ENGR-E 511; ENGR-E 399

# "Machine Learning for Signal Processing"

Module 01: Lecture 03:

# Optimization

## Minje Kim

Department of Intelligent Systems Engineering

Email: minje@indiana.edu

Website: http://minjekim.com

Research Group: http://saige.sice.indiana.edu

Meeting Request: http://doodle.com/minje
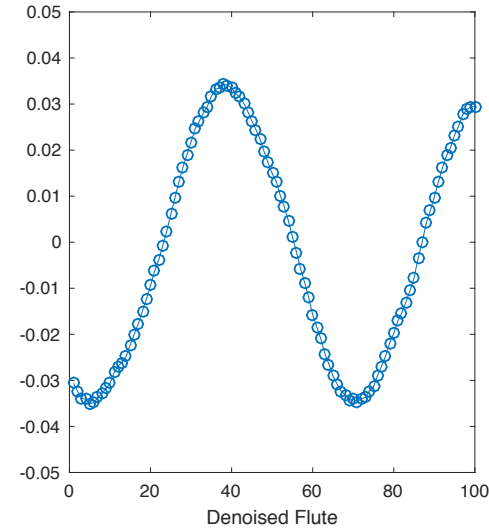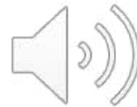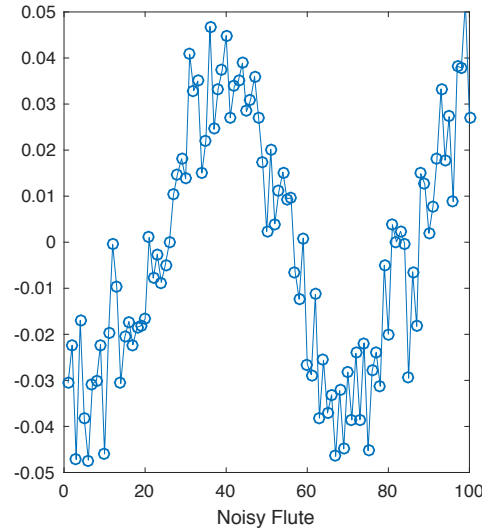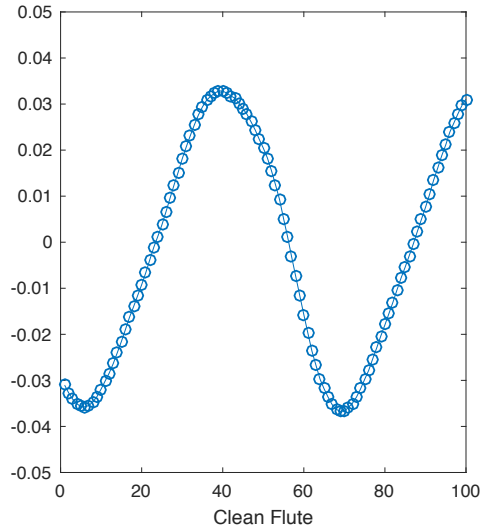
INDIANA UNIVERSITY
## SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Audio Denoising

## - Least mean square

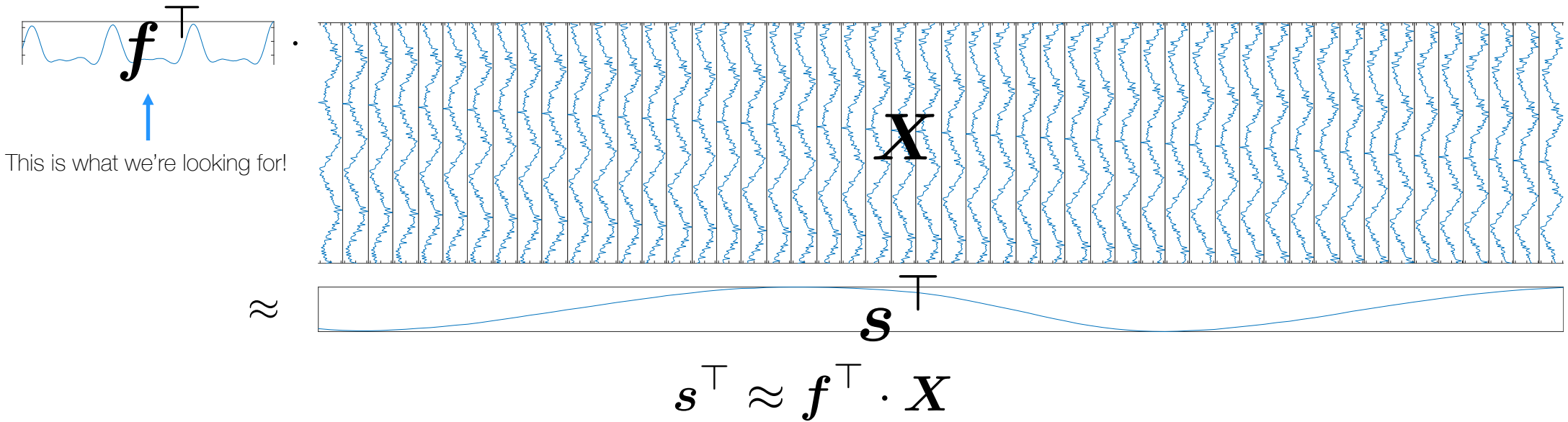○  We want to find out a filter that denoises the noisy input

# Audio Denoising

## - Least mean square

○ Let's assume that there is a filter that does this job

$$s[n] \approx f * x[n] = \sum_{\tau=0}^{T-1} f[\tau]x[n-\tau]$$

○ In other words



This is what we're looking for!

$$s^\top \approx f^\top \cdot X$$

# Audio Denoising

## - Least mean square

○ What we find: the filter that minimizes

  □ the error between the source and the reconstruction

$$\arg \min_{\boldsymbol{f}} \mathcal{E}\left(\boldsymbol{s}^\top \middle|| \boldsymbol{f}^\top \boldsymbol{X}\right)$$

○ Let's use the **mean squared error**

$$\mathcal{E}\left(\boldsymbol{s}^\top \middle|| \boldsymbol{f}^\top \boldsymbol{X}\right) = \frac{1}{N}(\boldsymbol{s}^\top - \boldsymbol{f}^\top \boldsymbol{X}) \cdot (\boldsymbol{s}^\top - \boldsymbol{f}^\top \boldsymbol{X})^\top = \frac{1}{N}(\boldsymbol{s}^\top - \hat{\boldsymbol{s}}^\top) \cdot (\boldsymbol{s} - \hat{\boldsymbol{s}}) \qquad (\text{where } \hat{\boldsymbol{s}}^\top = \boldsymbol{f}^\top \boldsymbol{X})$$

$$= \frac{1}{N}\boldsymbol{d}^\top \boldsymbol{d} \;=\; \frac{1}{N}\sum_n d_n^2 \;=\; \frac{1}{N}\sum_n (s_n - \hat{s}_n)^2 \qquad (\text{where } \boldsymbol{d} = \boldsymbol{s} - \hat{\boldsymbol{s}})$$

$$= \frac{1}{N}\sum_n \left(s_n - \sum_\tau f_\tau X_{\tau n}\right)^2$$

$$= \frac{1}{N}\sum_n \left(s_n - \sum_{\tau \neq \tau'} f_\tau X_{\tau n} - f_{\tau'} X_{\tau' n}\right)^2 = \frac{1}{N}\sum_n (C_n - f_{\tau'} X_{\tau' n})^2 \qquad (\text{where } C_n = s_n - \sum_{\tau \neq \tau'} f_\tau X_{\tau n})$$

$$= \frac{1}{N}\sum_n X_{\tau' n}^2 f_{\tau'}^2 - 2 C_n f_{\tau'} + C_n^2 = \frac{1}{N}\left(\sum_n X_{\tau' n}^2\right) f_{\tau'}^2 - 2\left(\sum_n C_n\right) f_{\tau'} + \left(\sum_n C_n^2\right)$$

○ **The mean squared error is a quadratic function (of one of your filter values)!**

# Audio Denoising

## - Least mean square

○ Some linear algebra

$$\frac{\partial \boldsymbol{a}^\top \boldsymbol{x}}{\partial \boldsymbol{x}} = \boldsymbol{a} \qquad \frac{\partial \sum_i a_i x_i}{\partial x_i} = a_i \qquad \frac{\partial \boldsymbol{x}^\top \boldsymbol{a}}{\partial \boldsymbol{x}} = \boldsymbol{a}$$

$$\frac{\partial \boldsymbol{x}^\top \boldsymbol{A}\boldsymbol{A}^\top \boldsymbol{x}}{\partial \boldsymbol{x}} = \boldsymbol{A}\boldsymbol{A}^\top \boldsymbol{x} + \boldsymbol{A}\boldsymbol{A}^\top \boldsymbol{x}$$

○ Partial derivatives

$$\mathcal{E}(\boldsymbol{s}^\top || \boldsymbol{f}^\top \boldsymbol{X}) = \frac{1}{N}(\boldsymbol{s}^\top - \boldsymbol{f}^\top \boldsymbol{X}) \cdot (\boldsymbol{s}^\top - \boldsymbol{f}^\top \boldsymbol{X})^\top$$

$$\frac{\partial \mathcal{E}}{\partial \boldsymbol{f}} = \frac{\partial (\boldsymbol{s}^\top \boldsymbol{s} - \boldsymbol{f}^\top \boldsymbol{X}\boldsymbol{s} - \boldsymbol{s}^\top \boldsymbol{X}^\top \boldsymbol{f} + \boldsymbol{f}^\top \boldsymbol{X}\boldsymbol{X}^\top \boldsymbol{f})}{\partial \boldsymbol{f}}$$

$$= \frac{\partial (-2\boldsymbol{f}^\top \boldsymbol{X}\boldsymbol{s} + \boldsymbol{f}^\top \boldsymbol{X}\boldsymbol{X}^\top \boldsymbol{f})}{\partial \boldsymbol{f}}$$

$$= -2\boldsymbol{X}\boldsymbol{s} + 2\boldsymbol{X}\boldsymbol{X}^\top \boldsymbol{f} \quad = 0$$

○ Therefore: $\quad \underset{\boldsymbol{f}}{\arg\min} \, \mathcal{E}(\boldsymbol{s}^\top || \boldsymbol{f}^\top \boldsymbol{X}) = (\boldsymbol{X}\boldsymbol{X}^\top)^{-1}\boldsymbol{X}\boldsymbol{s}$

# Audio Denoising

## - Least mean square

○ Now what?

    □ We learned a filter that can denoise a particular kind of signal

        • i.e. A flute note corrupted by Gaussian noise

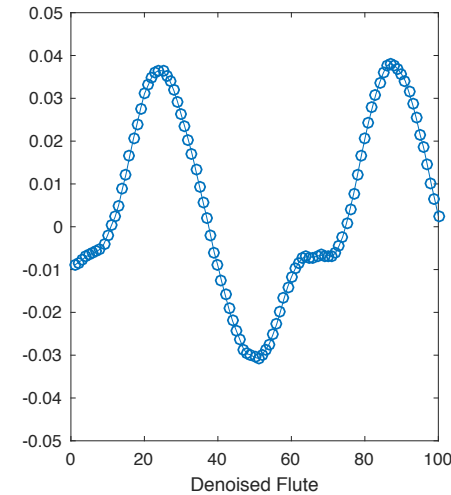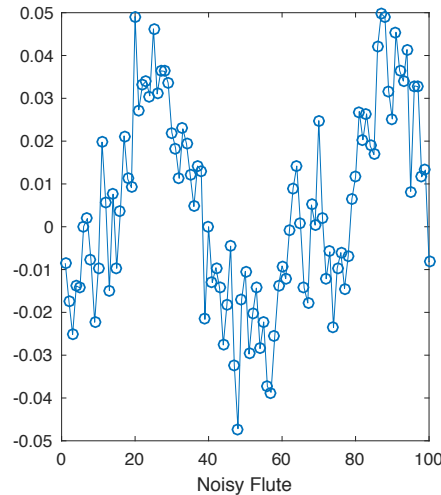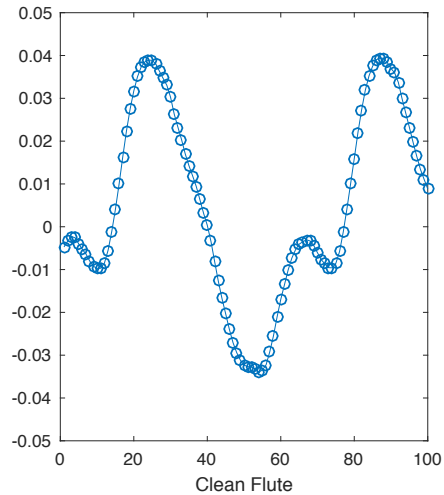    □ Testing on another flute note contaminated with Gaussian noise



Clean Flute

Noisy Flute

Denoised Flute

# Image Denoising

## - Gradient descent

○ 2D denoising filter

□ Images are in gray scale

- Black: 0
- White: 1

○ If you flatten the filter and patches, the procedure boils down to
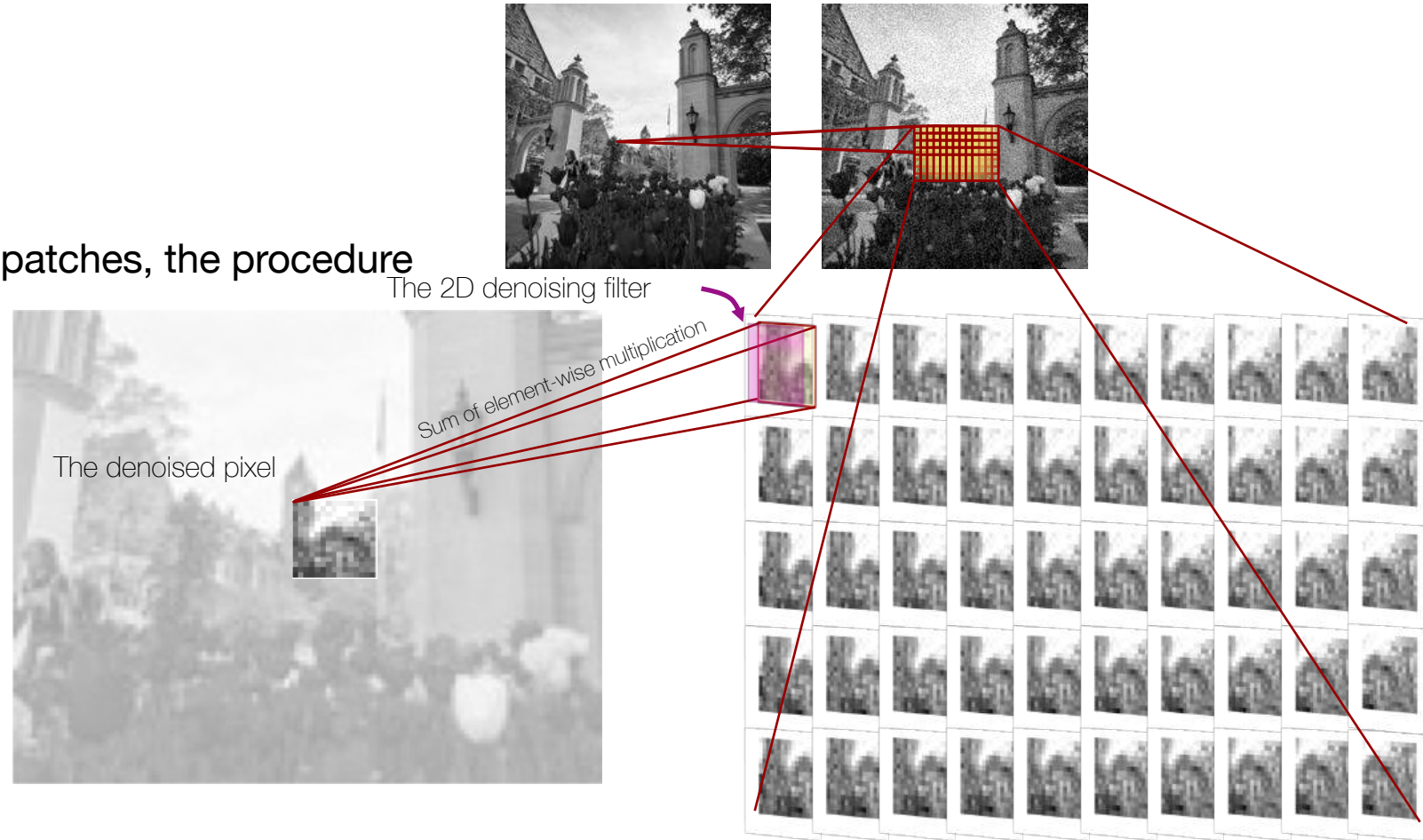
□ An inner product

The 2D denoising filter

Sum of element-wise multiplication

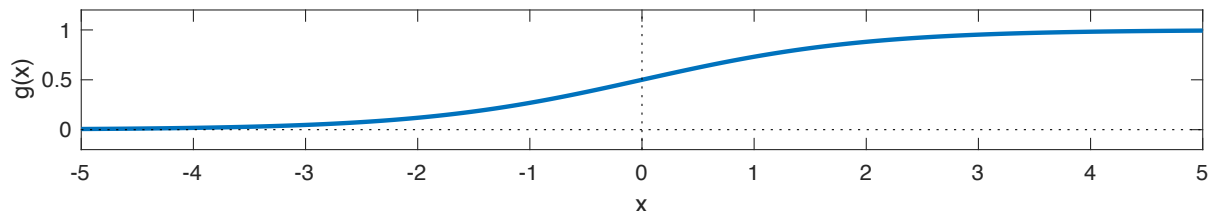The denoised pixel

# Image Denoising

## - Gradient descent

○ What I don't like about the error function I used before:

$$\mathcal{E}(s^\top \| f^\top X) = \frac{1}{N}(s^\top - f^\top X) \cdot (s^\top - f^\top X)^\top$$

 □ The reconstructed images pixels should be neither negative nor bigger than 1
 □ No way to control this in the original error function

○ Logistic function can take care of it

$$g(x) = \frac{1}{1 + e^{-x}}$$



 □ On a vector input:

$$g(\boldsymbol{x}) = \begin{bmatrix} g(x_1) \\ g(x_2) \\ \vdots \\ g(x_D) \end{bmatrix}$$

○ The new error function:  $\mathcal{E}(s^\top \| g(f^\top X)) = \frac{1}{N}(s^\top - g(f^\top X)) \cdot (s^\top - g(f^\top X))^\top$
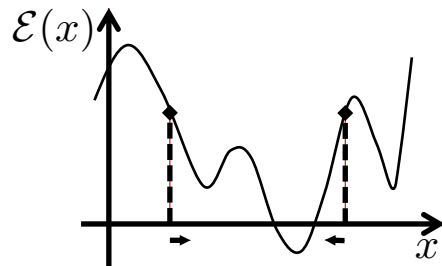
# Image Denoising

## - Gradient descent

- ○ Is this a quadratic function of $f$ ?

$$\mathcal{E}(s^\top \| g(f^\top X)) = \frac{1}{N}(s^\top - g(f^\top X)) \cdot (s^\top - g(f^\top X))^\top = \frac{1}{N}\left(s^\top - \frac{1}{1+\exp(-f^\top X)}\right) \cdot \left(s^\top - \frac{1}{1+\exp(-f^\top X)}\right)^\top$$

  - □ Maybe not. I have no idea.

- ○ Now we need a different strategy: **gradient descent**



- ○ Partial differentiation $\frac{\partial \mathcal{E}}{\partial f} = -\frac{2}{N}X\left\{(s^\top - g(f^\top X)) \odot g'(f^\top X)\right\}^\top$

- ○ Minimization
  - □ Negative gradient: $\nabla f^{(i)} = -\frac{\partial \mathcal{E}}{\partial f^{(i)}}$
  - □ Update rule:
    $$f^{(i+1)} \leftarrow f^{(i)} + \rho \nabla f^{(i)}$$

# Image Denoising

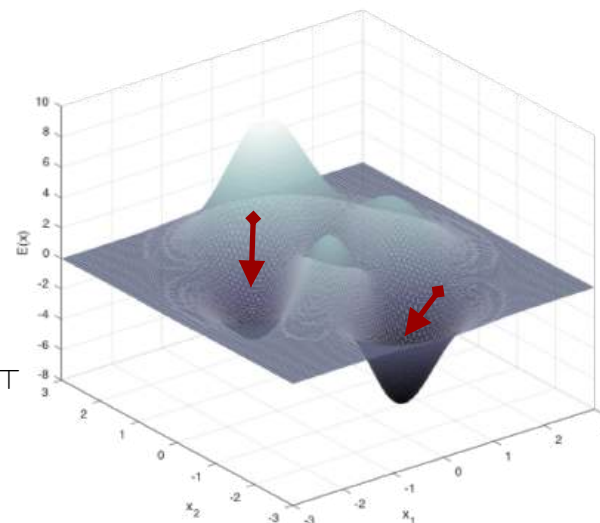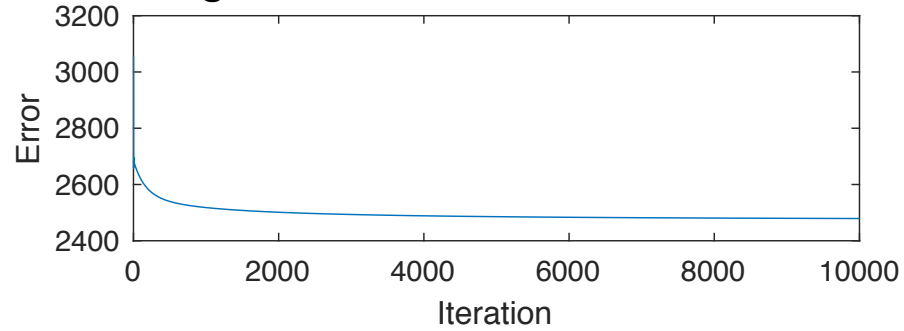## - Gradient descent

○ Convergence



○ Learned filter
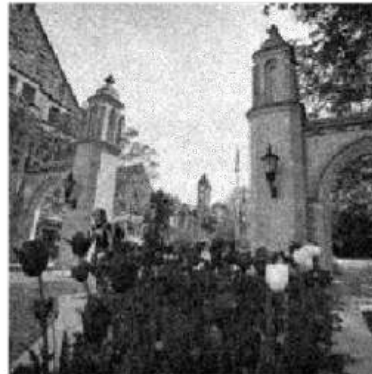


○ Denoising results (training)
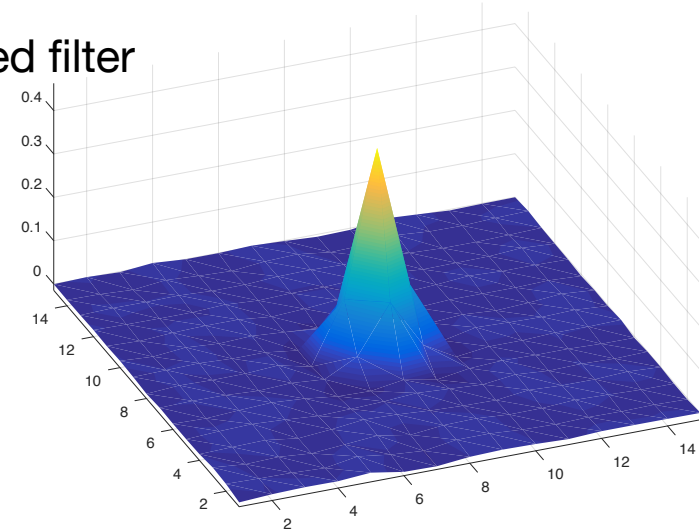


Original

Noisy

Denoised
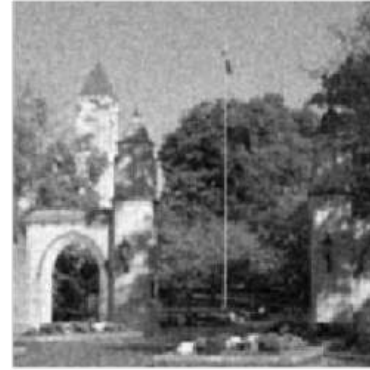
# Image Denoising

## - Gradient descent

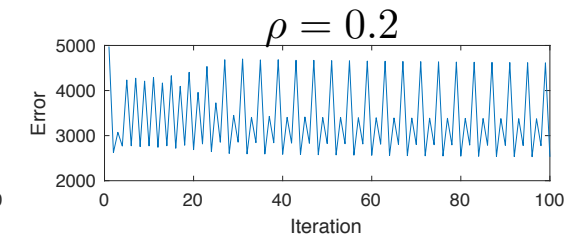○ Denoising results (test)



Original          Noisy          Denoised
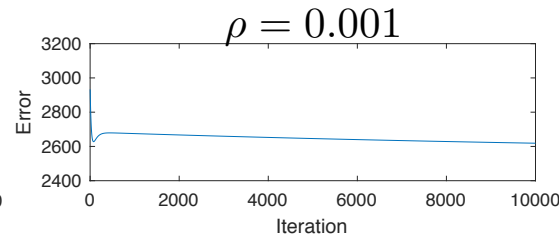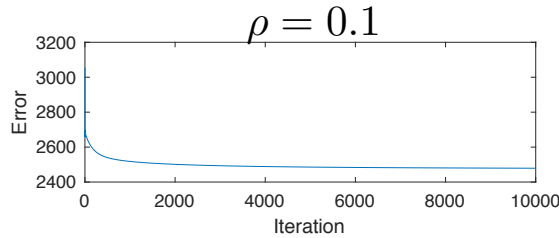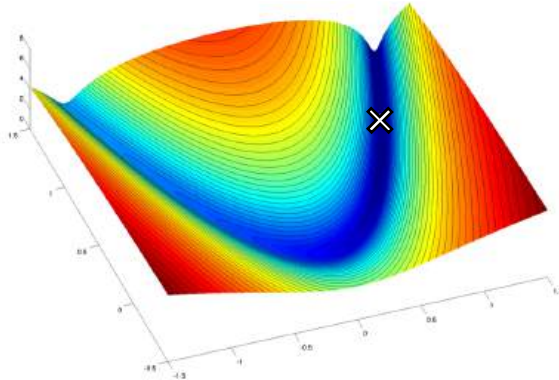
# Image Denoising

## - Gradient descent

○ Limitations of the gradient descent method $f^{(i+1)} \leftarrow f^{(i)} + \rho \nabla f^{(i)}$
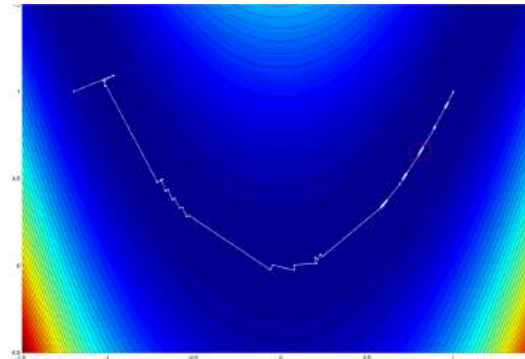
    □ Step size



    □ Zigzag movement



Rosenbrock function         Trajectory of x         Zigzag pattern

# Image Denoising

## - Newton's method

○ Gradient descent is based on an imperfect approximation

    □ First-order Taylor series expansion
$$f(\theta) \approx f(\theta^*) + f'(\theta^*)(\theta - \theta^*)$$

    □ Not accurate enough for non-linear objective functions

    □ Why?

       • The gradient is changing its slope

    □ Second-order expansion
$$f(\theta) \approx f(\theta^*) + f'(\theta^*)(\theta - \theta^*) + \frac{1}{2}f''(\theta^*)(\theta - \theta^*)^2$$

$$f(\theta) \approx f(\theta^*) + \left( f'(\theta^*)(\theta - \frac{1}{2}\theta^* f'(\theta^*)(\theta - \theta^*) \right)(\theta - \theta^*)$$

$$\left( f'(\theta^*) + \frac{1}{2}f''(\theta^*)(\theta - \theta^*) \right)$$

$$\left( f'(\theta^*) + f''(\theta^*)(\theta - \theta^*) \right)$$

$f(\theta)$

$f(\theta^*) + f'(\theta^*)(\theta - \theta^*)$

$f'(\theta^*)$

$f'(\theta^*)(\theta - \theta^*)$

$f(\theta^*)$

$\theta^*$    $\theta$

# Image Denoising

## - Newton's method

○ The new update rule with Newton's method

 □ Let's derive the update rule from the Taylor's series
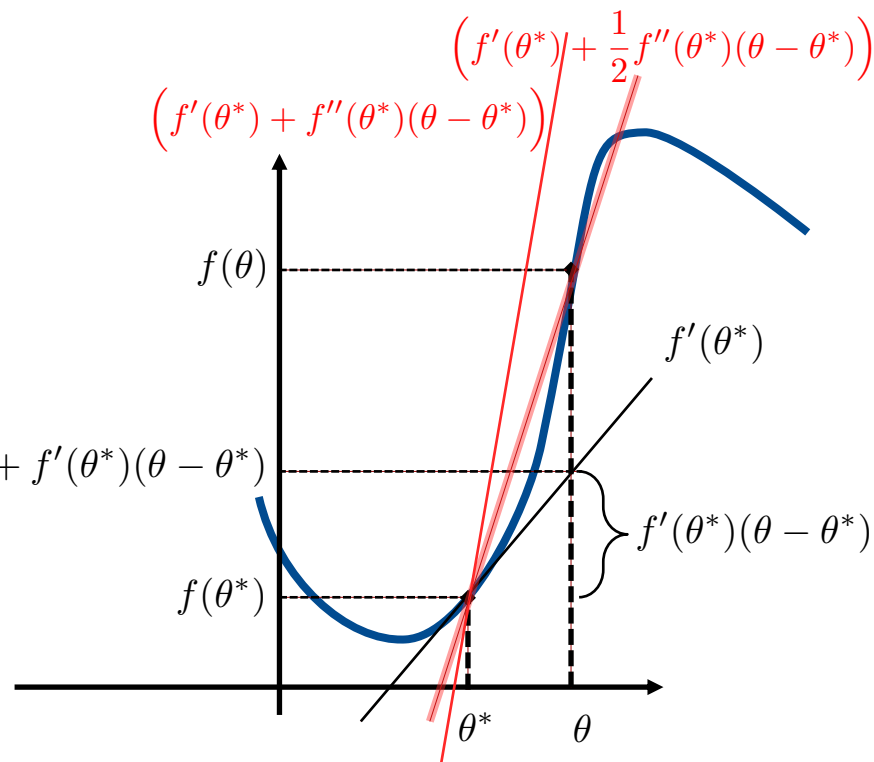
$$f(\theta^{(i+1)}) \approx \tilde{f}(\theta^{(i+1)}) = f(\theta^{(i)}) + f'(\theta^{(i)})\Delta\theta + \frac{1}{2}f''(\theta^{(i)})\Delta\theta^2$$

 □ We know that $\tilde{f}(\theta^{(i+1)})$ is a quadratic function of $\Delta\theta$

 • So, the stationary point is the local minimum

$$\frac{\partial\tilde{f}(\theta^{(i+1)})}{\partial\theta} = f'(\theta^{(i)}) + f''(\theta^{(i)})\Delta\theta = 0 \qquad \theta^{(i+1)} \leftarrow \theta^{(i)} + \Delta\theta = \theta^{(i)} - \frac{f'(\theta^{(i)})}{f''(\theta^{(i)})}$$

 □ No step size (not always)

 □ Finds the global minimum if the objective function is quadratic

○ Multidimensional case

 □ Gradient descent: $\quad \boldsymbol{f}^{(i+1)} \leftarrow \boldsymbol{f}^{(i)} + \rho\nabla\boldsymbol{f}^{(i)} \qquad \nabla\boldsymbol{f}^{(i)} = -\dfrac{\partial\mathcal{E}}{\partial\boldsymbol{f}^{(i)}}$

 □ Newton's method: $\quad \boldsymbol{f}^{(i+1)} \leftarrow \boldsymbol{f}^{(i)} - (\boldsymbol{H}^{(i)})^{-1}\nabla\boldsymbol{f}^{(i)} \qquad H_{mn}^{(i)} = \dfrac{\partial^2\mathcal{E}}{\partial f_m^{(i)}\partial f_n^{(i)}}$

 • Problematic if there are too many parameters

# Image Denoising

## - Newton's method

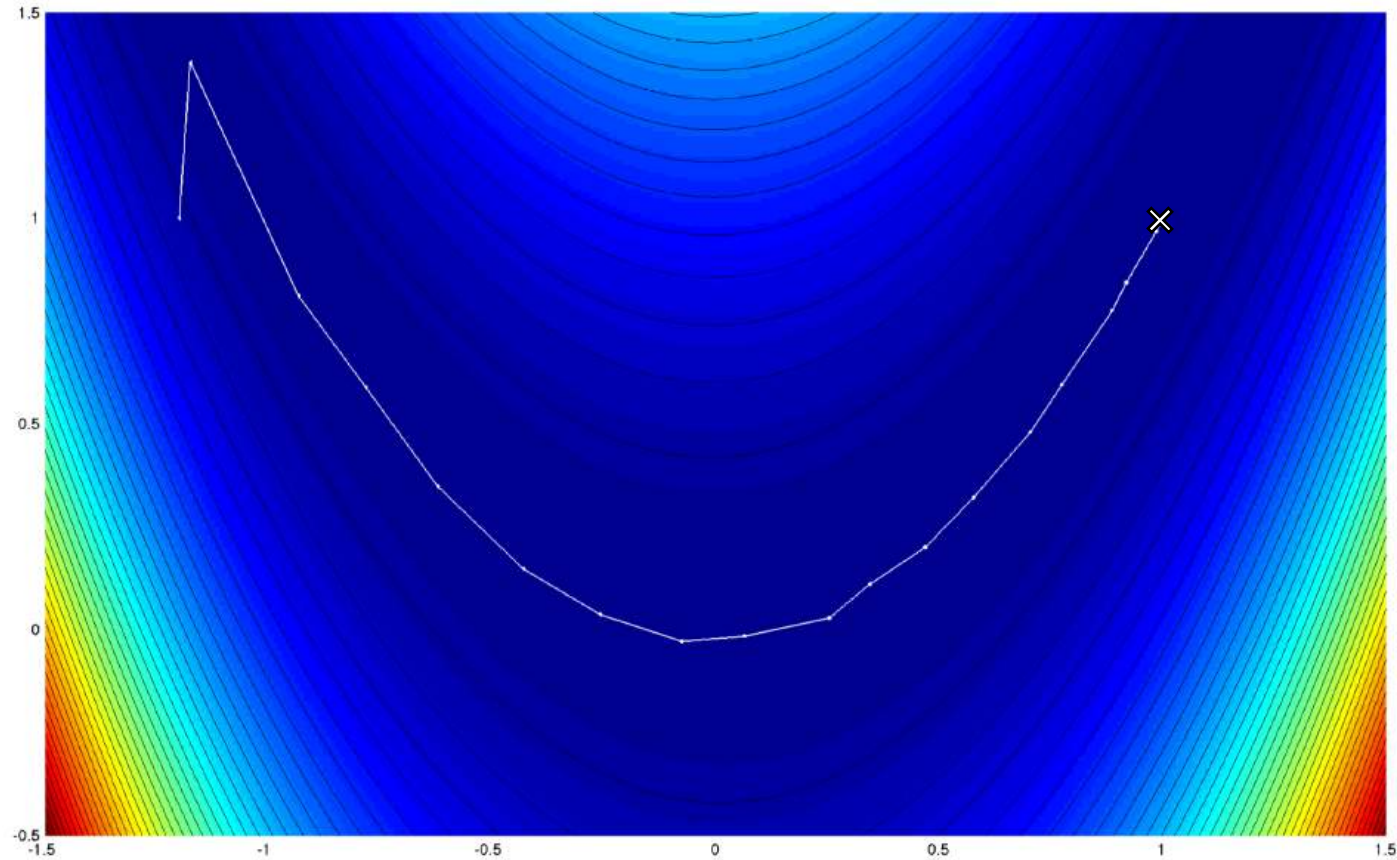○ Newton's method for optimizing the Rosenbrock function

# Image Denoising

## - Newton's method

○ You don't want to calculate the inverse-Hessian matrix at every iteration

  □ Especially if you have too many parameters

  □ Inversion is expensive

  □ There are a lot of approximation techniques

○ Quasi Newton's methods

  □ BFGS calculates an approximation $\quad G^{(i)} \approx \left( H^{(i)} \right)^{-1}$

$$G^{(i+1)} \leftarrow G^{(i)} + \frac{pp^\top}{p^\top v} - \frac{(G^{(i)}v)v^\top G^{(i)}}{v^\top G^{(i)} v} + (v^\top G^{(i)} v)uu^\top$$

$$p = f^{(i+1)} - f^{(i)}$$

$$v = \nabla f^{(i+1)} - \nabla f^{(i)}$$

$$u = \frac{p}{p^\top v} - \frac{G^{(i)}v}{v^\top G^{(i)} v}$$
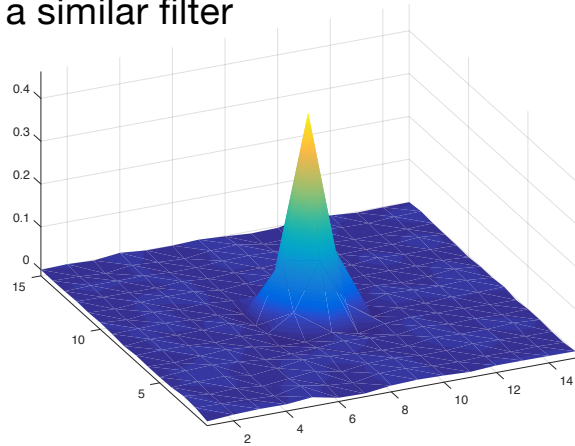
  □ We use gradients and parameters to approximate Hessian
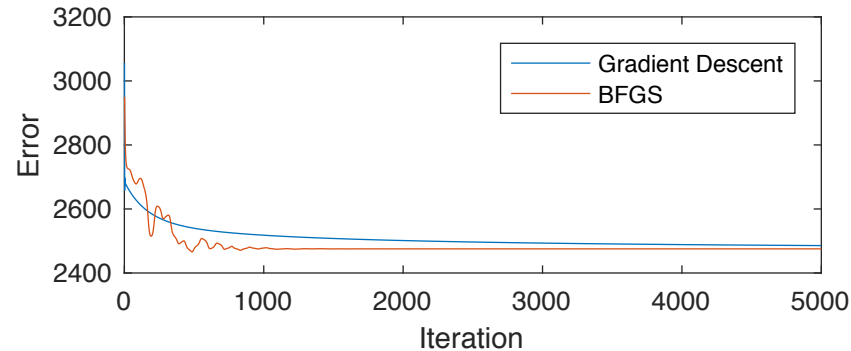
# Image Denoising

## - Newton's method

○ BFGS results

  □ Similar denoising performance

  □ Learns a similar filter



Original · Noisy · Denoised

○ Convergence

# Constrained Optimization

## - Maximum likelihood estimation with Lagrange multiplier

○ You're analyzing Prof. K's writing to build a model out of it

○ Somebody claims that a newly discovered document looks like written by Prof. K

　□ You'll use the model to determine whether the spurious article is a forgery or not

○ You collect many articles from Prof. K's journal as a training set and count the number of three keywords:

　□ "Research": 19 times, "Burger": 15 times, "Spinach": 7 times

○ Now you find the multinomial distribution that best fits the data

　□ Multinomial distribution

$$P(X_1 = x_1, X_2 = x_2, \cdots, X_K = x_K; N, p_1, p_2, \cdots, p_K) = \frac{N!}{x_1! x_2! \cdots x_K!} p_1^{x_1} p_2^{x_2} \cdots p_K^{x_K}$$

$$0 \le p_k \le 1 \qquad \sum_{k=1}^{K} p_k = 1 \qquad N = \sum_{k=1}^{K} x_k$$

# Constrained Optimization

## - Maximum likelihood estimation with Lagrange multiplier

○ The objective function:

$$\underset{p_1,p_2,\cdots,p_K}{\arg\max} \frac{N!}{x_1!x_2!\cdots x_K!}p_1^{x_1}p_2^{x_2}\cdots p_K^{x_K}$$
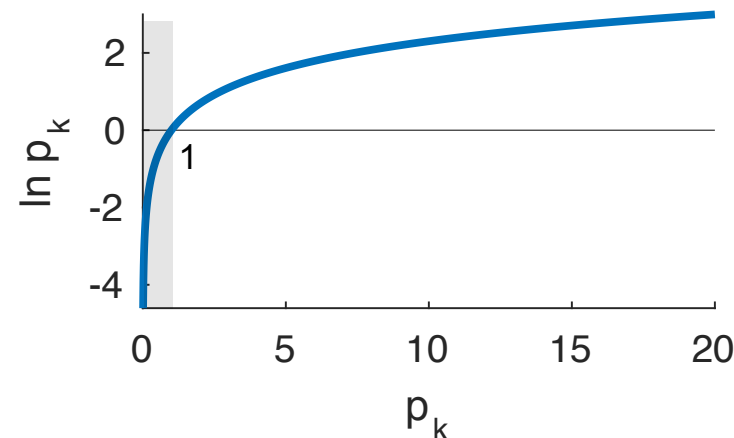
□ Take logarithm for convenience

$$\underset{p_1,p_2,\cdots,p_K}{\arg\max} \mathcal{LL}(p_1,p_2,\cdots,p_k) = \underset{p_1,p_2,\cdots,p_K}{\arg\max} \ln\left\{\frac{N!}{x_1!x_2!\cdots x_K!}p_1^{x_1}p_2^{x_2}\cdots p_K^{x_K}\right\} = \underset{p_1,p_2,\cdots,p_K}{\arg\max} \sum_{k=1}^{K}x_k\ln p_k + Constan$$

○ Along a particular dimension $p_k$ we know that the log function does not have a stationary point

○ But, we do know that $0 \leq p_k \leq 1$

□ Then, is this the solution? $p_1 = p_2 = \cdots = p_K = 1$

□ No, because we forgot another constraint: $\sum_{k=1}^{K}p_k = 1$

# Constrained Optimization
## - Maximum likelihood estimation with Lagrange multiplier

○ We saw that MLE is not enough to solve the problem
  □ Let's deal with this **equality constraint**

○ Let $\boldsymbol{p}^* = [p_1^*, p_2^*, \cdots, p_K^*]^\top$ be our solution to the optimization problem

○ And, the equality constraint $\quad g(\boldsymbol{p}) = \sum_k p_k - 1$
  □ Since $\boldsymbol{p}^*$ is the solution,
$$g(\boldsymbol{p}^*) = 0$$

○ Let's add a vector with very small values
  □ that still meet the constraint $\qquad g(\boldsymbol{p}^* + \boldsymbol{\epsilon}) = 0$

○ Also, because of the Taylor series expansion,
$$g(\boldsymbol{p}^* + \boldsymbol{\epsilon}) = 0 \approx g(\boldsymbol{p}^*) + \boldsymbol{\epsilon}^\top \nabla g(\boldsymbol{p}^*)$$

○ Therefore, $\nabla g(\boldsymbol{p}^*)$ is orthogonal to $\boldsymbol{\epsilon}$, and consequently to the constrained surface, too

○ Meanwhile, since $\boldsymbol{p}^*$ is the solution, $\mathcal{LL}(\boldsymbol{p}^*)$ is the maximum
  □ If we add $\boldsymbol{\epsilon}$, then $\boldsymbol{p}^* + \boldsymbol{\epsilon}$ will decrease the objective function value
  □ Hence, $\nabla \mathcal{LL}(\boldsymbol{p}^*)$ is orthogonal to the constrained surface, too

# Constrained Optimization

## - Maximum likelihood estimation with Lagrange multiplier

○ What we just found is that $\nabla\mathcal{LL}(\boldsymbol{p}^*)$ and $\nabla g(\boldsymbol{p}^*)$ are parallel

$$\nabla\mathcal{LL}(\boldsymbol{p}^*) = \lambda\nabla g(\boldsymbol{p}^*) \Leftrightarrow \nabla\mathcal{LL}(\boldsymbol{p}^*) + \lambda\nabla g(\boldsymbol{p}^*) = 0$$

○ Now, let's define the final objective function with Lagrange multiplier

$$\nabla\mathcal{LLL}(\boldsymbol{p}, \lambda) = \mathcal{LL}(\boldsymbol{p}) + \lambda g(\boldsymbol{p})$$

□ Therefore, $\boldsymbol{p}$ at the stationary point meets the parallel condition

□ On top of that, we still have the equality constraint to make use of

$$g(\boldsymbol{p}^*) = 0 \Leftrightarrow \sum_k p_k^* = 1$$

○ Let's get back to the MLE problem:

$$\arg\max_{\lambda,\boldsymbol{p}} \mathcal{LLL}(\lambda, \boldsymbol{p}) = \arg\max_{\lambda,\boldsymbol{p}} \sum_{k=1}^{K} x_k \ln p_k + \lambda\Big(\sum_k p_k - 1\Big)$$

□ We can derive a series of equations:

$$\frac{\partial\mathcal{LLL}(\lambda, \boldsymbol{p})}{p_k} = \frac{x_k}{p_k} + \lambda = 0 \Leftrightarrow x_k = -\lambda p_k \quad \Leftrightarrow \sum_k x_k = -\lambda\sum_k p_k \quad \Leftrightarrow N = -\lambda$$

□ Finally: $\quad p_k^* = \dfrac{x_k}{N}$

# What I didn't cover

○ Line search

○ Momentum

○ Inequality constraint

○ Conjugate gradient

○ Stochastic gradient descent

○ You name it

# Reading material

○ Textbook appendix C

○ Jorge Nocedal and Stephen Wright, "Numerical Optimization" http://iucat.iu.edu/catalog/11876818

# Thank You!