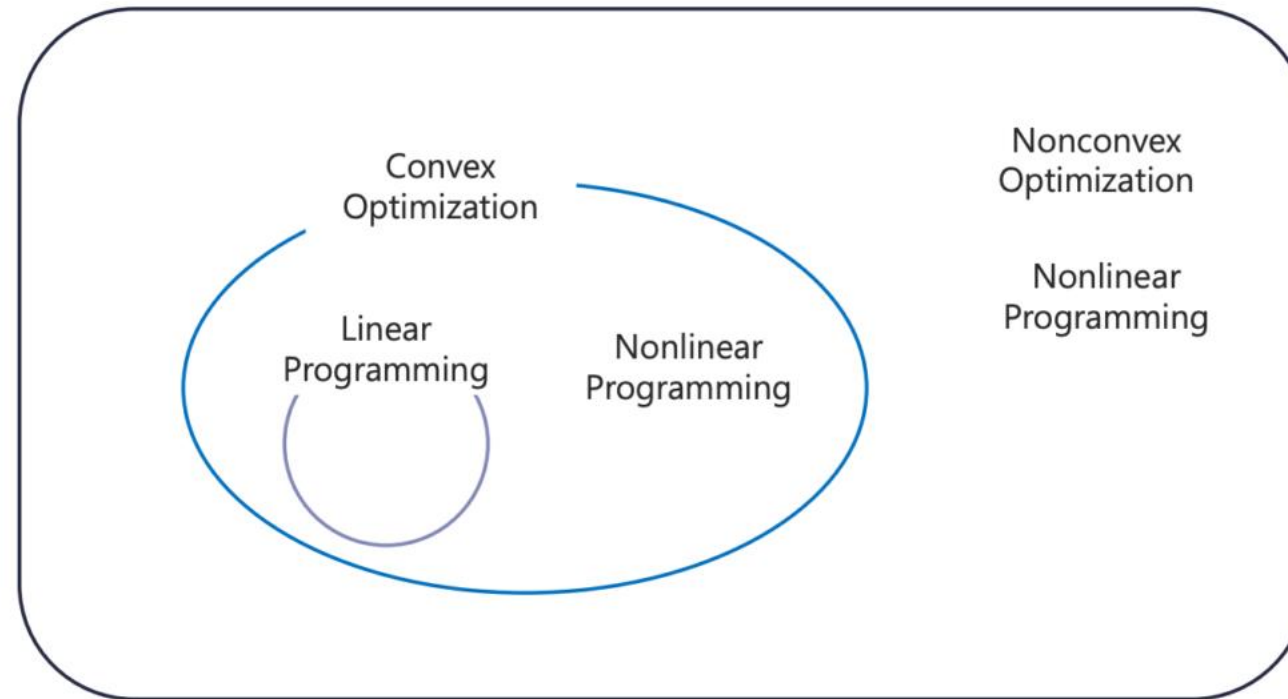


Introduction to Optimization in Deep Learning

(Sejong RCV 임근택)

Convex Optimization?

Optimization Problems



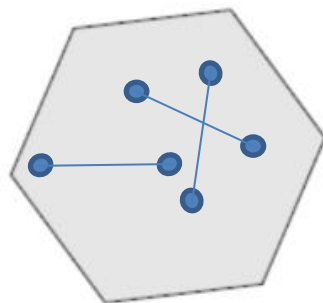
Convex Sets & Function

Convex Set

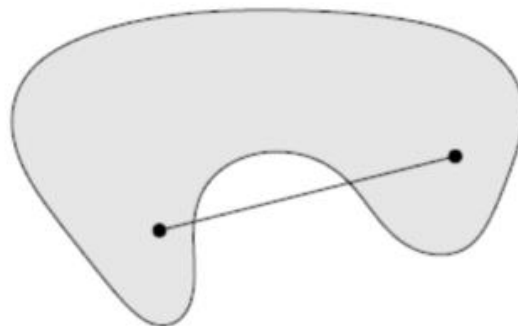
Affine Set : $x_1, x_2 \in C, \theta \in R, \theta x_1 + (1 - \theta)x_2$

Convex Set : $x_1, x_2 \in C, \theta \in R, \theta x_1 + (1 - \theta)x_2, 0 \leq \theta \leq 1$

◇ **Examples of convex and nonconvex sets**



Convex



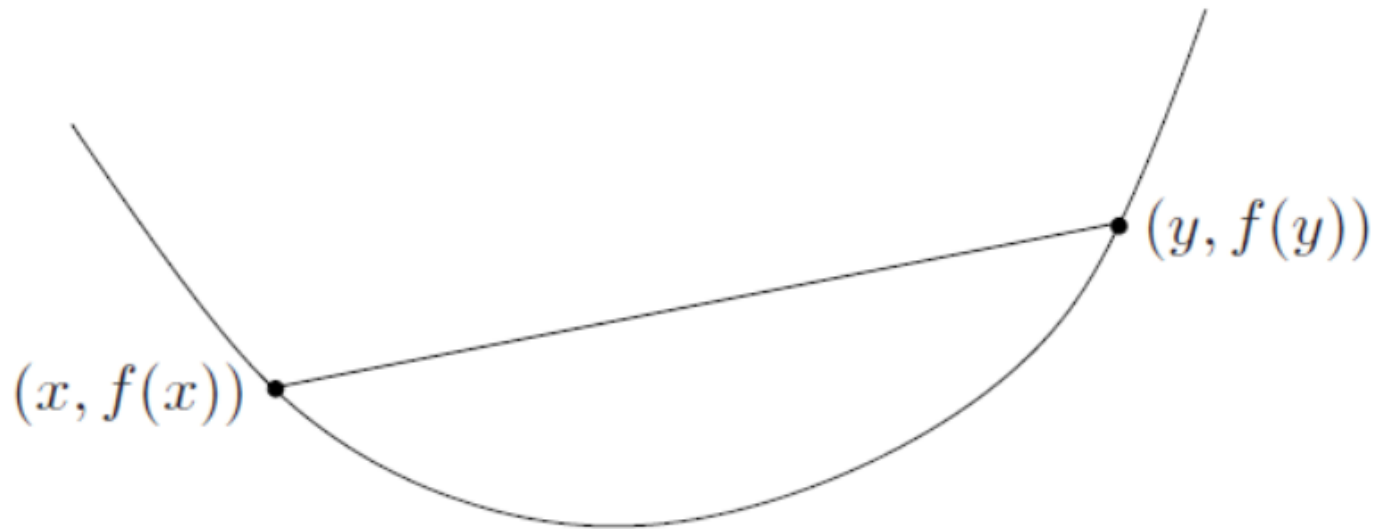
Nonconvex



Nonconvex

Convex Function

Convex Function : $f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y), 0 \leq \theta \leq 1$



Convex Optimization

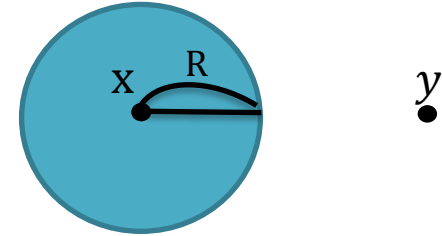
$$\begin{array}{ll}\text{minimize} & f_0(x) \\ \text{subject to} & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & a_i^T x = b_i, \quad i = 1, \dots, p,\end{array}$$

- The objective function must be convex
- The inequality constraint function must be convex
- The equality constraint function must be affine
- The feasible set of a convex optimization problem is convex

Global optimal = local optimal

- Suppose that x is locally optimal point and x is feasible

$$f_0(x) = \inf\{f_0(z) \mid z \text{ is feasible, } \|z - x\|_2 \leq R\}$$



- Suppose that x is not globally optimal there is a feasible point y such that $\|y - x\|_2 > R$

$$f_0(y) < f_0(x)$$

- Consider the point z given by

$$z = (1 - \theta)x + \theta y, \quad \theta = \frac{R}{2\|y - x\|_2}$$

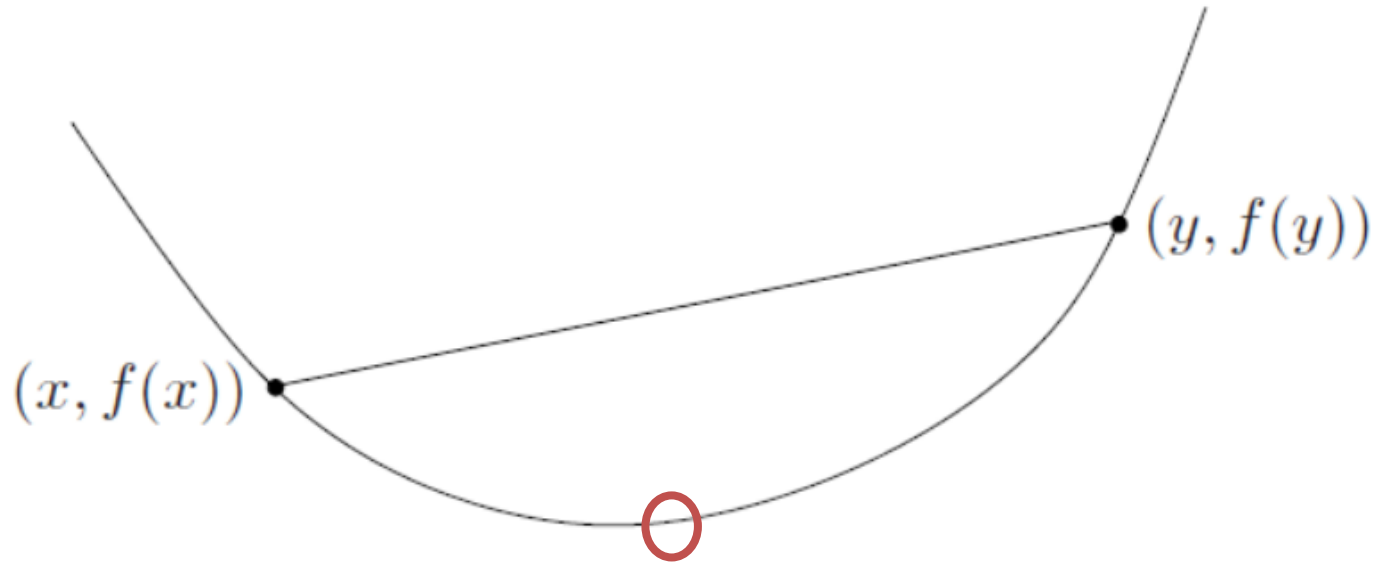
- By Convexity of the function we can have

$$f_0(z) \leq (1 - \theta)f_0(x) + \theta f_0(y) \leq f_0(x)$$

Unconstrained Convex Optimization

Unconstrained Convex Optimization

Minimize $f(x)$



we just need to find $\nabla_x f_0(x) = 0$

Example : Least Square

$$\text{minimize } \| Ax - b \|_2^2$$

$$\| Ax - b \|_2^2 = (Ax - b)^T (Ax - b)$$

$$(Ax - b)^T (Ax - b) = (Ax)^T (Ax) - (Ax)^T b - b^T (Ax) + b^T b$$

$$(Ax)^T (Ax) - (Ax)^T b - b^T (Ax) + b^T b = x^T (A^T A) x - 2(A^T b)^T x + b^T b$$

$$f(x) = x^T (A^T A) x - 2(A^T b)^T x + b^T b$$

$$\nabla f(x) = 2(A^T A) x - 2(A^T b) = 0$$

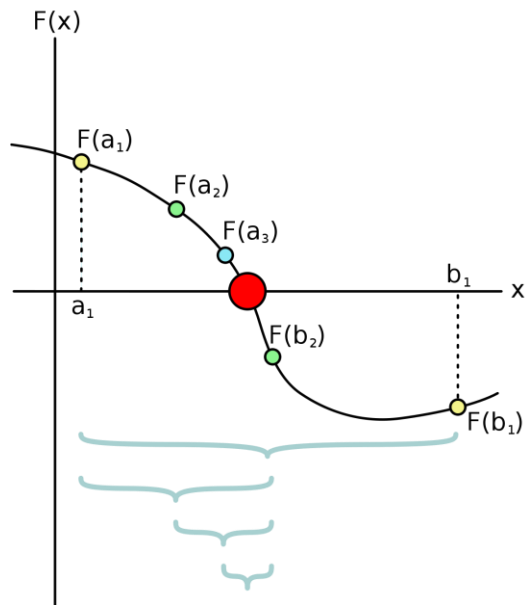
$$x = (A^T A)^{-1} A^T b$$



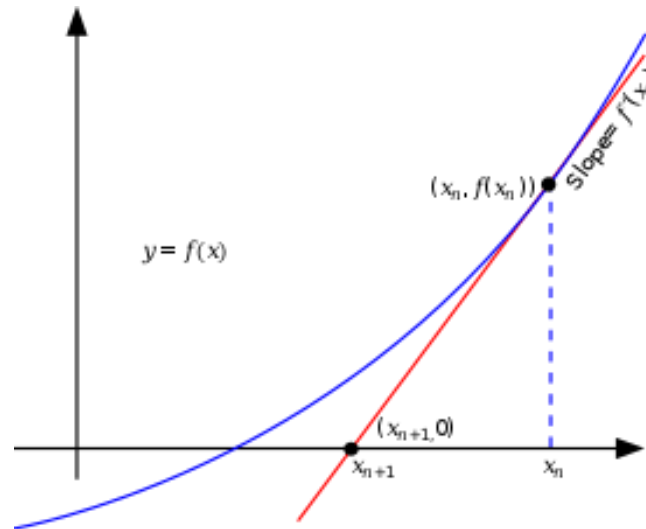
Vector Differentiation

Unconstrained Convex Optimization (Numerical)

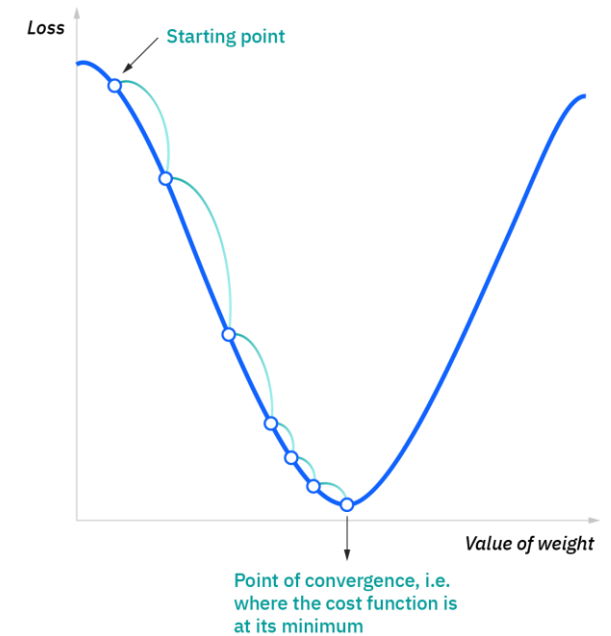
we just need to find $\nabla_x f_0(x) = 0$



Bisection method



Newton's method



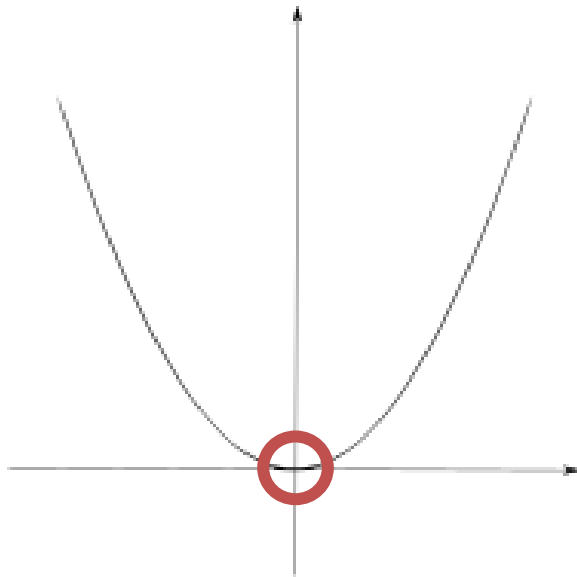
Gradient Descent

- To calculate roots of gradient, the exact solution may not be calculated.
- So, sometimes we use **numerical approach**

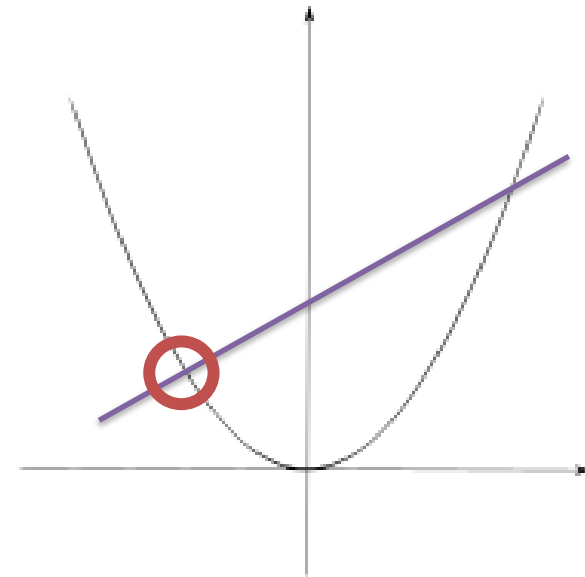
Constrained Convex Optimization

Constrained Convex Optimization

$$\begin{array}{ll}\text{minimize} & f_0(x) \\ \text{subject to} & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & a_i^T x = b_i, \quad i = 1, \dots, p,\end{array}$$



Unconstrained



Constrained

Lagrange multiplier

An optimization problem

$$\begin{aligned}
 \text{(P1)} \quad & \text{minimize} && f_0(x) \\
 & \text{subject to} && f_i(x) \leq 0, \quad i = 1, \dots, m \\
 & && h_i(x) = 0, \quad i = 1, \dots, p
 \end{aligned}$$

$$\mathcal{D} = \left(\bigcap_{i=0}^m \text{dom } f_i \right) \cap \left(\bigcap_{i=1}^p \text{dom } h_i \right)$$

- ◆ We do not assume this problem is convex.

Lagrangian $L : \mathbf{R}^n \times \mathbf{R}^m \times \mathbf{R}^p \rightarrow \mathbf{R}$

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x)$$

- ◆ λ_i : Lagrange multiplier associated with the i -th inequality constraint
- ◆ ν_i : Lagrange multiplier associated with the i -th equality constraint
- ◆ λ, \mathbf{v} : Lagrange multiplier vectors or dual variables associated with **Problem (P1)**

Lagrange Dual Function

The Lagrange dual function

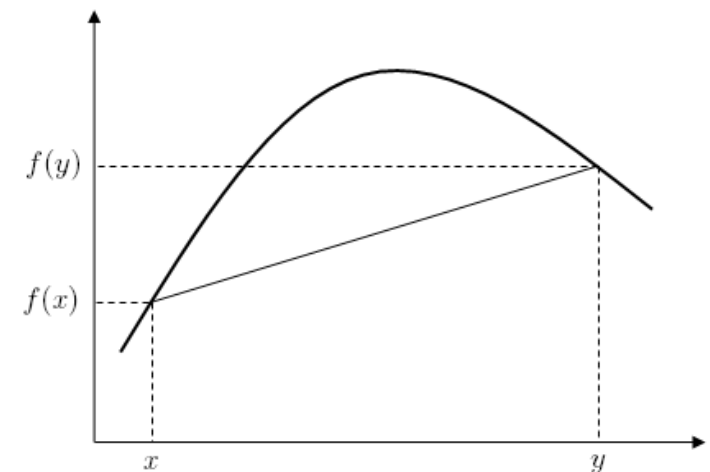
$$g(\lambda, \nu) = \inf_{x \in \mathcal{D}} L(x, \lambda, \nu) = \inf_{x \in \mathcal{D}} \left(f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x) \right)$$

- ◆ The Lagrange dual function is always **concave** in (λ, ν) even when **Problem (P1)** is not convex since it is the pointwise infimum of a family of affine functions of (λ, ν) .

Lower bound on the optimal value p^*

- ◆ The dual function yields lower bound on the optimal value of **Problem (P1)**.
- ◆ For any $\lambda \succeq 0$ and any ν , we have

$$g(\lambda, \nu) \leq p^*.$$



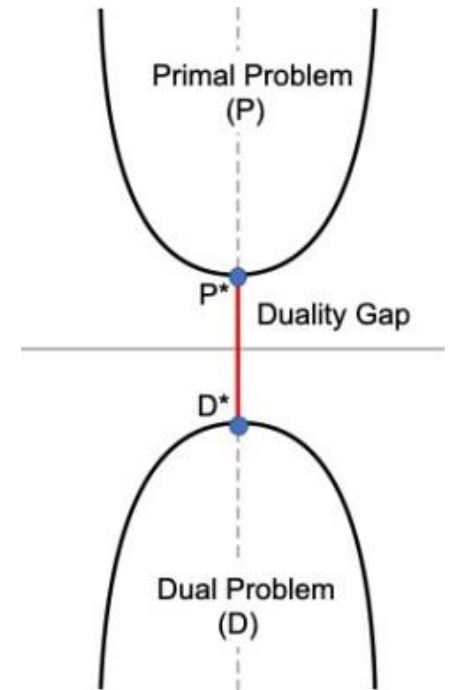
Concave function

Lagrange Dual Problem

The Lagrange dual problem

$$\begin{array}{ll} \text{(P2)} & \text{maximize} \quad g(\lambda, \nu) \\ & \text{subject to} \quad \lambda \succeq 0 \end{array}$$

- ◆ **Problem (P1)** is called the primal problem.
- ◆ The Lagrange dual problem is a convex optimization problem even when the primal problem is not convex.



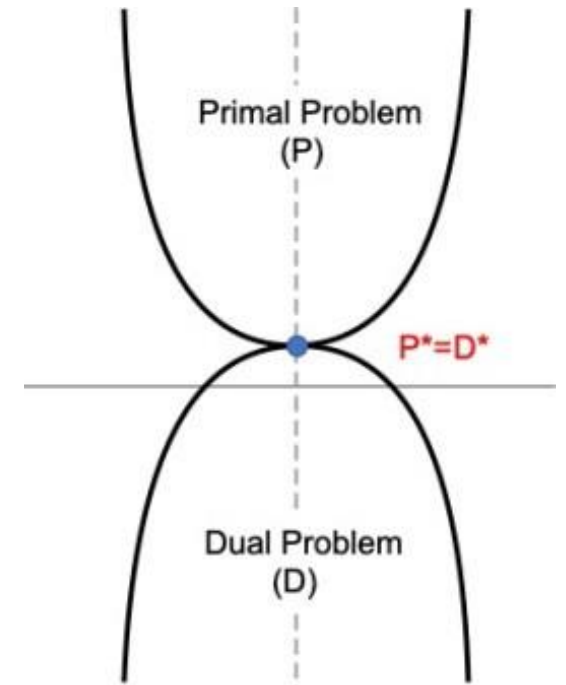
Duality Gap

Strong Duality

Strong duality

- ◆ If the equality $d^* = p^*$ holds, i.e., the optimal duality gap is zero, then we say that strong duality holds.
- ◆ Strong duality does not, in general, hold.
- ◆ If the primal problem is convex, we usually (but not always) have strong duality.

$$\begin{array}{ll}\text{minimize} & f_0(x) \\ \text{subject to} & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & Ax = b\end{array}$$



Strong duality

Slater's condition

◆ Slater's condition (one of constraint qualifications)

- ◆ There exists a *strictly feasible* $x \in \mathbf{relint} \mathcal{D}$, (\mathcal{D} : domain of the problem)

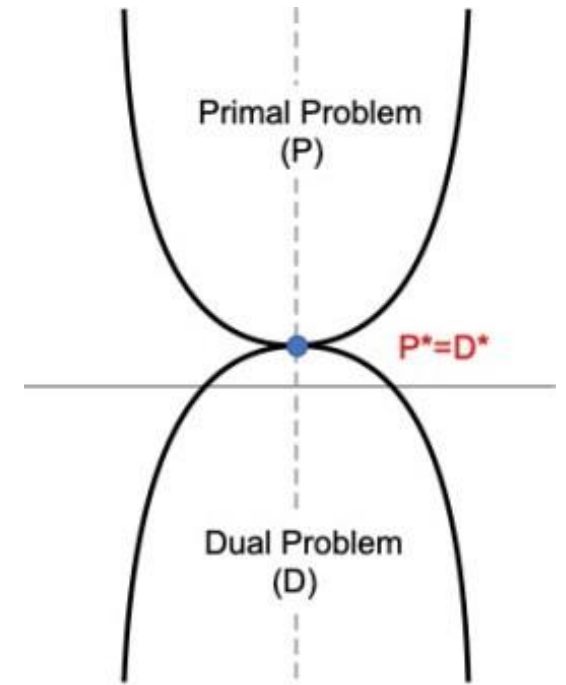
$$f_i(x) < 0, \quad i = 1, \dots, m, \quad Ax = b.$$

◆ Weaker form of Slater's condition

- ◆ If the first k constraint functions f_1, \dots, f_k are affine, then strong duality holds provided the following weaker condition holds.
- ◆ There exist an $x \in \mathbf{relint} \mathcal{D}$ with

$$f_i(x) \leq 0, \quad i = 1, \dots, k, \quad f_i(x) < 0, \quad i = k + 1, \dots, m,$$

$$Ax = b.$$



Strong duality

Karush-Kuhn-Tucker condition(KKT condition)

KKT optimality condition for **convex problems**

- ◆ When the primal problem is convex, the KKT conditions are also sufficient.
- ◆ Let $\tilde{x}, \tilde{\lambda}, \tilde{\nu}$ be any points that satisfy the KKT conditions.

$$f_i(\tilde{x}) \leq 0, \quad i = 1, \dots, m$$

$$h_i(\tilde{x}) = 0, \quad i = 1, \dots, p$$

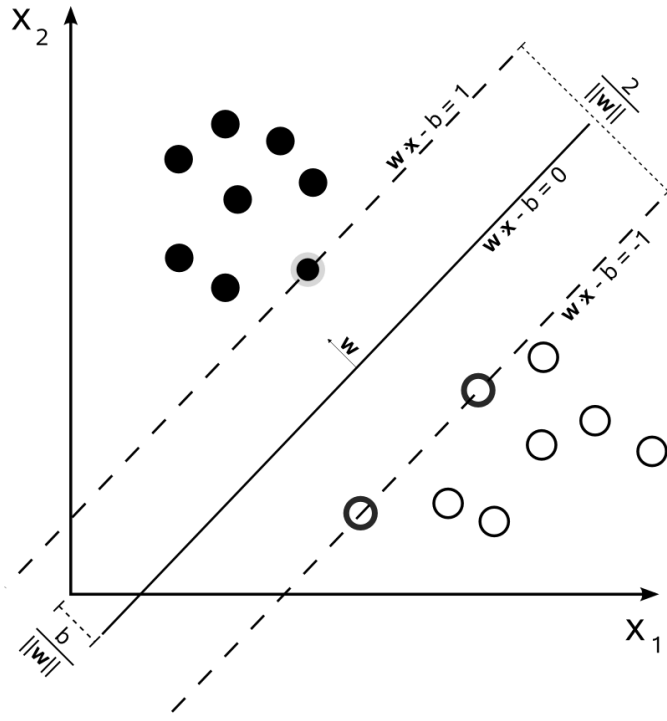
$$\tilde{\lambda}_i \geq 0, \quad i = 1, \dots, m$$

$$\tilde{\lambda}_i f_i(\tilde{x}) = 0, \quad i = 1, \dots, m$$

$$\nabla f_0(\tilde{x}) + \sum_{i=1}^m \tilde{\lambda}_i \nabla f_i(\tilde{x}) + \sum_{i=1}^p \tilde{\nu}_i \nabla h_i(\tilde{x}) = 0$$

- ◆ Then, \tilde{x} and $(\tilde{\lambda}, \tilde{\nu})$ are primal and dual optimal, with zero duality gap.

Example : Support Vector Machine(SVM)



$$\min_{w,b} \frac{1}{2} \|w\|_2^2$$

$$\text{subject to } -\{y_i(w^T x_i + b) - 1\} \leq 0$$

$$L(w, b, \lambda) = \frac{1}{2} \|w\|_2^2 - \sum_{i=1}^N \lambda_i \{y_i (w^T x_i + b) - 1\}$$

Derive Dual Function

$$g(\lambda) = \inf_{w,b} \left(\frac{1}{2} \|w\|_2^2 - \sum_{i=1}^N \lambda_i \{y_i (w^T x_i + b) - 1\} \right)$$

Example : Support Vector Machine (Dual Function)

$$\frac{\partial L(w, b, \lambda)}{\partial w} = w - \sum_{i=1}^N \lambda_i y_i x_i = 0 \rightarrow w = \sum_{i=1}^N \lambda_i y_i x_i$$

$$\frac{\partial L(w, b, \lambda)}{\partial b} = - \sum_{i=1}^N \lambda_i y_i = 0 \rightarrow \sum_{i=1}^N \lambda_i y_i = 0$$

Substitute

$$g(\lambda) = \inf_{w, b} \left(\frac{1}{2} \|w\|_2^2 - \sum_{i=1}^N \lambda_i \{y_i (w^T x_i + b) - 1\} \right)$$

$$g(\lambda) = \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j$$

Example : Support Vector Machine (Dual Problem)

$$\max_{\lambda} \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j$$

$$\text{subject to } 0 \leq \lambda_i \text{ and } \sum_{i=1}^N \lambda_i y_i = 0$$

Example : Support Vector Machine (KKT condition)

$$\min_{w,b} \frac{1}{2} \| w \|^2$$

$$\text{subject to } -\{y_i(w^T x_i + b) - 1\} \leq 0$$

$$\text{Stationarity : } \nabla_x f_0(x) + \sum_{i=1}^N \lambda_i \nabla_x f_i(x) = 0$$

$$\text{Dual feasibility : } 0 \leq \lambda_i \quad i = 1, \dots, N$$

$$\frac{\partial L(w, b, \lambda)}{\partial w} = w - \sum_{i=1}^N \lambda_i y_i x_i = 0 \rightarrow w = \sum_{i=1}^N \lambda_i y_i x_i$$

$$\frac{\partial L(w, b, \lambda)}{\partial b} = -\sum_{i=1}^N \lambda_i y_i = 0 \rightarrow \sum_{i=1}^N \lambda_i y_i = 0$$

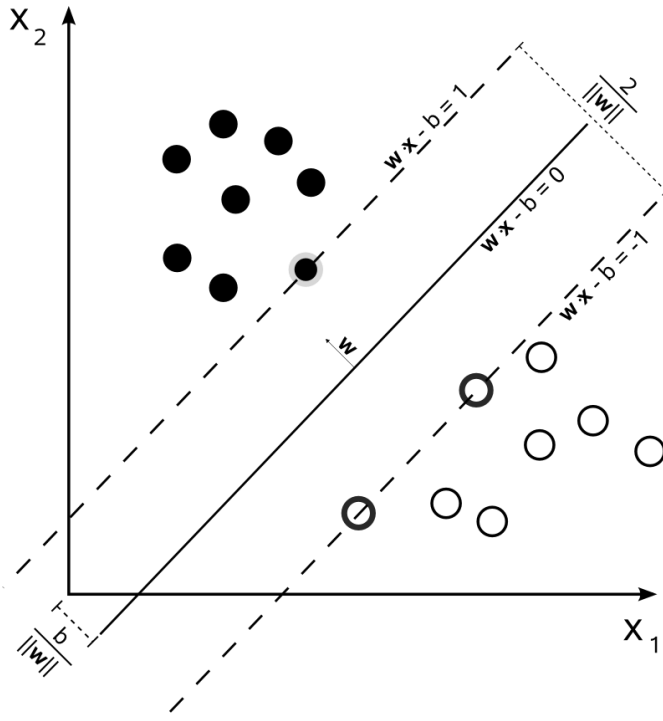
$$\text{Primal feasibility : } 1 - y_i(w^T x_i + b) \leq 0$$

$$\text{complementary slackness : } \lambda_i (y_i(w^T x_i + b) - 1) = 0$$

Example : Support Vector Machine (KKT condition)

Primal feasibility : $1 - y_i(w^T x_i + b) \leq 0$

complementary slackness : $\lambda_i (y_i(w^T x_i + b) - 1) = 0$



if $\lambda_i \neq 0 \rightarrow y_i (w^T x_i + b) = 1$ (x_i is support vector)

$$w = \sum_{i=1}^N \lambda_i y_i x_i \quad b = \frac{1}{y_i} - w^T x_i$$

Gradient Descent

Loss function

A **loss function** tells how good our current classifier is

Low loss = good classifier

High loss = bad classifier

(Also called: **objective function**;
cost function)

Given a dataset of examples

$$\{(x_i, y_i)\}_{i=1}^N$$

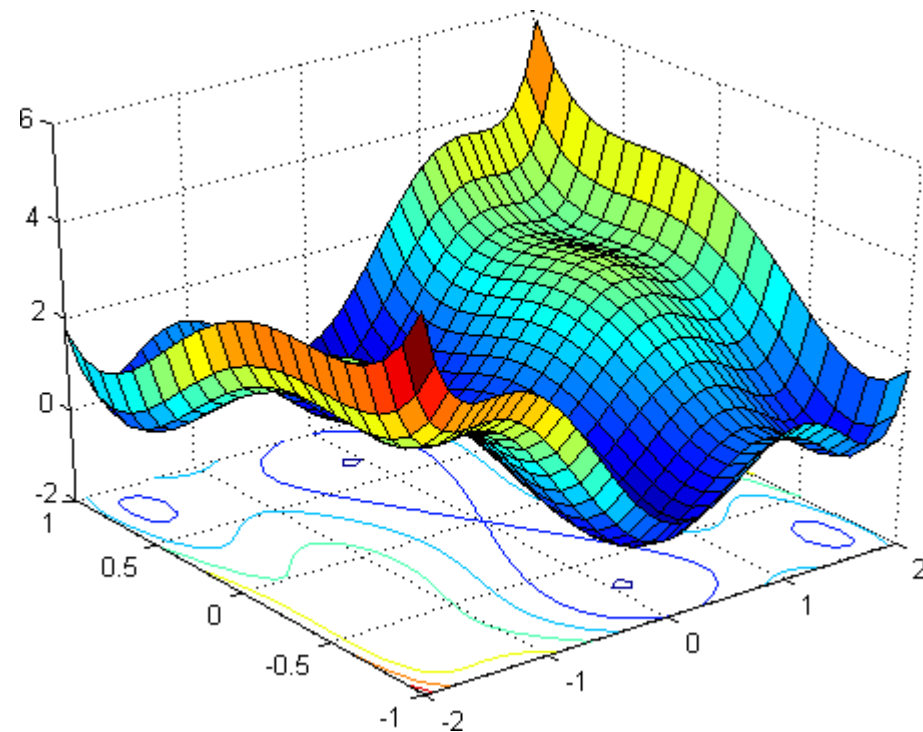
Where x_i is a vector (image) and
 y_i is an integer-valued label

Loss for a single example is

$$L_i = L(f(x_i, W), y_i)$$

Loss for the dataset is average of per-example losses:

분류기의
Loss → $\phi(W) = \frac{1}{N} \sum_{i=1}^N L(f(x_i, W), y_i)$

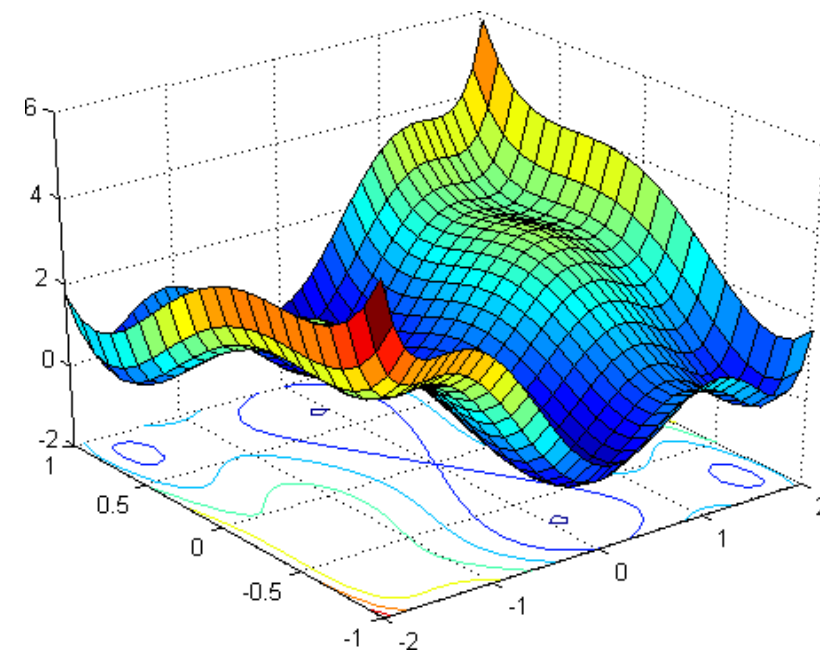
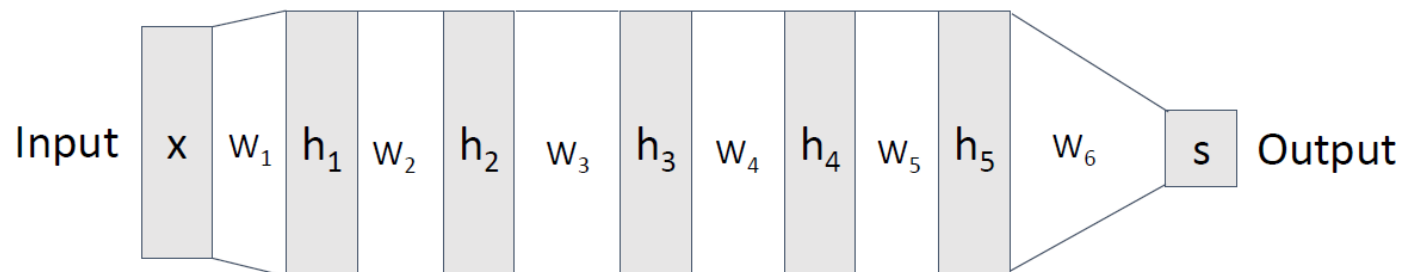


Gradient Descent

Goal:

find the \mathbf{w} minimizing some loss function L .

$$\arg \min_{\mathbf{w} \in \mathbb{R}^N} L(\mathbf{w})$$



Gradient Descent

Method: at each step, move in direction of negative gradient

```
w0 = initialize()                                #initialize
for iter in range(numIters):
    g =  $\nabla_w L(\mathbf{w})$                         # eval gradient
    w = w + -stepsize*g                        # update w
return w
```



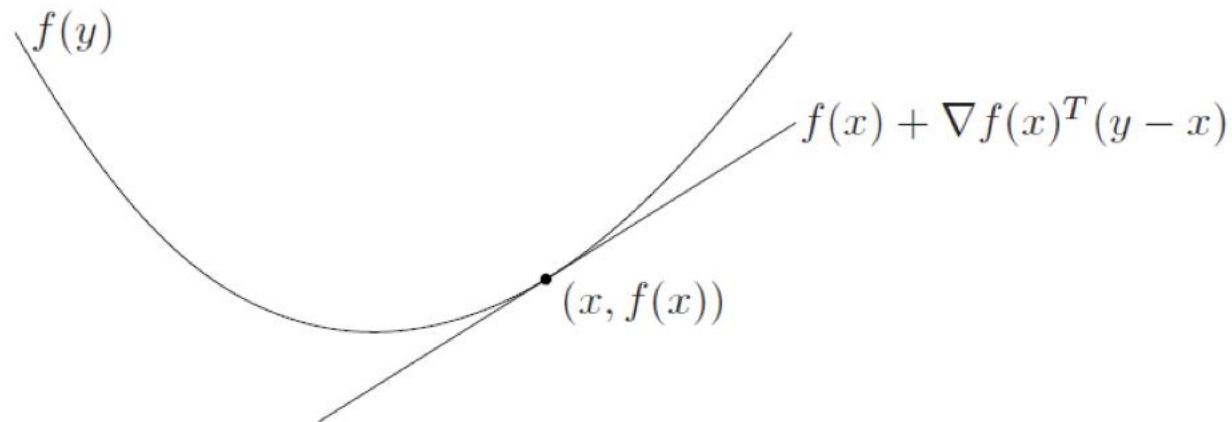
Learning rate

Convergence Analysis of Gradient Descent

Differentiable Convex Function

Lemma 1.2. Let f be twice continuously differentiable. Then f is convex if either of the following hold

- **First-order condition** $f(x) + \nabla f(x)^T(y - x) \leq f(y)$
- **Second-order condition** $0 \preceq \nabla^2 f(x)$



L-Smoothness

A differential function f is said to be L -smooth if its gradients are Lipschitz continuous, that is

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|. \quad (9)$$

Lemma 1.3. Let $f : \mathbb{R}^d \mapsto \mathbb{R}$ be a twice differentiable function. If f is L -smooth then the following holds

$$\langle \nabla^2 f(x)v, v \rangle \leq L\|v\|_2^2, \quad \forall x, v \in \mathbb{R}^d, \quad (10)$$

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2}\|y - x\|_2^2. \quad (11)$$

Convergence for convex and smooth functions

6.1.1 Convergence of gradient descent with fixed step size

Theorem 6.1 *Suppose the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and differentiable, and that its gradient is Lipschitz continuous with constant $L > 0$, i.e. we have that $\|\nabla f(x) - \nabla f(y)\|_2 \leq L\|x - y\|_2$ for any x, y . Then if we run gradient descent for k iterations with a fixed step size $t \leq 1/L$, it will yield a solution $f^{(k)}$ which satisfies*

$$f(x^{(k)}) - f(x^*) \leq \frac{\|x^{(0)} - x^*\|_2^2}{2tk}, \quad (6.1)$$

where $f(x^*)$ is the optimal value. Intuitively, this means that gradient descent is guaranteed to converge and that it converges with rate $O(1/k)$.

Convergence for convex and smooth functions

$$f(y) = f(x) + \nabla f(x)^T (y - x) + \frac{1}{2} \nabla^2 f(x) \|y - x\|_2^2$$
$$f(y) \leq f(x) + \nabla f(x)^T (y - x) + \frac{1}{2} L \|y - x\|_2^2$$

Now let's plug in the gradient descent update by $x^+ = x - t \nabla f(x)$

$$f(x^+) \leq f(x) + \nabla f(x)^T (x - t \nabla f(x) - x) + \frac{1}{2} L \|x - t \nabla f(x) - x\|_2^2$$

$$f(x^+) \leq f(x) - \nabla f(x)^T (t \nabla f(x)) + \frac{1}{2} L \|t \nabla f(x)\|_2^2$$

$$f(x^+) \leq f(x) - t \|\nabla f(x)\|_2^2 + \frac{1}{2} L t^2 \|\nabla f(x)\|_2^2$$

$$f(x^+) \leq f(x) - (1 - \frac{1}{2} L t) t \|\nabla f(x)\|_2^2$$

Convergence for convex and smooth functions

$$f(x^+) \leq f(x) - (1 - \frac{1}{2}Lt)t \|\nabla f(x)\|_2^2$$

Using $t \leq 1/L$, we know that $-1(1 - \frac{1}{2}Lt) \leq \frac{1}{2}L(1/L) - 1 = -\frac{1}{2}$

$$f(x^+) \leq f(x) - \frac{1}{2}t \|\nabla f(x)\|_2^2$$

Next we can bound $f(x^+)$ in terms of $f(x^*)$. Since f is convex, we can write

$$f(x) + \nabla f(x)^T(x^* - x) \leq f(x^*)$$

$$f(x) \leq f(x^*) + \nabla f(x)^T(x - x^*)$$

$$f(x^+) \leq f(x^*) + \nabla f(x)^T(x - x^*) - \frac{1}{2}t \|\nabla f(x)\|_2^2$$

Convergence for convex and smooth functions

$$f(x^+) \leq f(x^*) + \nabla f(x)^T(x - x^*) - \frac{1}{2}t \|\nabla f(x)\|_2^2$$

$$f(x^+) - f(x^*) \leq \frac{1}{2t}(2t \nabla f(x)^T(x - x^*) - t^2 \|\nabla f(x)\|_2^2)$$

$$f(x^+) - f(x^*) \leq \frac{1}{2t}(2t \nabla f(x)^T(x - x^*) - t^2 \|\nabla f(x)\|_2^2 - \|x - x^*\|_2^2 + \|x - x^*\|_2^2)$$

Since the square of $\|x - t \nabla f(x) - x^*\|_2^2$ yields $\|x - x^*\|_2^2 - 2t \nabla f(x)^T(x - x^*) + t^2 \|\nabla f(x)\|_2^2$

$$f(x^+) - f(x^*) \leq \frac{1}{2t}(\|x - x^*\|_2^2 - \|x - t \nabla f(x) - x^*\|_2^2)$$

Notice that by definition we have $x^+ = x - t \nabla f(x)$

$$f(x^+) - f(x^*) \leq \frac{1}{2t}(\|x - x^*\|_2^2 - \|x^+ - x^*\|_2^2)$$

Convergence for convex and smooth functions

$$f(x^+) - f(x^*) \leq \frac{1}{2t} (\|x - x^*\|_2^2 - \|x^+ - x^*\|_2^2)$$

This inequality holds for x^+ on every iteration of gradient descent. Summing over iterations, we get :

$$\sum_{i=1}^k f(x^{(i)}) - f(x^*) \leq \sum_{i=1}^k \frac{1}{2t} (\|x^{(i-1)} - x^*\|_2^2 - \|x^{(i)} - x^*\|_2^2)$$

$$\sum_{i=1}^k f(x^{(i)}) - f(x^*) \leq \frac{1}{2t} (\|x^{(0)} - x^*\|_2^2 - \|x^{(k)} - x^*\|_2^2)$$

$$\sum_{i=1}^k f(x^{(i)}) - f(x^*) \leq \frac{1}{2t} (\|x^{(0)} - x^*\|_2^2)$$

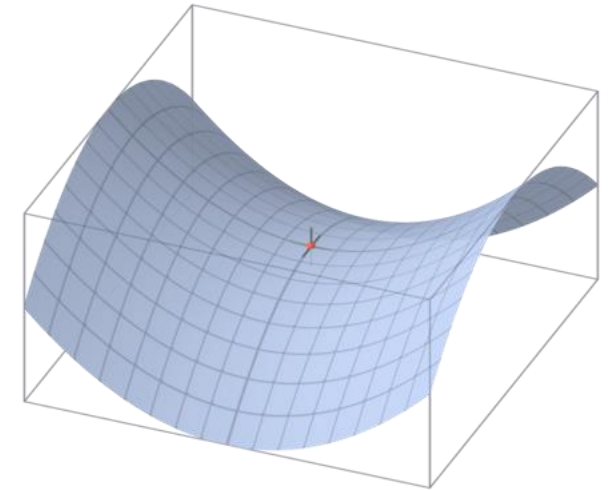
$$f(x^{(k)}) - f(x^*) \leq \frac{1}{k} \sum_{i=1}^k f(x^{(i)}) - f(x^*) \leq \frac{\|x^{(0)} - x^*\|_2^2}{2tk}$$

QED

Advanced Gradient Descent

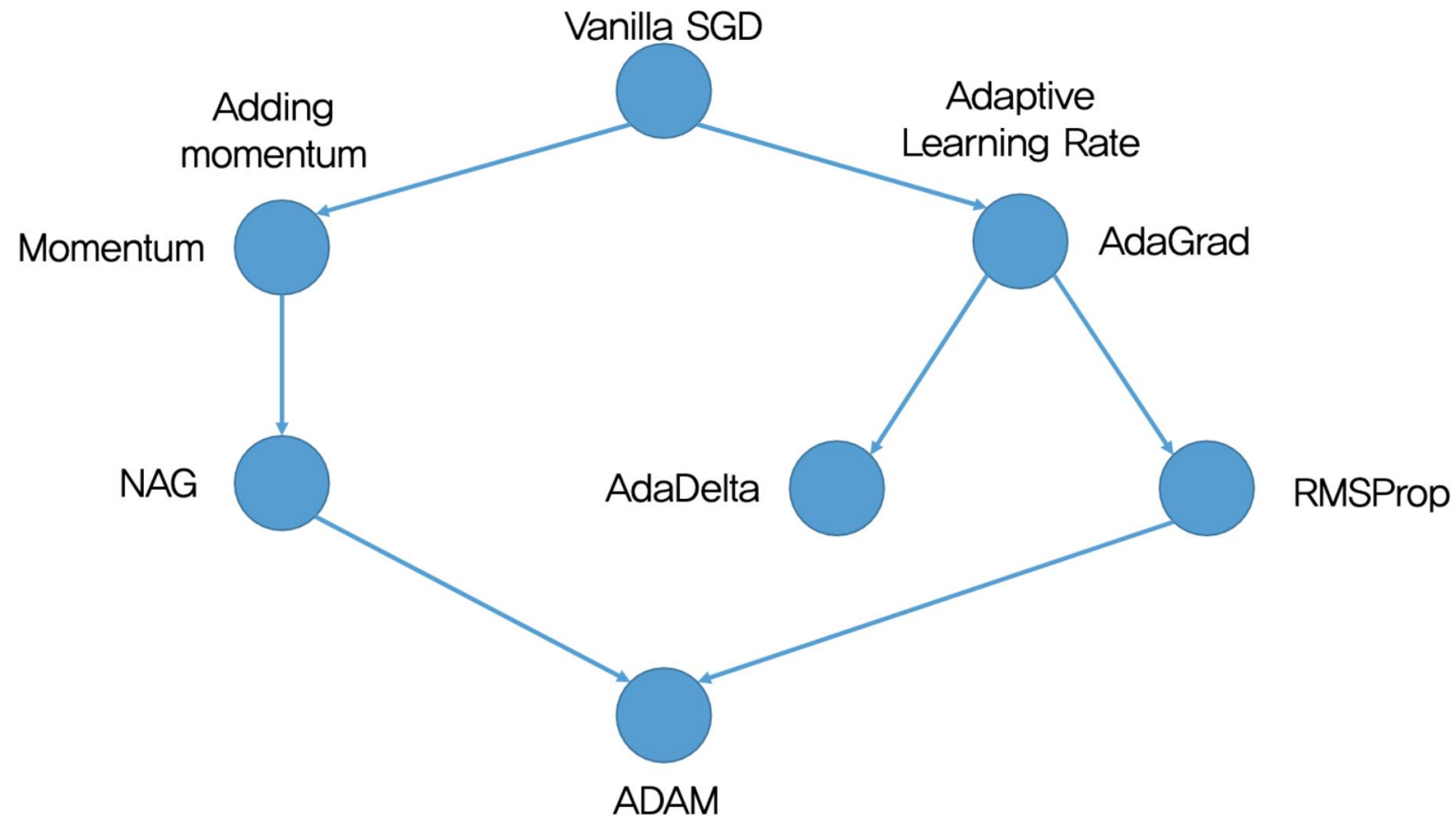
Challenges of Gradient Descent in Non-Convex

- Choosing a proper learning rate can be difficult
- The same learning rate applies to all parameter update
- It is difficult to avoid getting trapped in their numerous local minima, saddle point



Saddle Point

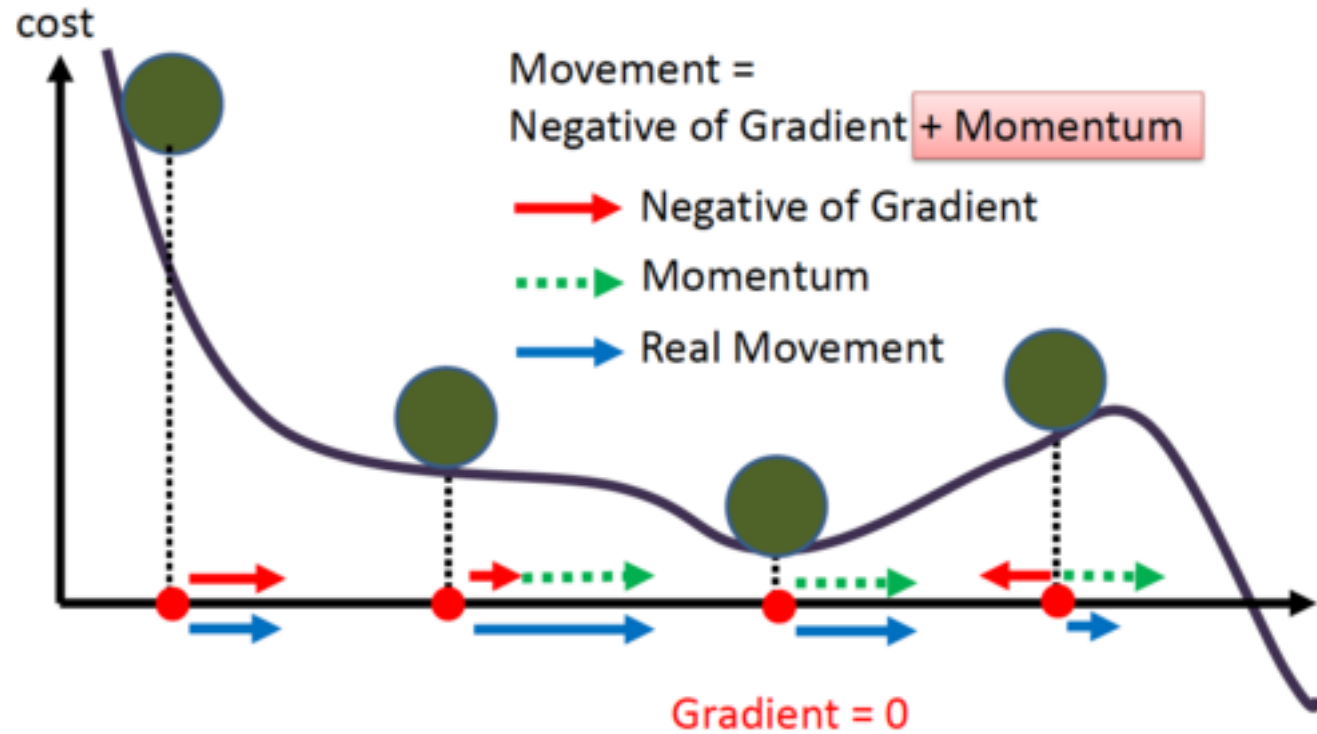
Overview of GD algorithms



Momentum

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta)$$

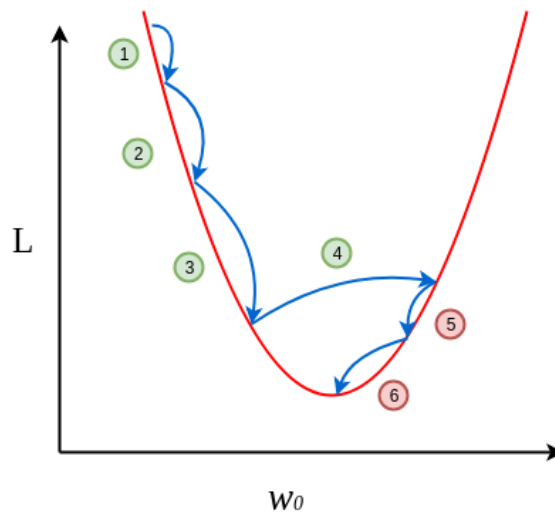
$$\theta = \theta - v_t$$



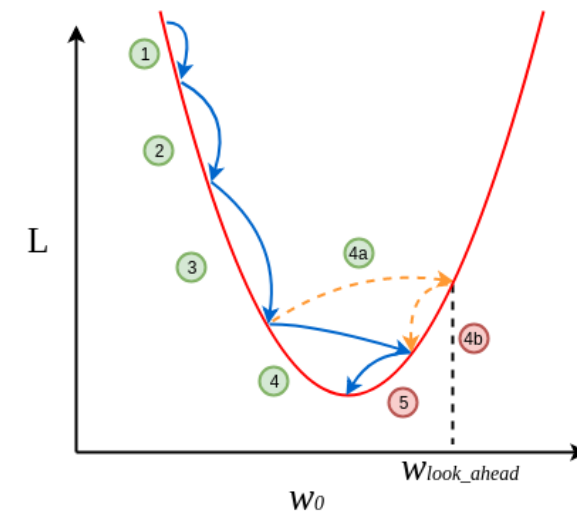
Nesterov Momentum

$$v_t = \gamma v_{t-1} + \eta \nabla J_{\theta}(\theta - \gamma v_{t-1})$$

$$\theta = \theta - v_t$$



(a) Momentum-Based Gradient Descent



(b) Nesterov Accelerated Gradient Descent

$$\text{Green Circle} \Rightarrow \frac{\partial L}{\partial w_0} = \frac{\text{Negative}(-)}{\text{Positive}(+)}$$

$$\text{Red Circle} \Rightarrow \frac{\partial L}{\partial w_0} = \frac{\text{Negative}(-)}{\text{Negative}(-)}$$

AdaGrad

$$G_t = G_{t-1} + (\nabla_{\theta} J(\theta))^2$$

$$G_t = [G_t^{(1)}, G_t^{(2)}, \dots, G_t^{(n-1)}, G_t^{(n)}]$$

$$G_t = [\sum_{\tilde{t}=0}^t (\frac{\partial J(\theta_{\tilde{t}})}{\partial \theta_{\tilde{t}}^{(1)}})^2, \dots, \sum_{\tilde{t}=0}^t (\frac{\partial J(\theta_{\tilde{t}})}{\partial \theta_{\tilde{t}}^{(n)}})^2]$$

$$\theta_t - \frac{\eta}{\sqrt{G_t} + \varepsilon} \odot \nabla_{\theta_t} J(\theta_t)$$

- Large update of a parameter -> small learning rate
- Small update of a parameter -> large learning rate

RMSprop

$$G_t = \gamma G_{t-1} + (1 - \gamma) (\nabla_{\theta} J(\theta))^2$$

$$\theta_t \leftarrow \frac{\eta}{\sqrt{G_t} + \varepsilon} \odot \nabla_{\theta_t} J(\theta_t)$$

Algorithm 8.5 The RMSProp algorithm

Require: Global learning rate ϵ , decay rate ρ .

Require: Initial parameter θ

Require: Small constant δ , usually 10^{-6} , used to stabilize division by small numbers.

Initialize accumulation variables $\mathbf{r} = 0$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Accumulate squared gradient: $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$

 Compute parameter update: $\Delta \theta = -\frac{\epsilon}{\sqrt{\delta + \mathbf{r}}} \odot \mathbf{g}$. ($\frac{1}{\sqrt{\delta + \mathbf{r}}}$ applied element-wise)

 Apply update: $\theta \leftarrow \theta + \Delta \theta$

end while

Adam

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) \nabla_{\theta} J(\theta)$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) (\nabla_{\theta} J(\theta))^2$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t$$

Algorithm 1: *Adam*, our proposed algorithm for stochastic optimization. See section 2 for details, and for a slightly more efficient (but less clear) order of computation. g_t^2 indicates the elementwise square $g_t \odot g_t$. Good default settings for the tested machine learning problems are $\alpha = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$. All operations on vectors are element-wise. With β_1^t and β_2^t we denote β_1 and β_2 to the power t .

Require: α : Stepsize

Require: $\beta_1, \beta_2 \in [0, 1)$: Exponential decay rates for the moment estimates

Require: $f(\theta)$: Stochastic objective function with parameters θ

Require: θ_0 : Initial parameter vector

$m_0 \leftarrow 0$ (Initialize 1st moment vector)

$v_0 \leftarrow 0$ (Initialize 2nd moment vector)

$t \leftarrow 0$ (Initialize timestep)

while θ_t not converged **do**

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep t)

$m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate)

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (Update biased second raw moment estimate)

$\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$ (Compute bias-corrected first moment estimate)

$\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$ (Compute bias-corrected second raw moment estimate)

$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$ (Update parameters)

end while

return θ_t (Resulting parameters)

Machine learning is just Mathematics
