Recap:

minimize $f_0(x)$  Subject to

$$f_i(x) \leq 0, \quad i = 1, 2, \cdots m$$
$$h_i(x) = 0, \quad i = 1, 2, \cdots P$$

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

$$P^* = \min_{x} \max_{\substack{\lambda_i \geq 0 \\ \nu_i}} L(x, \lambda, \nu)$$

$$d^* = \max_{\substack{\lambda_i \geq 0 \\ \nu_i}} \min_{x} L(x, \lambda, \nu)$$

KKT Conditions :

If $P^* = d^*$, then if $x^*$ be Primal optimal and $(\lambda^*, \nu^*)$ be dual optimal, then

$$\nabla_x f_0(x^*) + \sum_{i=1}^{m} \lambda_i^* \nabla f_i(x^*) + \sum_{i=1}^{P} \nu_i^* \nabla h_i(x^*) = 0$$

$$f_i(x^*) \leq 0$$
$$h_i(x^*) = 0$$
$$\lambda_i^* \geq 0$$
$$\lambda_i^* f_i(x^*) = 0$$

If $f_i$ are convex and $h_i$ are affine, and $x$, $\tilde{\lambda}$, $\tilde{\nu}$ are any points that satisfy KKT conditions

then $\tilde{x}$ ~~is~~ and $(\tilde{\lambda}, \tilde{\nu})$ are Primal and dual optimal with zero duality gap.

SVM:

Let $y_i = f(\bar{x}(i))$

minimize $\frac{1}{2} \|w\|_2^2$

subject to $y_i (w \cdot \bar{x}(i) + w_0) \geq 1, \quad i = 1, 2, \cdots n$

Dual:

maximize $\sum_{i=1}^{n} \lambda_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} y_i \lambda_i (\bar{x}(i) \cdot \bar{x}(j))$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad y_j \lambda_j$

subject to $\sum_{i=1}^{n} \lambda_i y_i = 0, \quad \lambda_i \geq 0 \quad i = 1, 2, \cdots n.$

Moreover

$w^* = \sum_{i=1}^{n} \lambda_i^* y_i \bar{x}(i)$

$w_0^* = y_i - w^* \cdot \bar{x}(i)$ for any $i$ s.t.
$\qquad\qquad\qquad\qquad\qquad \lambda_i^* > 0$

# Gradient Descent :

we want to minize $f(\underline{x})$ ( )

Starting from $\underline{x}_0$ ( )

~~Taylor Series Expansion :~~

$f(\underline{x}_0 + h\underline{u})$, $\quad$ $\underline{u}$ unit vector

$\qquad\qquad\qquad\qquad$ $h > 0$ Small Step

Taylor Series expansion :

$\cancel{f(x_0)}$

$$\underbrace{f(\underline{x}_0 + h\underline{u})}_{} \approx f(x_0) + h\,\underline{u}^T \nabla_x f(x_0)$$

$\qquad\qquad\qquad \downarrow$

minimize

$$u_0 = -\frac{\nabla_x f(x_0)}{\|\nabla_x f(x_0)\|}$$

$$\underline{x}_1 = x_0 + h\underline{u}_0$$

Initialize $\underline{x}_0$
$\quad$ Repeat : ~~until~~

$$\underline{x}_{t+1} = \underline{x}_t - r\,\nabla_{\underline{x}} f(\underline{x}_t)$$

Stop if $\|\nabla_x f(\underline{x}_{t+1})\| \leq \eta$

# Stochastic Gradient descent:

Suppose objective function can be written as Sum over training examples:

$$MSE(\underline{w}) = \sum_{i=1}^{n} (y_i - w \cdot \bar{x}(i))^2$$

$$= \sum_{i=1}^{n} g_i(w)$$

Initialize $w_0$:
    Shuffle data $\bar{X}(1), \cdots \bar{X}(n)$
    for $i = 1, 2, \cdots n$:
        $w := w - r \nabla_w g_i(w)$

Advantages:
- Faster than gradient descent
- online learning

Disadvantages
- ~~Might May converge to local minima~~
- performs frequent updates with high variance, hence might keep overshooting.

$$\boxed{8.30 \rightarrow}$$