# Convex Optimization

Hu SiXing, Hakki Can Karaimer, Pan An, Philipp Keck

National University of Singapore

January 29, 2016

## Linear Regression Example

There should be a picture here.

## Ordinary Least Squares

Input: points $(x_i, y_i)$

Regression line: $y = mx + b$

Objective:

$\min_{m,b} \sum_i (y_i - mx_i - b)^2$

$(\vec{x_i}, y_i)$

$y = \vec{w} \cdot \vec{x} + b$

$\min_{\vec{w}} \sum_i (y_i - \vec{w} \cdot \vec{x_i} - b)^2$

- Easily Solved: $\vec{w}^*(X^\mathsf{T} X) - 1 X^\mathsf{T} \vec{y}$
- But what if $\dim \vec{x}$ is large?
- What about other similar regressions?

## Convex Optimization Problems

- OrdinaryLinearRegression: $\min\limits_{\vec{w}} \sum\limits_{i} (y_i - \vec{w} \cdot \vec{x_i})^2$

- General: $\min\limits_{x} f(x)$ where $f(x)$ is convex

- Set $C$ is convex $\iff \exists x, y \in C, 0 \leqslant t \leqslant 1 : tx + (1-t)y \in C$

- Function $f : \mathbb{R}^n \to \mathbb{R}$ is convex if $\mathbf{dom}\, f$ is convex and $\exists x, y \in \mathbf{dom}\, f, 0 \leqslant t \leqslant 1 :$

$$f(tx + (1-t)y) \leqslant tf(x) + (1-t)f(y)$$

- Unconstrained.

Supposed to be a picture here.

## Outlier Penalty

pic

## Capped Penalty

pic

Huber Penalty Function pic

## Unconstrained Optimization

- Minimize $f(x)$;
- Where $f : \mathbb{R}^n \to \mathbb{R}$ is convex and twice differentiable;
- No additional constraints;
- Assume that unique minimum $x^*$ exists.

## General Principle

- Objective: minimize $f(x)$
- Necessary and sufficient condition: $\nabla f(x^*) = 0$
  - Solve analytically
  - Iterative algorithms

### Iterative Algorithm:

$$x^{(0)}, x^{(1)}, ... \in \operatorname{dom} f$$

$$k \to \infty, f(x^{(k)}) < f(x^*)$$

### Descent Method:

$$x^{(k+1)} = x^{(k)} + t^{(k)} \Delta x^{(k)}, \textbf{s.t.} f(x^{(k+1)}) < f(x^{(k)})$$

## General Descent Method

### Descent Method:

$$x^{(k+1)} = x^{(k)} + t^{(k)}\Delta x^{(k)}, \textbf{s.t.} f(x^{(k+1)}) < f(x^{(k)}) \qquad (1)$$

### Algorithm:

Given $x^{(0)} \in \textbf{dom } f$;
repeat
    Determine a descent direction $\Delta x$;
    Choose a step size $t > 0$;
    Update $x^{(k+1)} = x^{(k)} + t^{(k)}\Delta x^{(k)}$;
until $\Delta x$ is within an acceptable range and is stable;;

### Noticing that f is convex:

$$\nabla f(x^{(k)})^\mathsf{T} \Delta x^{(k)} < 0 \qquad (2)$$

# General Descent Method

**Descent Method:**

$$x^{(k+1)} = x^{(k)} + t^{(k)}\Delta x^{(k)}, \textbf{s.t.} f(x^{(k+1)}) < f(x^{(k)}) \qquad (1)$$

**Theorem**

For a continuously differentiable function f:

$$f \text{ is convex} \gtrless f(x) \leqslant f(y) + f'(y)(x - y)$$

**Proof**

$$f(x^{(k+1)}) \geqslant f(x^{(k)}) + f'(x^{(k)})\Delta x^{(k)}$$
$$\nabla f(x^{(k)}) \leqslant f(x^{(()}k+1)) - f(x^{(k)}) < 0$$

**Noticing that f is convex:**

$$\nabla f(x^{(k)})^{\mathsf{T}}\Delta x^{(k)} < 0 \qquad (2)$$

## General Descent Method

### Descent Method:

$$x^{(k+1)} = x^{(k)} + t^{(k)}\Delta x^{(k)}, \textbf{s.t.} f(x^{(k+1)}) < f(x^{(k)}) \qquad (1)$$

### Algorithm:

Given $x^{(0)} \in \text{dom } f$;

repeat

    Determine a descent direction $\Delta x \Rightarrow \text{Gradient/SteepestDescent}$;

    Choose a step size $t > 0 \Rightarrow$                 $\text{LineSearchAlgo}$ ;

    Update $x^{(k+1)} = x^{(k)} + t^{(k)}\Delta x^{(k)}$;

until $\Delta x$ is within an acceptable range and is stable;;

### Noticing that f is convex:

$$\nabla f(x^{(k)})^\mathsf{T}\Delta x^{(k)} < 0 \qquad (2)$$

## Line Search

$x^{(k+1)} = x^{(k)} + t^{(k)}\Delta x^{(k)}, f(x^{(k+1)}) \leftarrow f(x^{(k)})$

Should be pics

Exact Line Search Method:

$$t = \underset{s \leqslant 0}{\operatorname{argmin}}\{f(x + s\Delta x)\}$$

## Line Search

- Armijo Condition:

$$f(x^{(k)} + t\Delta x^{(k)}) \leqslant f(x^{(k)}) + c_1\alpha\nabla f(x^{(k)})^\mathsf{T}\Delta x^{(k)}, c_1 > 0$$

- Wolfe Conditions (Including Armijo Condition):

$$\nabla f(x^{(k)} + t\Delta x^{(k)})^\mathsf{T} p^{(k)} \geqslant c_2\nabla f(x^{(k)})^\mathsf{T} p^{(k)}, 0 < c_1 < c_2 < 1$$

### Theorem:

Gradient descent will find local minimum if step size $\alpha$ satisfies Wolfe conditions.

Exact Line Search Method:

$$t = \operatorname*{argmin}_{s \leqslant 0}\{f(x + s\Delta x)\}$$

## Line Search

- Armijo Condition:

$$f(x^{(k)} + t\Delta x^{(k)}) \leqslant f(x^{(k)}) + c_1\alpha\nabla f(x^{(k)})^\mathsf{T}\Delta x^{(k)}, c_1 > 0$$

### Algorithm:

Given a descent direction $\Delta x$ for $f$ at
$x \in \operatorname{dom} f, \alpha \in (0, 0.5), \beta \in (0, 1), t = 1;$
repeat
$\quad|\quad t = \beta t;$
until $f(x^{(k)} + t\Delta x^{(k)}) \leqslant f(x^{(k)}) + c_1\alpha\nabla f(x^{(k)})^\mathsf{T}\Delta x^{(k)};;$

Exact Line Search Method:

$$t = \operatorname*{argmin}_{s \leqslant 0}\{f(x + s\Delta x)\}$$

## General Descent Method

- Gradient Descent Method
- Steepest Descent Method

$\Delta x$ satisfies:

$$\nabla f(x^{(k)})^{\mathsf{T}} \Delta x < 0$$

# Gradient Descent Method

### $\Delta x = -\nabla f(x)$

Given $x^{(0)} \in \mathrm{dom}\, f$;
repeat
  $\Delta x = -\nabla f(x^{(k)})$;
  Choose a step size $t > 0, [\mathit{LineSearch}]$;
  Update $x^{(k+1)} = x^{(k)} + t\Delta x$;
until $\Delta x$ is within an acceptable range and is stable;;

## Steepest Descent Method

$\Delta x = \Delta x_{sd}$

Taylor Series:

$$f(\mathbf{x} + \Delta \mathbf{x}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^{\mathsf{T}} \Delta \mathbf{x} + \frac{1}{2} \Delta \mathbf{x} \nabla f(\mathbf{x}) \Delta \mathbf{x}$$

$$f(\mathbf{x} + \mathbf{v}) \approx \hat{f}(\mathbf{x} + \mathbf{v}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^{\mathsf{T}} \mathbf{v}$$

$$x^{(k+1)} = x^{(k)} + t^{(k)} \Delta x^{(k)}, \mathbf{s.t.} f(x^{(k+1)}) < f(x^{(k)})$$

Where **v** is a descent direction if $\nabla f(\mathbf{x})^{\mathsf{T}} < 0$

## Steepest Descent Method

$$f(\mathbf{x} + \mathbf{v}) \approx \hat{f}(\mathbf{x} + \mathbf{v}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\mathsf{T} \mathbf{v}$$

Normalized Steepest Descent Direction:

$$\Delta \mathbf{x}_{\mathbf{nsd}} = \operatorname{argmin}\{\nabla f(\mathbf{x})^\mathsf{T} \mathbf{v} |\, \|\mathbf{v}\| = 1\}$$
$$= \operatorname{argmin}\{\nabla f(\mathbf{x})^\mathsf{T} \mathbf{v} |\, \|\mathbf{v}\| \leqslant 1\} \tag{3}$$

## Steepest Descent Method

Dual Norm, denoted $\|\cdot\|_*$, is defined as:

$$\|z\|_* = \sup\{z^\mathsf{T} x | \|x\| \leqslant 1\}$$

Unnormalized Steepest Descent Direction:

$$\Delta\mathbf{x} = \|\nabla f(x)\|_* \cdot \Delta\mathbf{x}_{nsd}$$

$$
\begin{aligned}
\nabla f(\mathbf{x})^\mathsf{T}\mathbf{v} &= \nabla f(\mathbf{x})^\mathsf{T}\Delta\mathbf{x}_{sd} \\
&= \|f(\mathbf{x})\|_* \nabla f(\mathbf{x})^\mathsf{T}\Delta\mathbf{x}_{nsd} \\
&= -\|\nabla f(\mathbf{x})\|_*^2
\end{aligned}
$$

### Proof

$$
\begin{aligned}
\Delta\mathbf{x}_{nsd} &= \operatorname{argmin}\{\nabla f(\mathbf{x})^\mathsf{T}\mathbf{v} | \|v\| = 1\} \\
&= -\operatorname{argmax}\{\nabla f(\mathbf{x})^\mathsf{T}\mathbf{v} | \|v\| \leqslant 1\}
\end{aligned}
$$

$$\|\nabla f(\mathbf{x})\|_*^2 = \sup\{\nabla f(\mathbf{x})^\mathsf{T}\mathbf{v} | \|v\| \leqslant 1\}$$

$$\Rightarrow \|\nabla f(\mathbf{x})\|_*^2 = -\nabla f(\mathbf{x})^\mathsf{T}\Delta x_{nsd}$$

## Steepest Descent Method

$\Delta\mathbf{x}_{nsd} = \text{argmin}\{\nabla f(\mathbf{x})^{\mathsf{T}}\mathbf{v} | \|\mathbf{v}\| \leqslant 1\}$
$\Delta\mathbf{x}_{sd} = \|\nabla f(\mathbf{x})\|_* \Delta\mathbf{x}_{nsd}$

### Steepest Descent Method

Given $x^{(0)} \in \text{dom} f$;
repeat
> Compute steepest descent direction $\Delta x_{sd}$;
> Choose a step size $t > 0$, [LineSearch];
> Update $x^{(k+1)} = x^{(k)} + t^{(k)}\Delta x_{sd}^{(k)}$;
until $\Delta x$ is within an acceptable range and is stable;

# Descent Method

## General

Given $x^{(0)} \in \text{dom } f$;
repeat
    Determine a descent direction $\Delta x$;
    Choose a step size $t > 0$;
    Update $x^{(k+1)} = x^{(k)} + t^{(k)} \Delta x^{(k)}$;
until $\Delta x$ is within an acceptable range and is stable;

## Gradient Descent

Given $x^{(0)} \in \text{dom } f$;
repeat
    $\Delta x = -\nabla f(x^{(k)})$;
    Choose a step size $t > 0, [LineSearch]$;
    Update $x^{(k+1)} = x^{(k)} + t \Delta x$;
until $\Delta x$ is within an acceptable range and is stable;

## Steepest Descent

Given $x^{(0)} \in \text{dom } f$;
repeat
    Compute steepest descent direction $\Delta x_{sd}$;
    Choose a step size $t > 0, [LineSearch]$;
    Update $x^{(k+1)} = x^{(k)} + t^{(k)} \Delta x_{sd}^{(k)}$;
until $\Delta x$ is within an acceptable range and is stable;

# Descent Method

## General

- $\Delta\mathbf{x}_{nsd} = \mathrm{argmin}\{\nabla f(\mathbf{x})^\mathsf{T}\mathbf{v}|\,\|\mathbf{v}\| \leqslant 1\}$
- $\Delta\mathbf{x}_{sd} = \|\nabla f(x)\|_* \cdot \Delta\mathbf{x}_{nsd}$
- If the norm $\|\cdot\|$ is Euclidean norm, $\Delta\mathbf{x} = -\nabla f(\mathbf{x})$

## Gradient Descent

Given $x^{(0)} \in \mathrm{dom}\,f$;
repeat
    $\Delta x = -\nabla f(x^{(k)})$;
    Choose a step size $t > 0, [LineSearch]$;
    Update $x^{(k+1)} = x^{(k)} + t\Delta x$;
until $\Delta x$ is within an acceptable range and is stable;

## Steepest Descent

Given $x^{(0)} \in \mathrm{dom}\,f$;
repeat
    Compute steepest descent direction $\Delta x_{sd}$;
    Choose a step size $t > 0, [LineSearch]$;
    Update $x^{(k+1)} = x^{(k)} + t^{(k)}\Delta x_{sd}^{(k)}$;
until $\Delta x$ is within an acceptable range and is stable;

## Convex Domain

Linear Regression method is applicable only if nonlinear
function is linear in terms of function parameters:

$$f(x; a) = \sum_{k=1}^{m} a_k h_k(x)$$

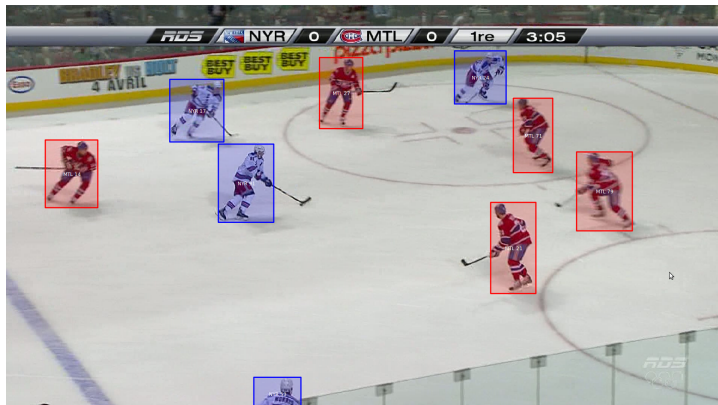Many nonlinear functions are not like that, for example:

$$f_1(x) = \frac{x^2}{a_1 + (x - a_2)}$$

$$f_2(x, y, z) = \frac{x^2}{a_1 + x^2} + \frac{y^2}{a_2 + y^2} + \frac{z^2}{a_3 + z^2}$$

## Condition Number

- The condition number of C gives a measure of its anisotropy or eccentricity.
- If the condition number of a set C is small (say, near one) it means that the set has approximately the same width in all directions, i.e., it is nearly spherical.
- If the condition number is large, it means that the set is far wider in some directions than in others.
- $\text{cond}(f) = \frac{\lambda_{max}(f)}{\lambda_{min}(f)}$
- $\lambda_{max}$ and $\lambda_{min}$ describes minimum and maximum eigenvalues in 2D.

# Image Processing — Lucas-Kanade

Classic examples are optical flow techniques like
Lucas-Kanade (VideoTracking), Horn-Schunck.

## Lucas-Kanade

### Goal of Lucas-Kanade

Minimize the sum of squared error between two images.

### Assumption

The displacement of the image contents between two nearby instants (frames) is small and approximately constant within a neighborhood of the point $p$ under consideration.

## Lucas-Kanade

> ### Optical Flow Equation (2 Dementional)
>
> For a pixel location $(x, y, t)$, the intensity has moved by $\Delta x, \Delta y, \Delta t$, the basic assumption can be represented as:
>
> $$I(x, y, t) = I(x + \Delta x, y + \Delta y, t + \Delta t)$$

## Lucas-Kanade

**Optical Flow Effect:**

For all pixels within a window centered at p:

$$I_x(q_i)V_x + I_y(q_i)V_y = -I_t(q_i)$$

Where $i = 1, 2, 3...n$.

**Abbreviations:**

$$A = [I_x(q_i)^\mathsf{T}, I_y(qi)^\mathsf{T}]$$

$$V = [v_x, v_y]^\mathsf{T}$$

$$b = [-I_t(q_i)]^\mathsf{T}$$

## Lucas-Kanade

### Lucas-Kanade Method Abstraction:

LK method tries to solve $2 \times 2$ system:

$$A^\mathsf{T} A V = A^\mathsf{T} b$$

A.K.A:

$$V = (A^\mathsf{T} A)^{-1} A^\mathsf{T} b$$

### Notice:

$V = [v_x, v_y]^\mathsf{T}$ is variable. Which means that the system does not know the actual velocity of the system.

## Lucas-Kanade

Goal of Lucas-Kanade Method:

To minimize $\|A^{\mathsf{T}}V - b\|^2$.

Basic LK Derivation for Models(Stuff to be Tracked):

$$E[v_x, v_y] = \Sigma[I(x + v_x, y + v_y) - T(x, y)]^2$$

Where $v_x, v_y$ is the hypothesized location of the model(s) to be tracked, and $T(x, y)$ model.

## Lucas-Kanade

### Key Step for Implementation of GD (Step 1):

Generalizing LK approach by introducing warp function W:

$$E[v_x, v_y] = \Sigma[I(W(x, y); P) - T(x, y)]^2$$

Generalizing is used to solve the problem where the constant flow of larger picture frames for a long time is a total waste of calculation power. Warp function examples are Affine and Projective.

The warping function are the convergence factor for steepest descent algorithm.

## Lucas-Kanade

### Key Step for Implementation of GD (Step 2):

The key to the derivation is Taylor series approximation:

$$I(W(x,y); P + \Delta P) \approx I(W([x,y]; P)) + \nabla I \frac{\partial W}{\partial P} \Delta P$$

- The approximation equation is actually the abstract of the basic assumption of optical flow described in the slides before.
- Derivation of this equation can be discussed in forum (Too long for slides).

## Lucas-Kanade

**Some Explainations:**

- Gradient image $\nabla I$
- Image error $I_E = T(x, y) - I(W[x, y]; P)$
- Jacobian matrix $\frac{\partial W}{\partial P}$
- Steepest image $I_S = \nabla I \frac{\partial W}{\partial P}$
- Hessian Matrix $\Sigma (\nabla I \frac{\partial W}{\partial P})^T (\nabla I \frac{\partial W}{\partial P})$
- Iteration step $\Delta P = \Sigma I_S^T I_E$

## Lucas-Kanade

### Algorithms:

- Warp image and get $I(W[x, y]; P)$;

- Get image error $I_E$;

- Warp gradient image $\nabla I$;

- Evaluate Jacobian;

- Compute steepest descent image $I_S = \nabla I \frac{\partial W}{\partial P}$;

- Compute Hessian matrix $\Sigma I_S^T I_S$;

- Get warping step $\Delta P = I_S I_E$;

- Update warping parameter $P = P + \Delta P$;

- Repeat until $\Delta P$ is negligible.

# APPLICATIONS – MACHINE LEARNING

## Generalized Utilization of Convex: Delta Rule

- The delta rule is derived by attempting to minimize the error in the output of the neural network through gradient descent.

- Gradient Descent optimization is the most basic principle for training neurons even with different activation functions.

- Delta rule, can also be modified, if possible, with steepest descent method.

# APPLICATIONS – MACHINE LEARNING

**Delta Rule:**
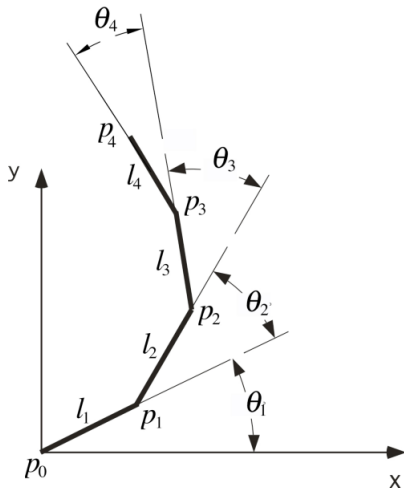
$$\Delta w_{ji} = \alpha(t_j - y_j)g'(h_j)x_j$$

Where $\alpha$ is the learning rate, $g(x)$ is the neuron's activation function. $t_j$ and $y_j$ is the target and actual output of the neuron. $h_j$ is the weighted sum of the neuron's inputs. And $x_i$ is the $i_{th}$ input.

**The above equation holds the following:**

$$h_j = \Sigma x_i w_{ji}$$

$$y_j = g(h_j)$$

# APPLICATIONS – INVERSE KINEMATICS

# APPLICATIONS – INVERSE KINEMATICS

### Goal of Inverse Kinematics

Given a position in the space, calculate a way for a robot hand to reach a place.

### Problem Abstract:

$$\vec{e} = R_1 T_1 R_2 T_2 R_3 T_3 R_4 T_4 \vec{e_0}$$

Where $T_i$ is a series of translation transformation and $R_i$ is a series of rotation translation.

# APPLICATIONS – INVERSE KINEMATICS

Abstraction for Convex Optimization:

$$\Delta\vec{\theta} = \alpha J^\mathsf{T}\vec{e}$$

.

The target for the optimization is to achieve $|\vec{e_p} - \vec{e_t}| = 0$, where $\vec{e_p}$ is th original position of the tip of the robotic arm and $\vec{e_t}$ is the target position. $J$ is the jacobian matrix in terms of $\vec{\theta}$, which is the vector of all the spatial angles of all joints. $\alpha$ is the convergence rate and $\vec{e}$ is the position derivation (step size).

# APPLICATIONS – INVERSE KINEMATICS

## About Inverse Kinematics

- Jacobian transpose is the implementation of gradient descent in the real physical world.

- It can actually achieve near linear solution for robotic arms with a fast convergence rate.