

# Bad Step Size



$\mu = 1$

## Scaling the Step Size

- We could try just changing the optimization direction to find somewhere “better”
- However most approaches assume that we simply made our step size “too big”
- Therefore, take the scaled step

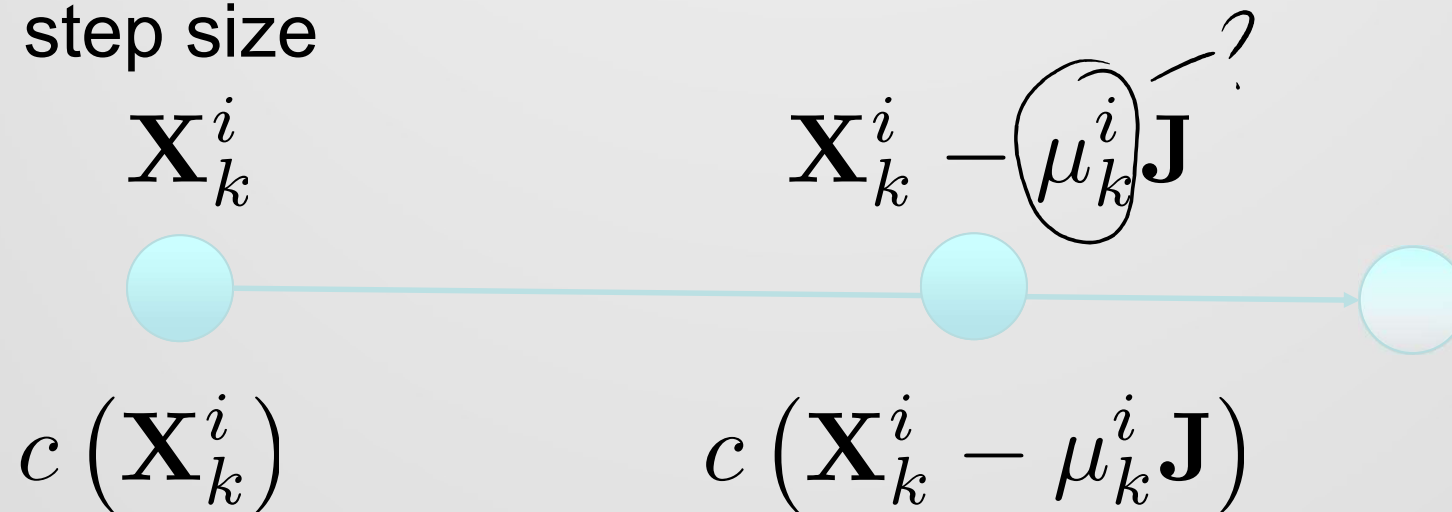
$$\mathbf{X}_k^{i+1} = \mathbf{X}_k^i - \mu_k^i \mathbf{J}(\mathbf{X}_k^i)$$

$$0 < \mu_k^i \leq \underbrace{1}_{\uparrow \downarrow}$$

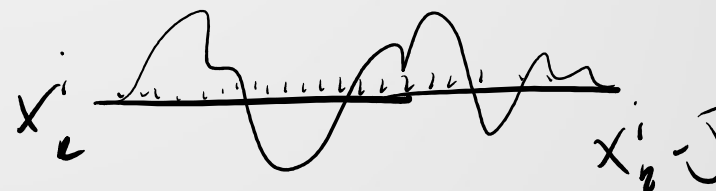


# Scaled Step Size

- We take a step along the gradient and change the step size



## Choosing the Step Size



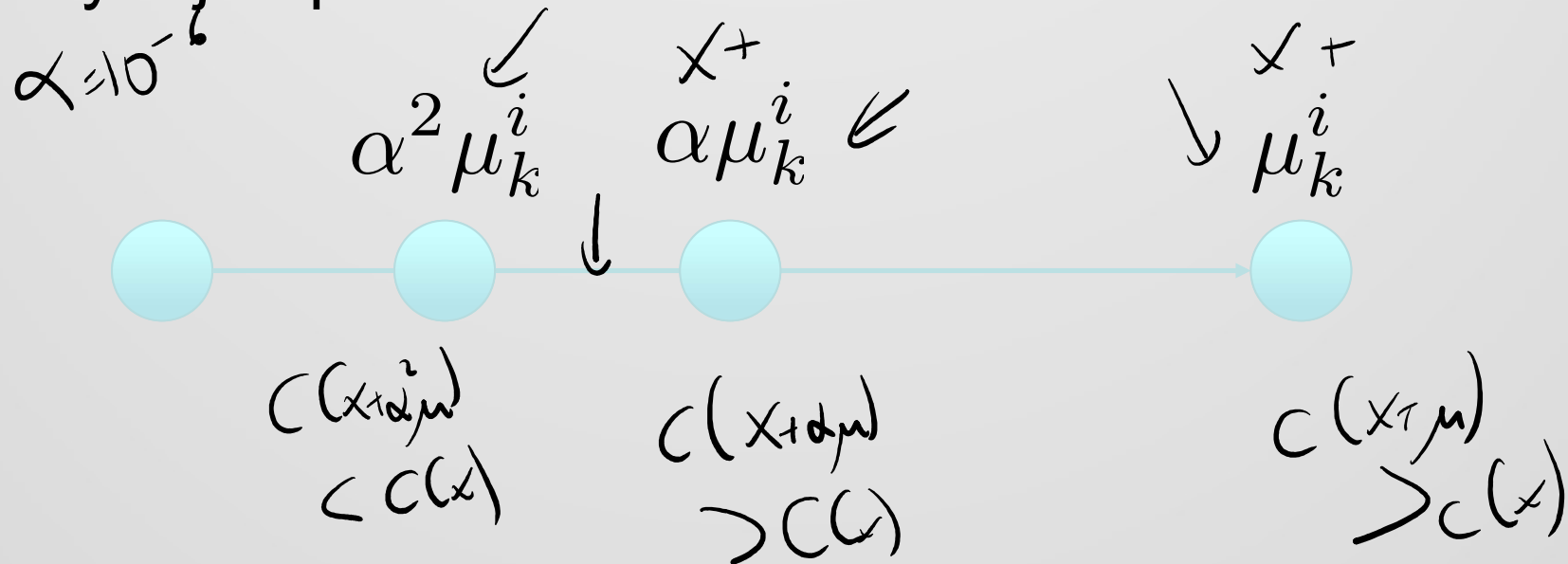
- The optimal solution is *exact line search*
- The step size is given by the solution

$$\arg \min_{\mu_k^i} c(\mathbf{X}_k^i - \mu_k^i \mathbf{J})$$

- However, the actual cost function is completely arbitrary

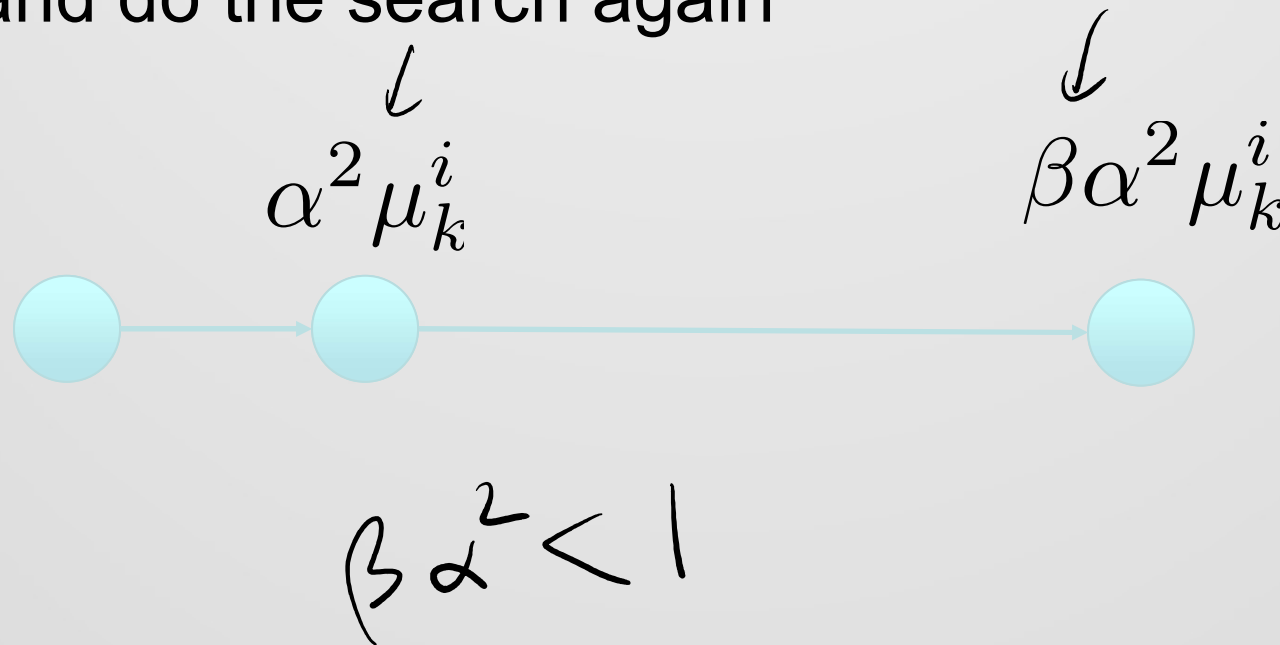
## Choosing the Step Size

- Most versions use some kind of *backtracking line search*
- Try a jump and scale back on it if it's bad



# Choosing the Step Size

- In the next step, scale the step size up by a factor  $\beta$  and do the search again

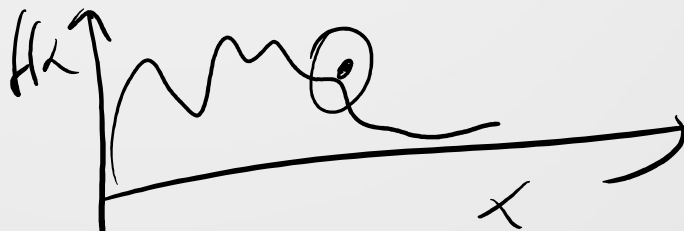


# Newton-Raphson

- Gradient descent can find the local minimum
- However, its convergence can be very slow
- There are lots of variants of nonlinear optimization algorithms which have been derived
- Here we'll look at Newton, Gauss-Newton and Levenberg-Marquardt



# Newton's Method



- The idea is to approximate the function by a quadratic and pick the step to optimize it
- The quadratic approximation is obtained from Taylor's series,

$$f(\underset{\uparrow}{x} + \underset{\uparrow}{\delta x}) \approx f(x) + f'(x)\delta x + \frac{1}{2}f''(x)\delta x^2$$

$\downarrow$                        $\downarrow$                        $\downarrow$

## Newton's Method

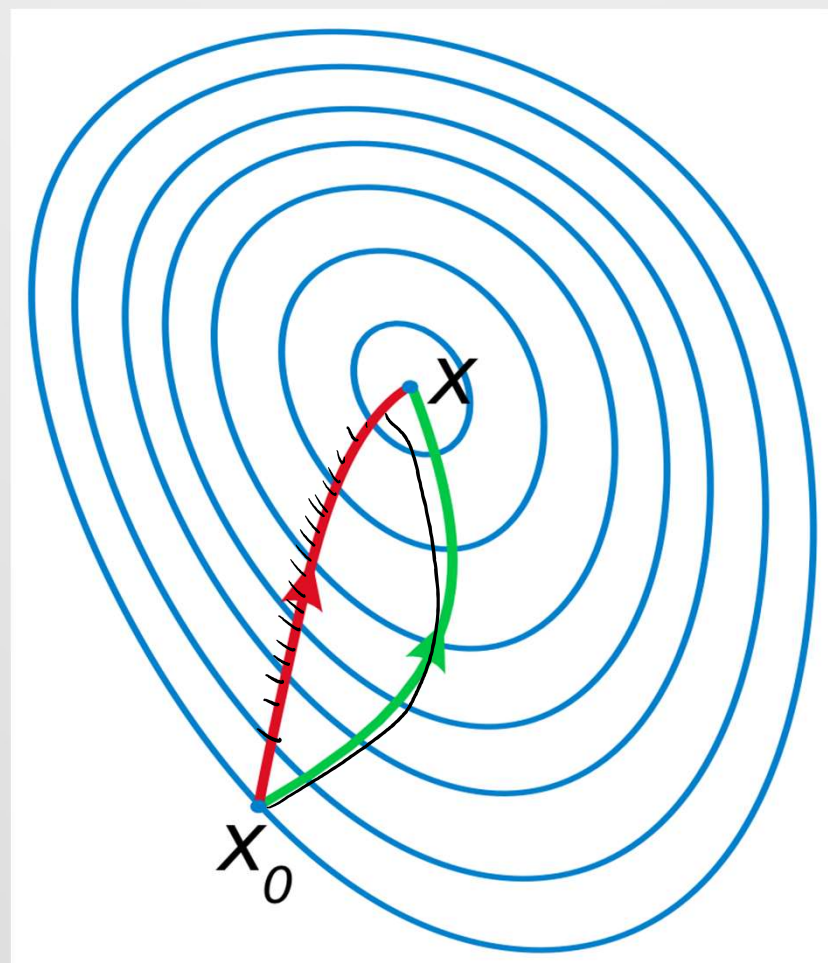
- Taking derivatives with respect to the step size,

$$\frac{d f(x + \delta x)}{d \delta x} \stackrel{=0}{=} \underline{f'(x)} + \underline{f''(x)\delta x}$$

- Since the left hand side is 0,

$$\delta x = - \frac{f'(x) \Leftarrow}{f''(x) \Leftarrow}$$

# Trajectory Taken by Newton's Method





# Newton's Method

$C = \text{scalar}$   
 $f$

$\frac{f'(x)}{f(x)}$   
 $\text{vector}$

- The multi-dimensional version is

$$\mathbf{X}_k^{i+1} = \mathbf{X}_k^i - \gamma \left[ \mathbf{H}_c(\mathbf{X}_k^i) \right]^{-1} \nabla c(\mathbf{X}_k^i)$$

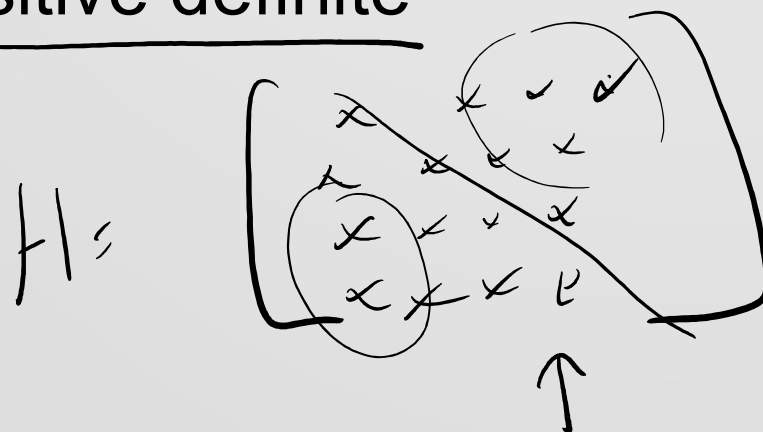
where

$$\mathbf{H}_c(\mathbf{X}_k^i) = \begin{bmatrix} \frac{\partial^2 c(\mathbf{X}_k^i)}{\partial x_1^2} & \dots & \frac{\partial^2 c(\mathbf{X}_k^i)}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 c(\mathbf{X}_k^i)}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 c(\mathbf{X}_k^i)}{\partial x_n^2} \end{bmatrix}^T$$

$\text{matrix}$

## Issues with Newton's Method

- It requires the calculation of the Hessian which can be *very expensive* to compute
- The method blows up if the Hessian is non positive definite



$$\text{eig}(H) = \lambda_1, \dots, \lambda_n$$

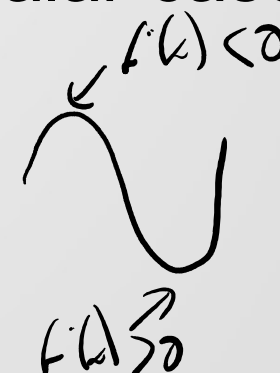
$$x_i > 0 \quad i = 1, \dots, n$$

$$\underline{x_i <}$$

# Non-Positiveness with the Hessian

- The easiest way to see this is with the scalar case
- We computed the turning point value

$$\delta x = - \frac{f'(x)}{f''(x)}$$



- However, the type of turning point depends on the second derivative

$$f''(x)$$

- If this is negative we have found a local maximum!

## Gauss-Newton

$$\nabla c = 0 \quad \epsilon \quad \epsilon^2 \quad 1/\epsilon^2 \leftarrow$$

- This uses an approximate value for the Hessian,

$$\mathbf{X}_k^{i+1} = \mathbf{X}_k^i - \gamma \left( \underbrace{\nabla c(\mathbf{X}_k^i)^\top \nabla c(\mathbf{X}_k^i)} \right)^{-1} \underbrace{\nabla c(\mathbf{X}_k^i)}$$

- The approximate value is the “square” of the Jacobian
- You don’t have to compute second derivatives
- This is guaranteed to be positive semidefinite

## Limitations with Gauss-Newton

- Gauss Newton is not guaranteed to find a local minimum
- We also have the problem that if the Jacobian is nearly singular, the step might be very inaccurately calculated, if it can be calculated at all
- Levenberg-Marquardt attempts to combine the “slow and safe” approach from Gradient Descent with “faster and less safe” ness from Gauss Newton

# Levenberg-Marquardt

- The method introduces a damping factor:

$$\mathbf{X}_k^{i+1} = \mathbf{X}_k^i$$

$$\lambda = 0 \quad \lambda \gg 0$$

$$-\gamma \left( \nabla c(\mathbf{X}_k^i)^\top \nabla c(\mathbf{X}_k^i) + \lambda \mathbf{I} \right)^{-1} \nabla c(\mathbf{X}_k^i)$$

- The damping factor is chosen dynamically, like the step size, to give the best result

# Levenberg-Marquardt Pseudo-Code

```

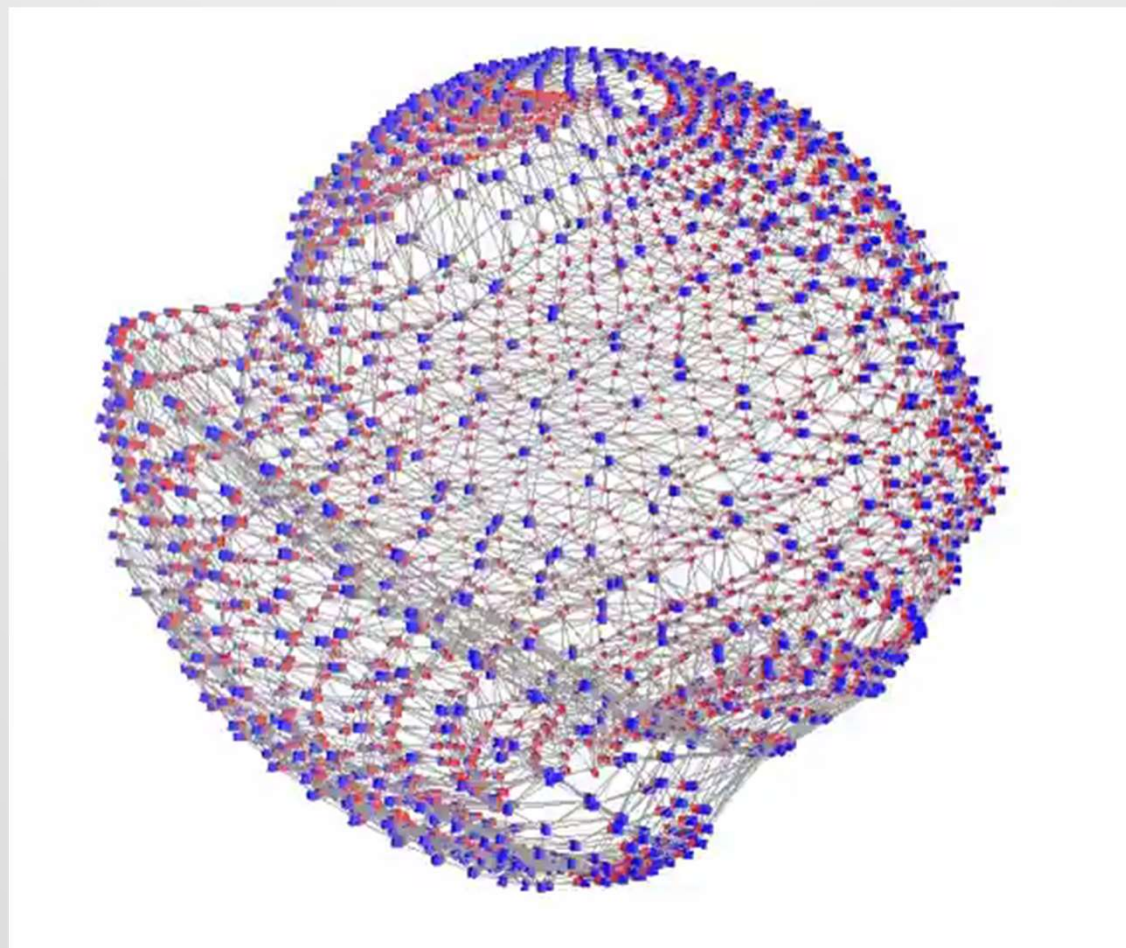
 $\lambda \leftarrow 10^{-3};$ 
Repeat
  1. Compute the gradient vector  $\mathbf{b}$  and the Hessian matrix  $\mathbf{H}$ ,
  2. Compute the increment  $\hat{\delta\theta} = -(\mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{b}$ ;
  3. If  $f(\theta + \hat{\delta\theta}) < f(\theta)$ 
    •  $\lambda \leftarrow \frac{\lambda}{10}$ ;
    Else
      • Repeat
        -  $\lambda \leftarrow \lambda \times 10$ ;
        - Compute the increment  $\hat{\delta\theta} = -(\mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{b}$ ;
      Until  $f(\theta + \hat{\delta\theta}) < f(\theta)$ ;
    Endif
  4.  $\theta \leftarrow \theta + \hat{\delta\theta}$ ;
Until convergence.
  
```

Handwritten notes:  $\nabla_C^T \nabla_C$  (pointing to Hessian matrix  $\mathbf{H}$ ),  $1/\lambda \mathbf{I}$  (pointing to  $(\mathbf{H} + \lambda \mathbf{I})^{-1}$ ), and a box around the inner repeat loop.

From ["HESSIAN MATRIX VS. GAUSS—NEWTON HESSIAN MATRIX", P. Chen, 2011](#)

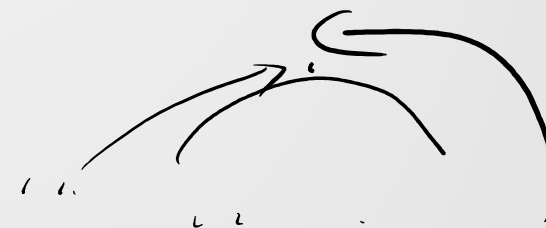
# Challenge Optimization Problem

TORO

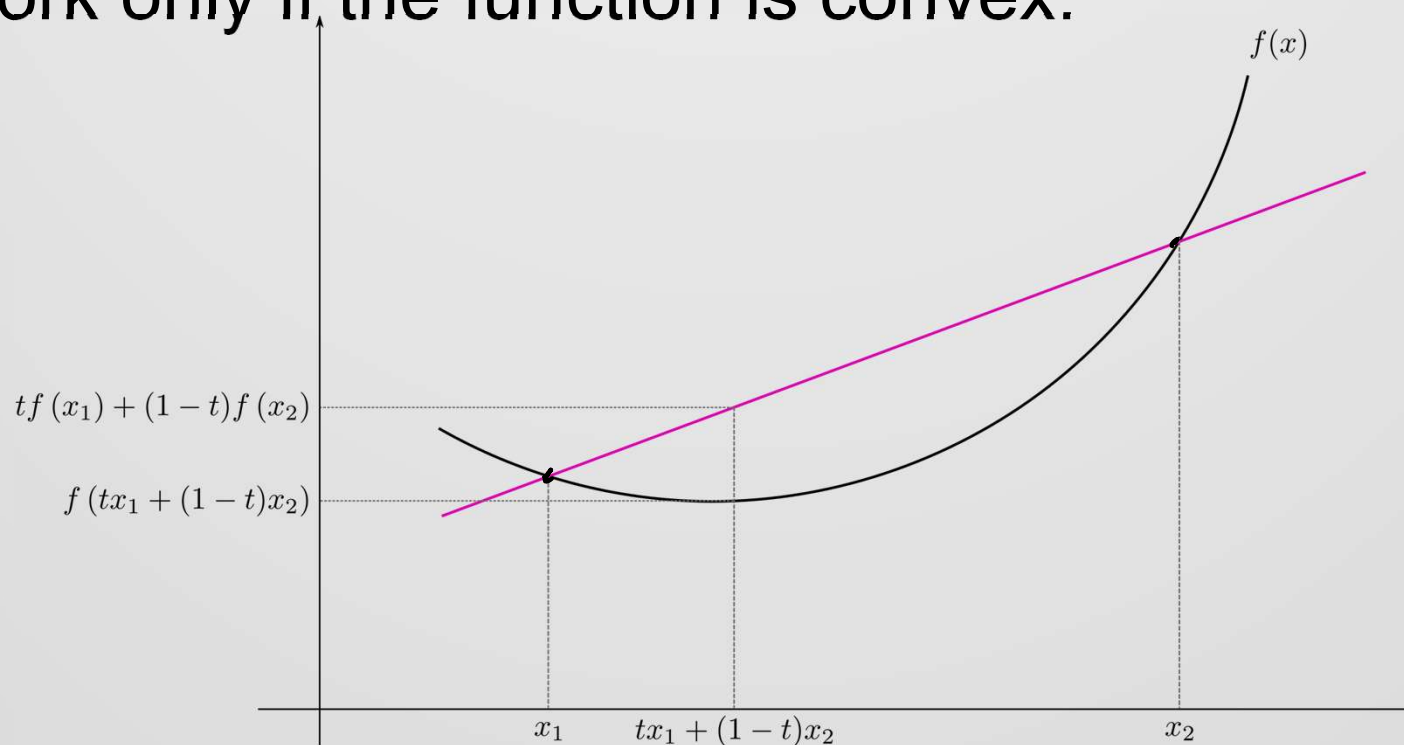




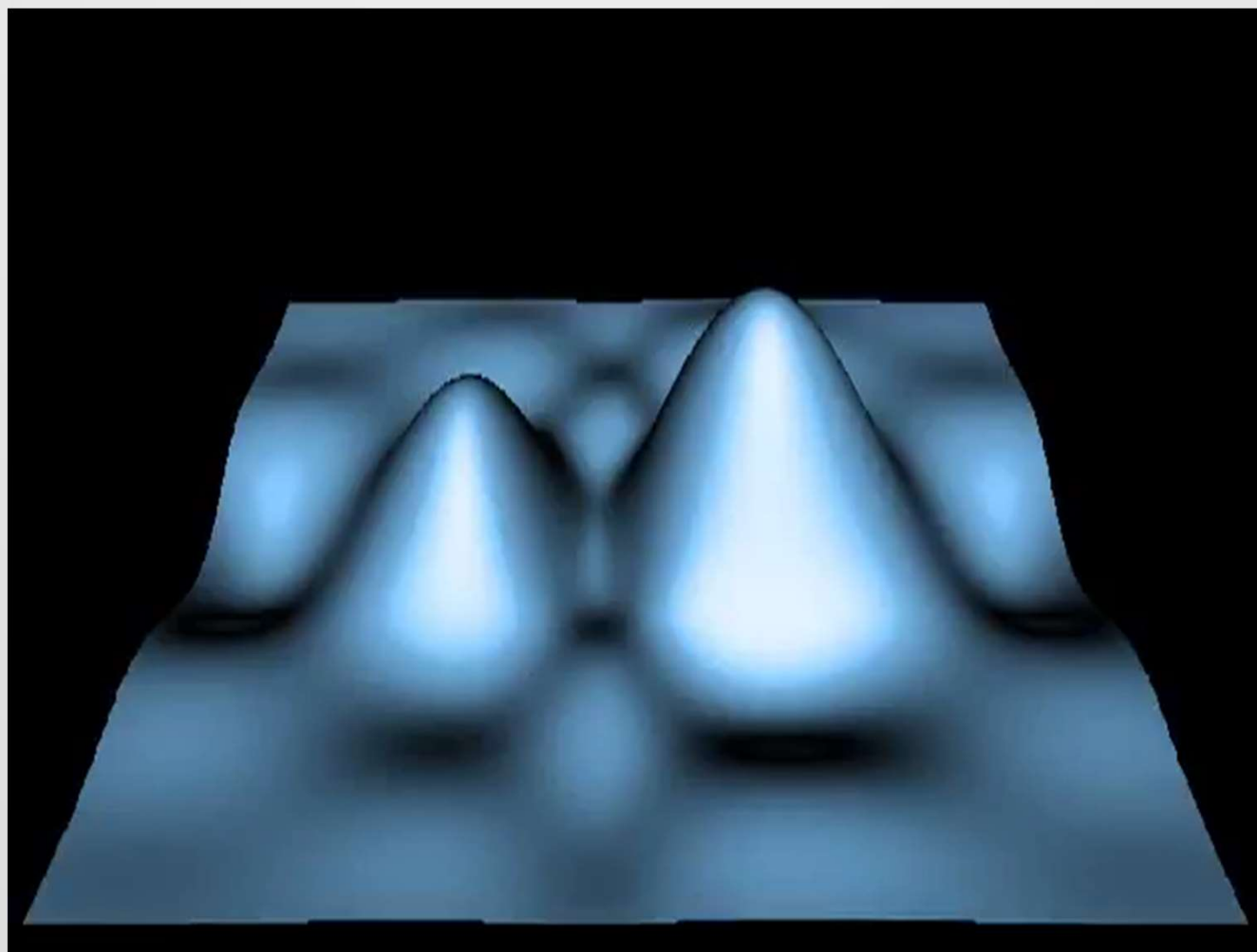
# Convexity and Optimization



- Gradient descent approaches are guaranteed to work only if the function is convex:



# Gradient Ascent and non-Convex Functions



## Addressing Issues of Non-Convexity

- Almost any real system is non-convex
- However, they have a “basin of attraction” for the solution
- The rule-of-thumb is to try to compute an initial solution which is as close as possible to the optimum
- There are new “certifiable” SLAM algorithms, but the maths is rather challenging

# Approximate Covariance Calculations

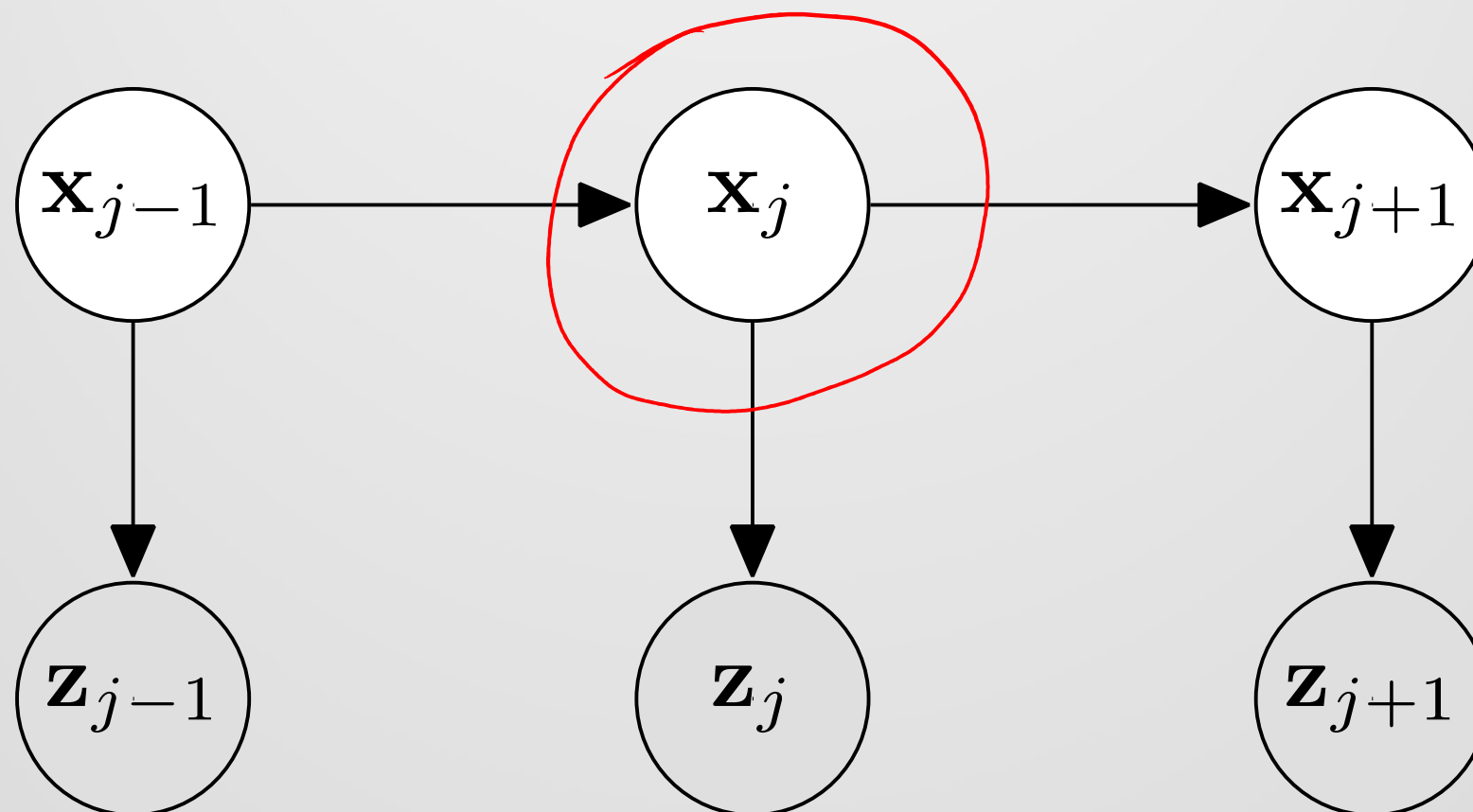
- *Motivation*
- *Interpreting the Graph*
- *Maximum Likelihood Estimation*
- *Optimization Algorithms*
- Approximate Covariance Calculations

## Computing the Covariance

$$E[(x_i - \hat{x}_i)^2]$$

- The covariance of a state estimate is *never used* when computing the MAP estimate
- However, in many cases we would like to estimate the covariance of the state estimate to determine things such as whether we have collected enough measurements, or to perform data association
- Computing it properly is challenging because we'd like to compute the expectation of a squared function

# Computing the Covariance

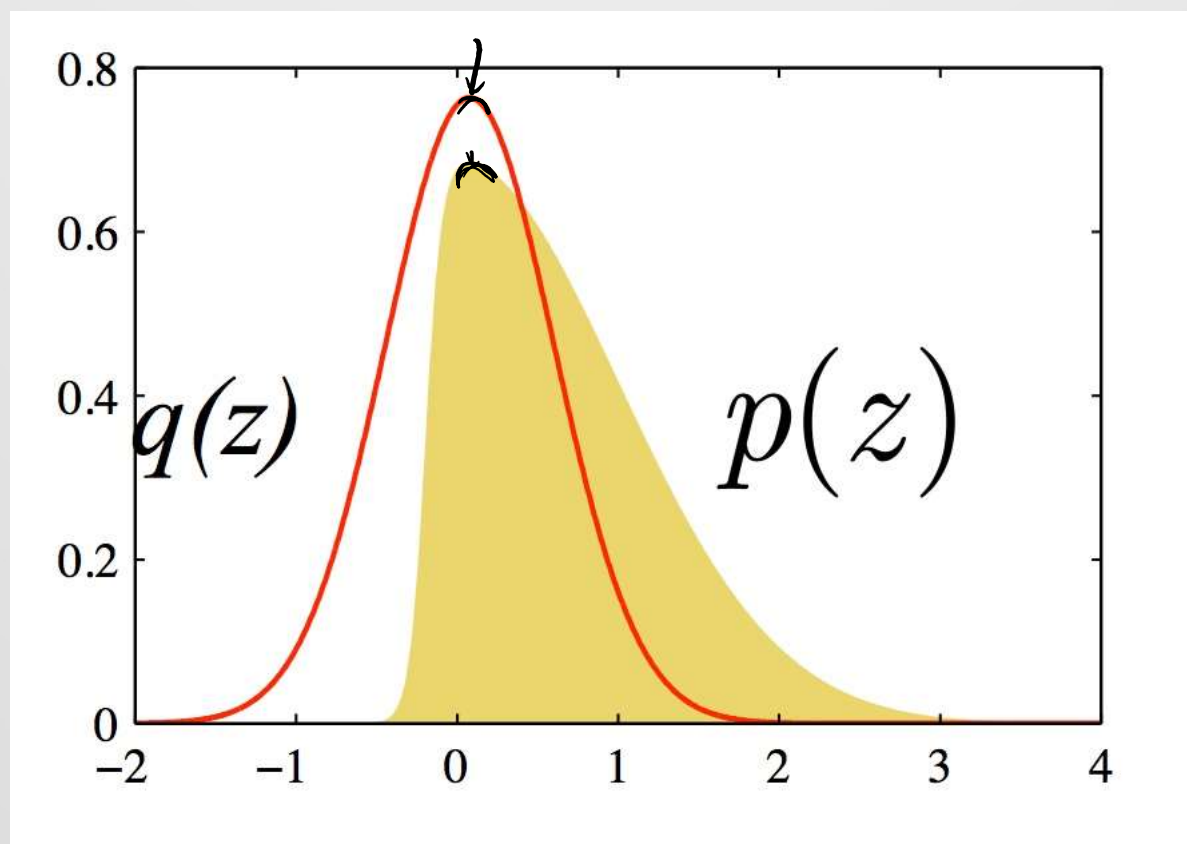


# Laplace Approximation



- The distribution around the MAP estimate is approximated by a Gaussian
- The mean of the Gaussian is the MAP estimate
- The covariance of the Gaussian is given by the local curvature of the distribution

# The Laplace Approximation





## The Laplace Approximation in a Graph

- For the graph, it turns out the curvature is actually given by the inverse of the Hessian
- Using various approximations, its value is given by

$$\rightarrow \left( \nabla c(\mathbf{X}_k^*)^\top \nabla c(\mathbf{X}_k^*) \right)^{-1} \quad \frac{\lambda \mathbf{I}}{\uparrow}$$

*marginals*

*( $\hat{x}, \rho$ )*

- This value was actually worked out as part of the Levenberg-Marquardt update step

# Summary

- To extract a point estimate from the graph, we use maximum a posterior (MAP) estimates
- These do not require calculation of the normalization constant and tend to put the estimate in the <sup>sensible</sup> ~~correct~~ place
- Computing the MAP estimate is a weighted nonlinear least squares estimation problem
- Recursive solutions exist to compute it

↑ Bundle adjustment ←