

Function Approximation Methods-IV

CS771: Introduction to Machine Learning
Purushottam Kar

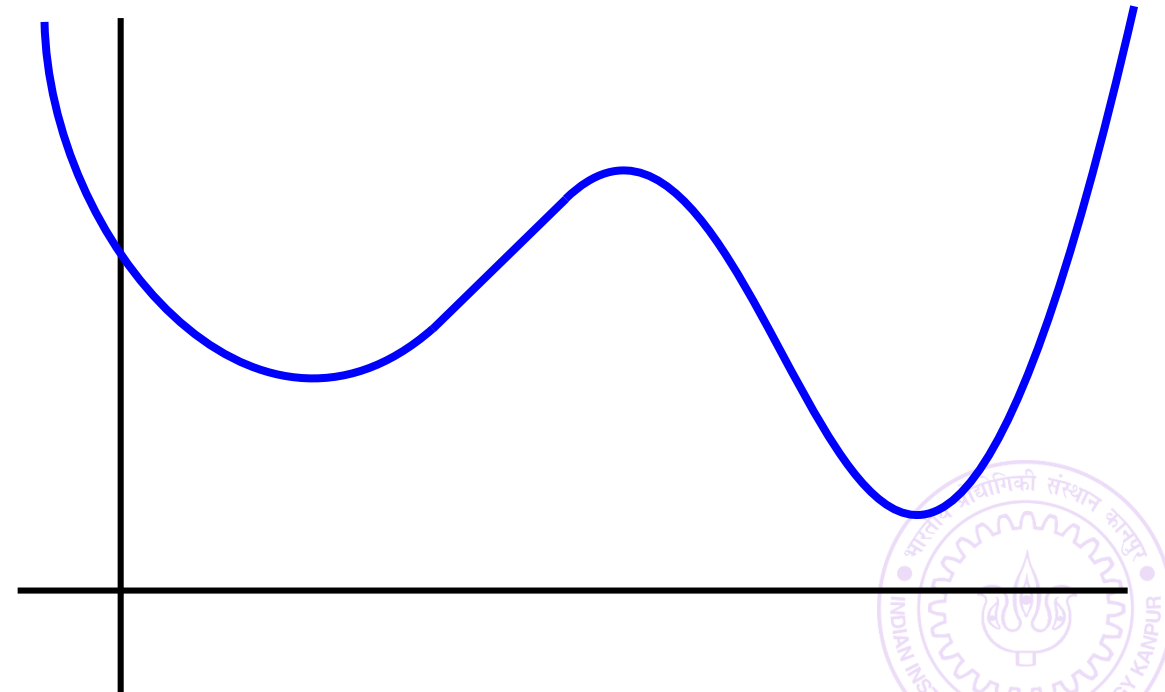


Stochastic Gradient Descent

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$

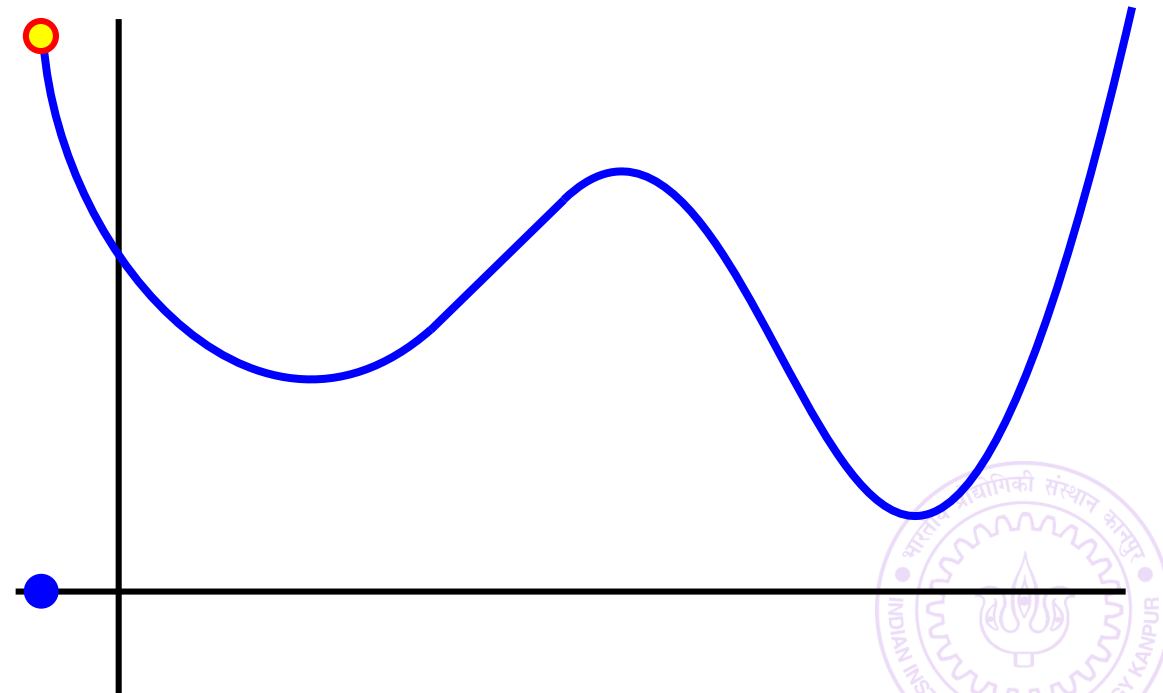
Stochastic Gradient Descent

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$



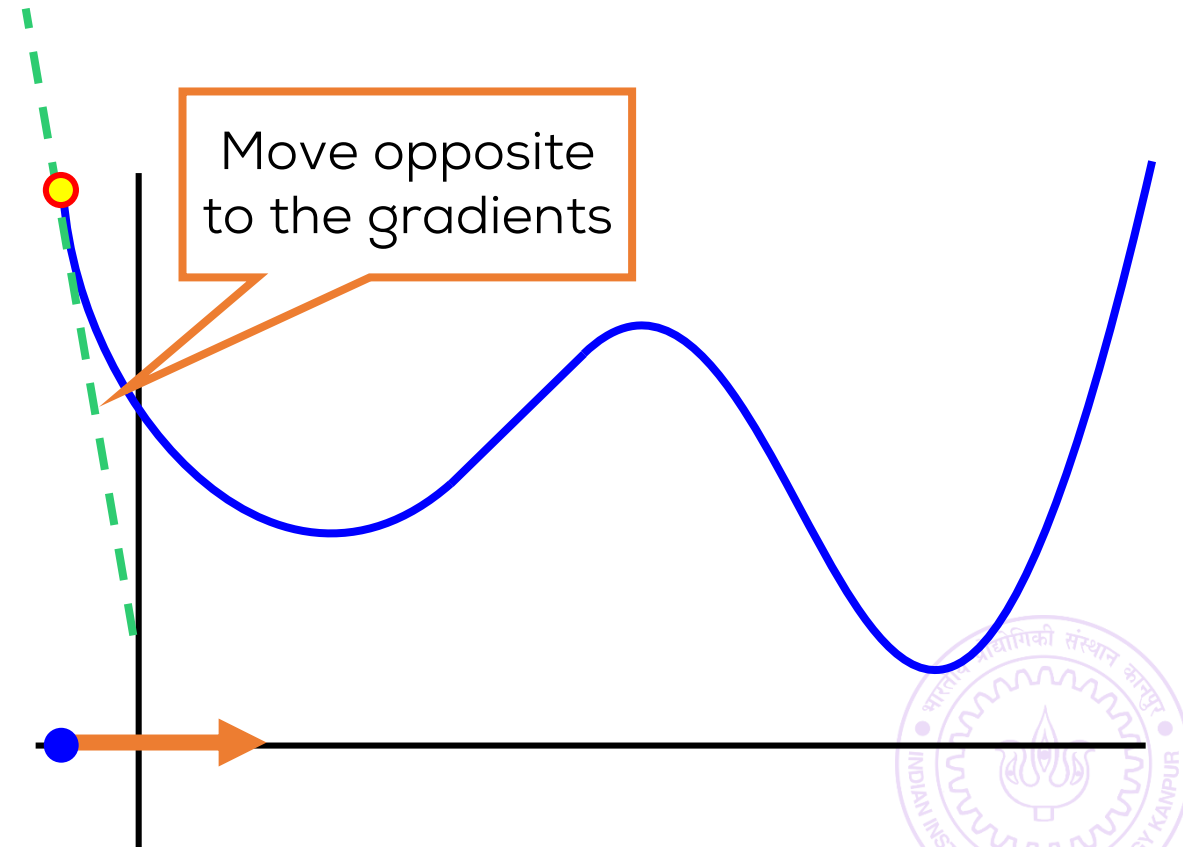
Stochastic Gradient Descent

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$



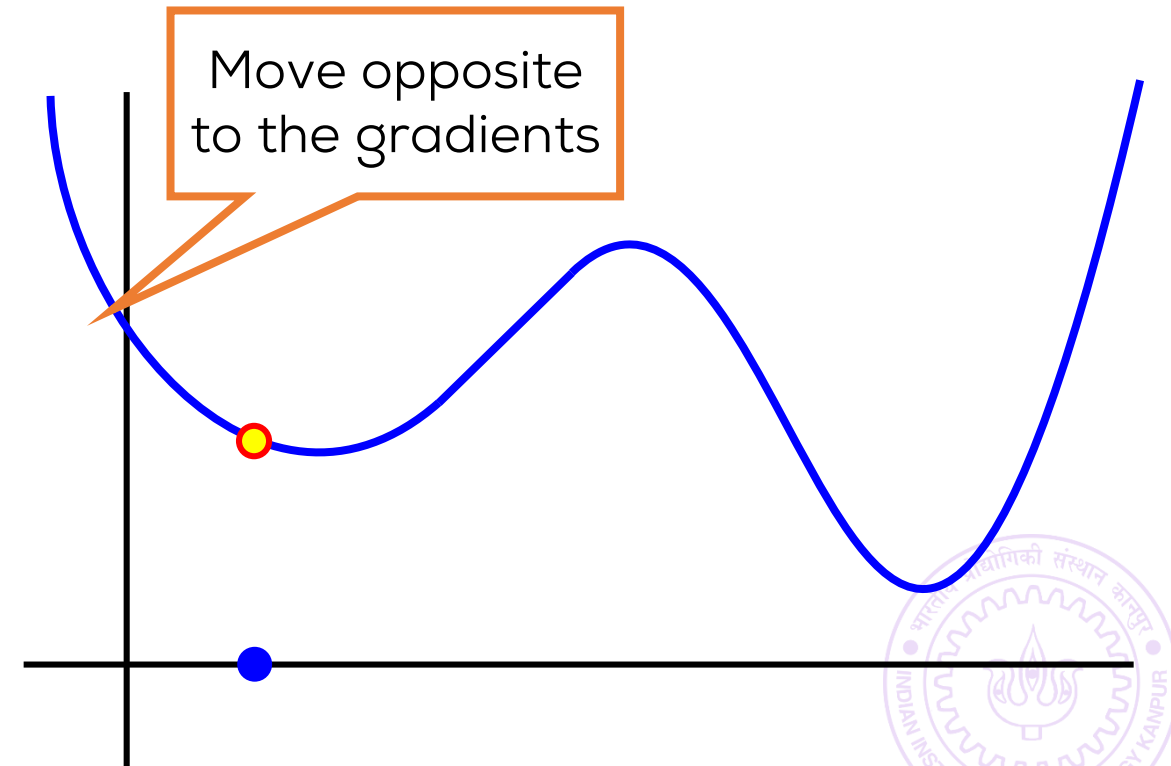
Stochastic Gradient Descent

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$



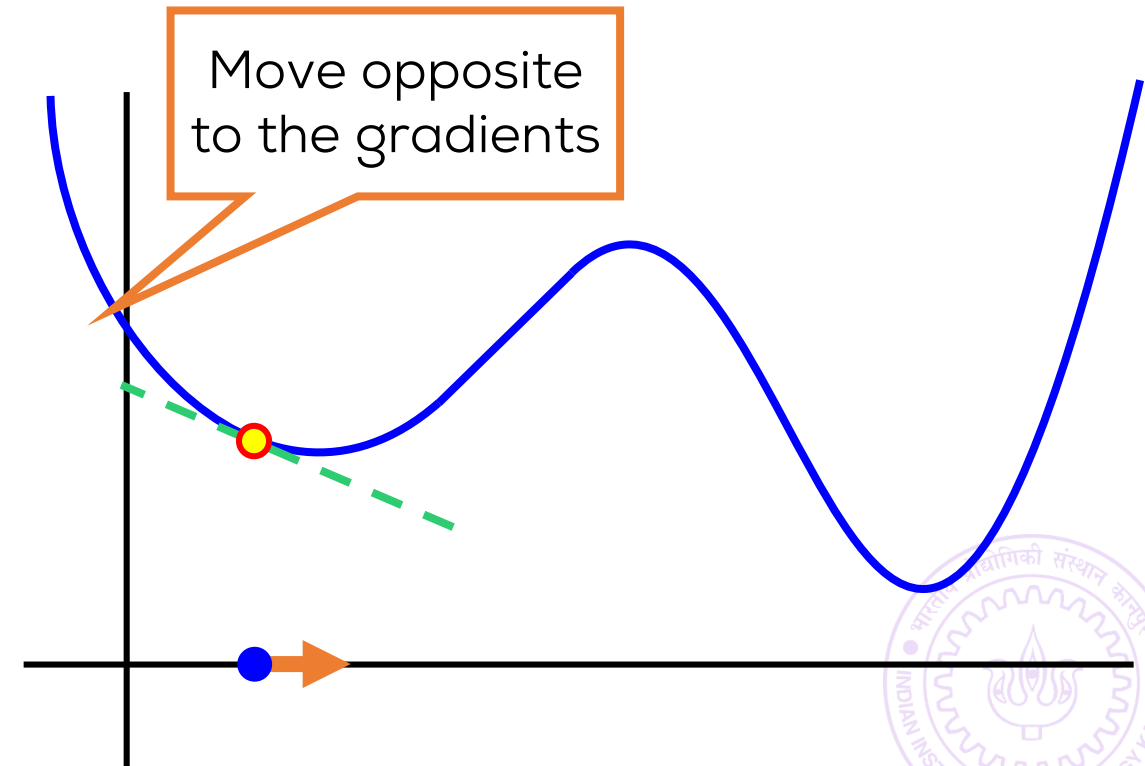
Stochastic Gradient Descent

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$



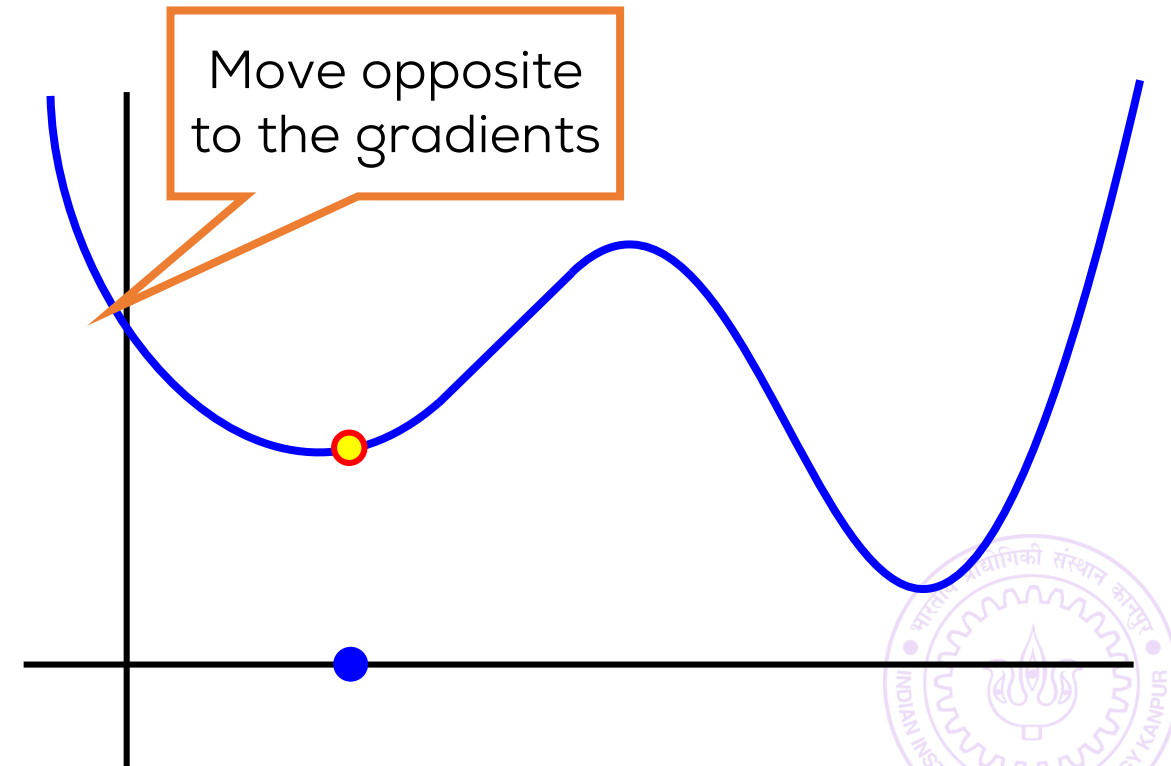
Stochastic Gradient Descent

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$



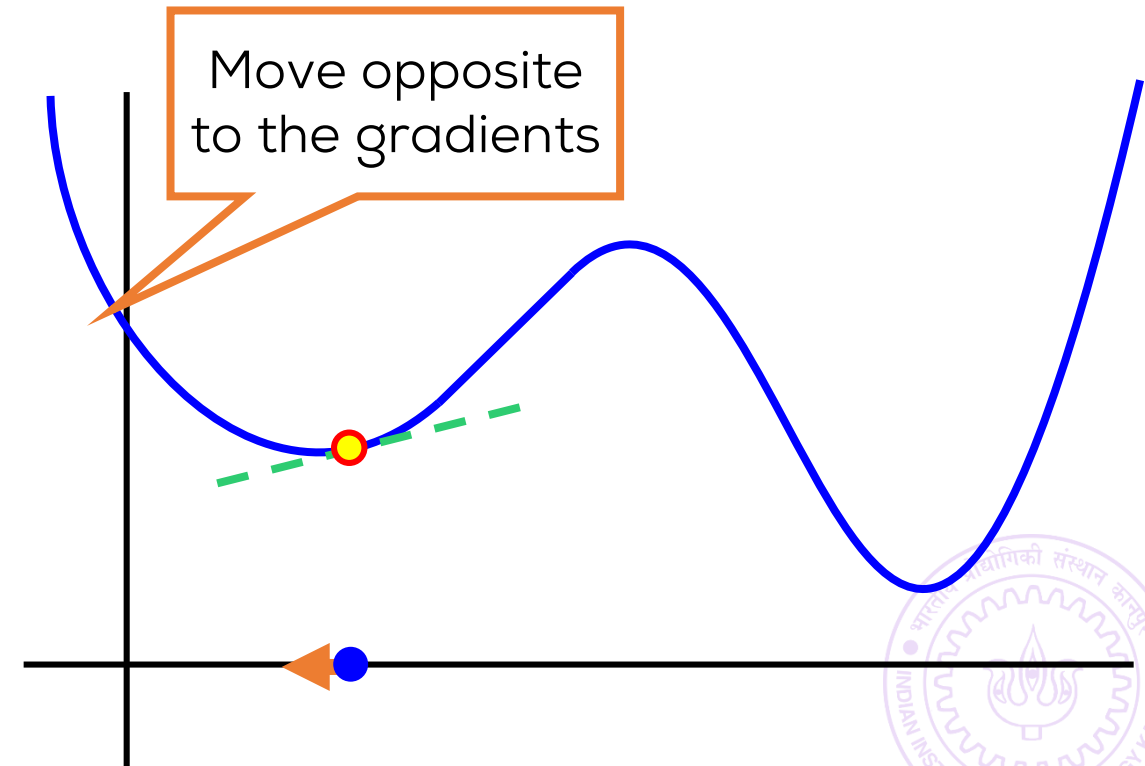
Stochastic Gradient Descent

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$



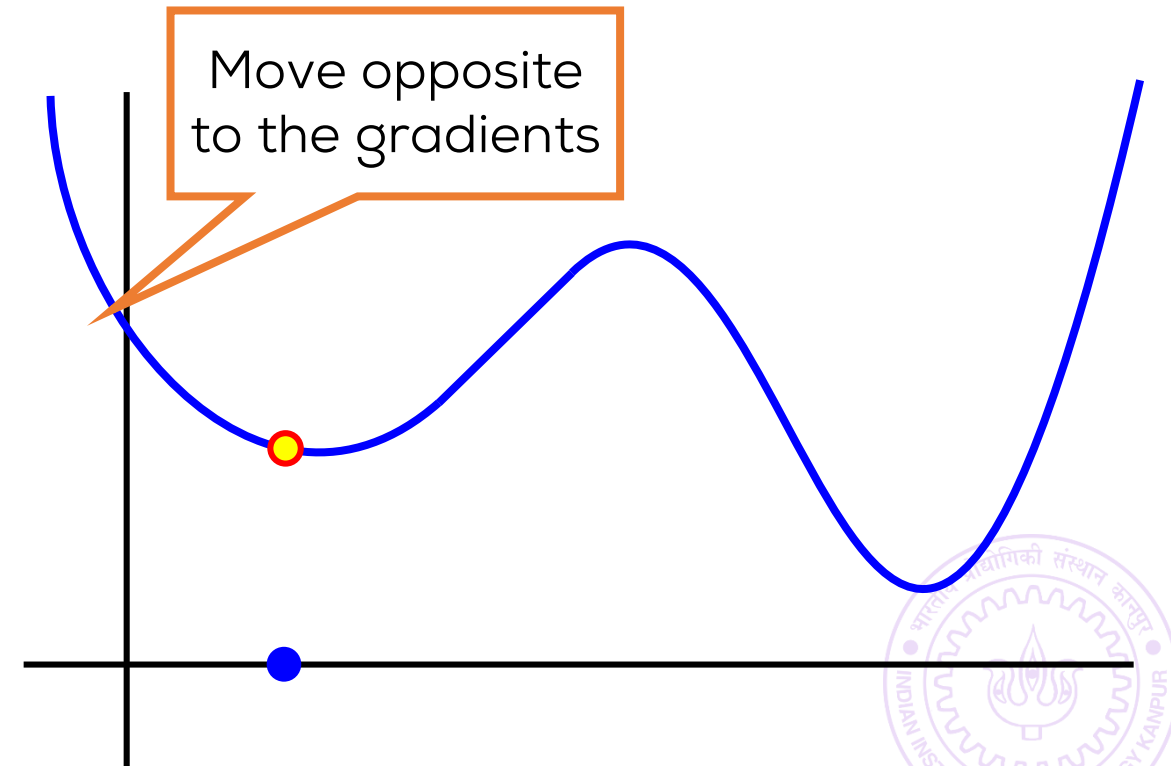
Stochastic Gradient Descent

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$



Stochastic Gradient Descent

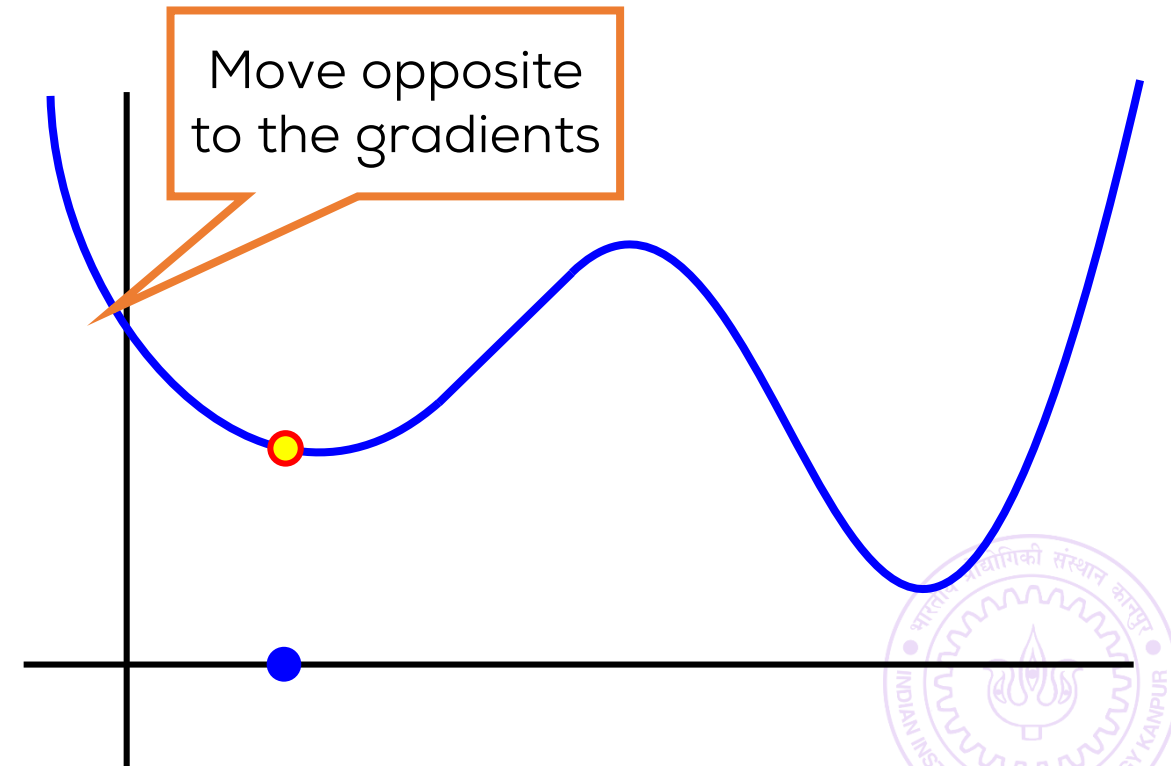
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$



Stochastic Gradient Descent

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$

$$\mathbf{g}^t = \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$



Stochastic Gradient Descent

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$

$$\mathbf{g}^t = \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

Stochastic Gradient Descent

$$f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\mathbf{g}^t = \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

Stochastic Gradient Descent

$O(nd)$ time

$$\nabla f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\mathbf{g}^t = \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

Stochastic Gradient Descent

$O(nd)$ time

$$\nabla f(\mathbf{w}) \approx \nabla f_i(\mathbf{w}) \approx \nabla \ell(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\mathbf{g}^t = \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

Stochastic Gradient Descent

$O(d)$ time!!

$$\nabla f(\mathbf{w}) \approx \nabla f_i(\mathbf{w}) \approx \nabla \ell(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\mathbf{g}^t = \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

Stochastic Gradient Descent

$O(d)$ time!!

$$\nabla f(\mathbf{w}) \approx \nabla f_i(\mathbf{w}) \approx \nabla \ell(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\mathbf{g}^t = \nabla f_i(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

Stochastic Gradient Descent

$O(d)$ time!!

$$\nabla f(\mathbf{w}) \approx \quad \nabla f_i(\mathbf{w}) \approx \quad \nabla \ell(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\mathbf{g}^t = \nabla f_i(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

STOCHASTIC GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Select a unif. random point $I_t \in [n]$
3. Let $\mathbf{g}^t \leftarrow \nabla f_{I_t}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

Stochastic Gradient Descent

$O(d)$ time!!

$$\nabla f(\mathbf{w}) \approx \nabla f_i(\mathbf{w}) \approx \nabla \ell(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\mathbf{g}^t = \nabla f_i(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

$$\mathbb{E}[\mathbf{g}^t | \mathbf{w}^t] = \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

STOCHASTIC GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Select a unif. random point $I_t \in [n]$
3. Let $\mathbf{g}^t \leftarrow \nabla f_{I_t}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

Stochastic Gradient Descent

$O(d)$ time!!

$$\nabla f(\mathbf{w}) \approx \nabla f_i(\mathbf{w}) \approx \nabla \ell(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\mathbf{g}^t = \nabla f_i(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

$$\mathbb{E}[\mathbf{g}^t | \mathbf{w}^t] = \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

Only $O(d)$ time
per iter!

STOCHASTIC GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Select a unif. random point $I_t \in [n]$
3. Let $\mathbf{g}^t \leftarrow \nabla f_{I_t}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

Stochastic Gradient Descent

$O(d)$ time!!

$$\nabla f(\mathbf{w}) \approx \nabla f_i(\mathbf{w}) \approx \nabla \ell(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\mathbf{g}^t = \nabla f_i(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

$$\mathbb{E}[\mathbf{g}^t | \mathbf{w}^t] = \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

Only $O(d)$ time
per iter!

STOCHASTIC GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Select a unif. random point $I_t \in [n]$
3. Let $\mathbf{g}^t \leftarrow \nabla f_{I_t}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

Step length
 $\eta_t = C, \frac{C}{t}, \frac{C}{\sqrt{t}}$

Stochastic Gradient Descent

$O(d)$ time!!

$$\nabla f(\mathbf{w}) \approx \nabla f_i(\mathbf{w}) \approx \nabla \ell(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\mathbf{g}^t = \nabla f_i(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

$$\mathbb{E}[\mathbf{g}^t | \mathbf{w}^t] = \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

Only $O(d)$ time
per iter!

STOCHASTIC GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Select a unif. random point $I_t \in [n]$
3. Let $\mathbf{g}^t \leftarrow \nabla f_{I_t}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

Step length
 $\eta_t = C, \frac{C}{t}, \frac{C}{\sqrt{t}}$

w.h.p. ϵ -optimal
solution in $O\left(\frac{1}{\epsilon^2}\right)$
iterations

Stochastic Gradient Descent

$O(d)$ time!!

$$\nabla f(\mathbf{w}) \approx \nabla f_i(\mathbf{w}) \approx \nabla \ell(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\mathbf{g}^t = \nabla f_i(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

$$\mathbb{E}[\mathbf{g}^t | \mathbf{w}^t] = \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

Only $O(d)$ time per iter!

STOCHASTIC GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Select a unif. random point $I_t \in [n]$
3. Let $\mathbf{g}^t \leftarrow \nabla f_{I_t}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

Step length
 $\eta_t = C, \frac{C}{t}, \frac{C}{\sqrt{t}}$

w.h.p. ϵ -optimal
solution in $O\left(\frac{1}{\epsilon^2}\right)$
iterations

$$\begin{aligned} f(\mathbf{w}^t) + r(\mathbf{w}^t) \\ \leq f(\mathbf{w}^*) + r(\mathbf{w}^*) + \epsilon \end{aligned}$$

Stochastic Gradient Descent

$O(d)$ time!!

$$\nabla f(\mathbf{w}) \approx \nabla f_i(\mathbf{w}) \approx \nabla \ell(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\mathbf{g}^t = \nabla f_i(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

ONLINE GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Select a unif. random point $I_t \in [n]$
3. Let $\mathbf{g}^t \leftarrow \nabla f_{I_t}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

Only $O(d)$ time per iter!

Step length
 $\eta_t = C, \frac{C}{t}, \frac{C}{\sqrt{t}}$

w.h.p. ϵ -optimal
solution in $O\left(\frac{1}{\epsilon^2}\right)$
iterations

$$f(\mathbf{w}^t) + r(\mathbf{w}^t) \leq f(\mathbf{w}^*) + r(\mathbf{w}^*) + \epsilon$$

Stochastic Gradient Descent

$O(d)$ time!!

$$\nabla f(\mathbf{w}) \approx \nabla f_i(\mathbf{w}) \approx \nabla \ell(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\mathbf{g}^t = \nabla f_i(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

ONLINE GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. Let $\mathbf{g}^t \leftarrow \nabla f_{I_t}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

Only $O(d)$ time per iter!

Step length
 $\eta_t = C, \frac{C}{t}, \frac{C}{\sqrt{t}}$

w.h.p. ϵ -optimal
solution in $O\left(\frac{1}{\epsilon^2}\right)$
iterations

$$\begin{aligned} f(\mathbf{w}^t) + r(\mathbf{w}^t) \\ \leq f(\mathbf{w}^*) + r(\mathbf{w}^*) + \epsilon \end{aligned}$$

Stochastic Gradient Descent

$O(d)$ time!!

$$\nabla f(\mathbf{w}) \approx \quad \nabla f_i(\mathbf{w}) \approx \quad \nabla \ell(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\mathbf{g}^t = \nabla f_i(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

ONLINE GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. Let $\mathbf{g}^t \leftarrow \nabla f(\mathbf{w}^t, \mathbf{z}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

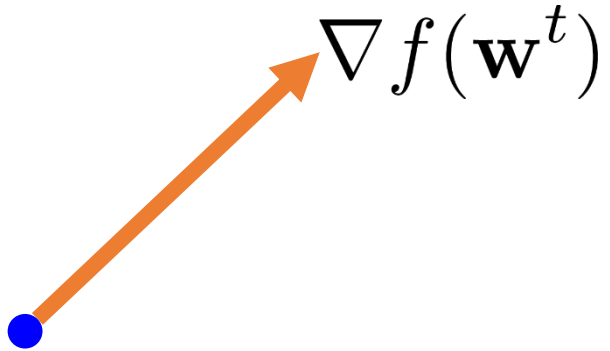
Only $O(d)$ time per iter!

Step length
 $\eta_t = C, \frac{C}{t}, \frac{C}{\sqrt{t}}$

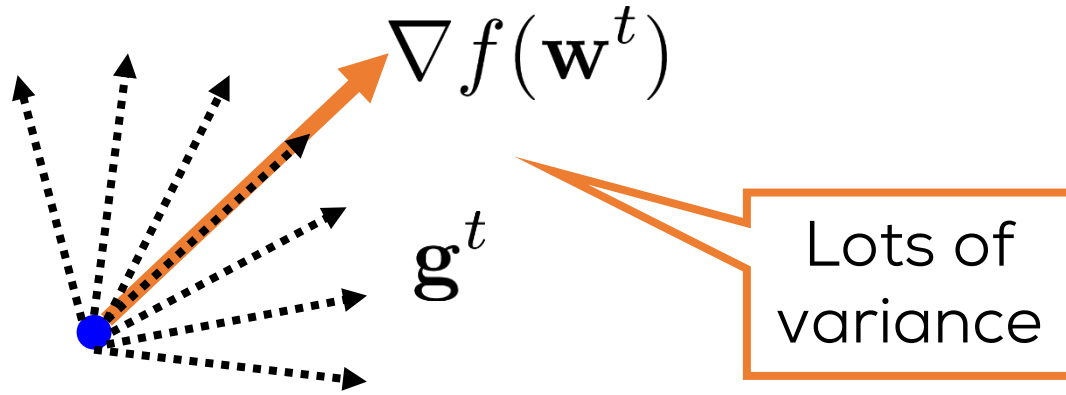
w.h.p. ϵ -optimal
solution in $O\left(\frac{1}{\epsilon^2}\right)$
iterations

$$\begin{aligned} f(\mathbf{w}^t) + r(\mathbf{w}^t) \\ \leq f(\mathbf{w}^*) + r(\mathbf{w}^*) + \epsilon \end{aligned}$$

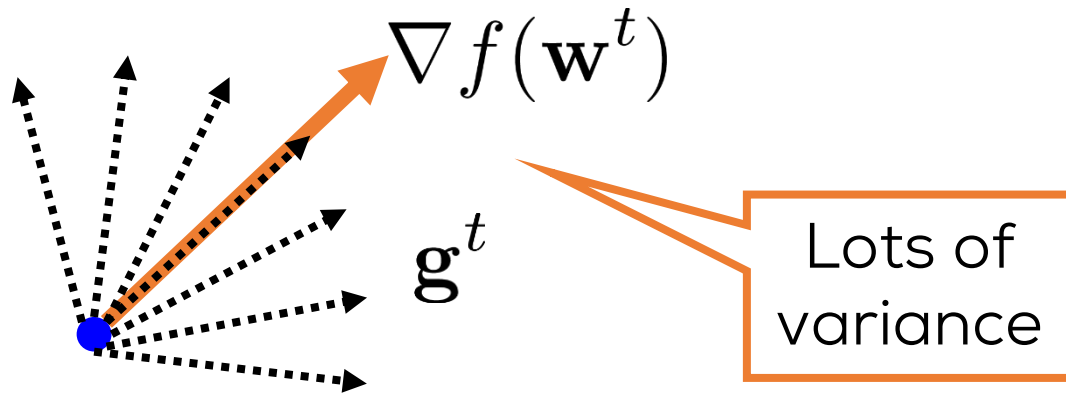
Mini-batch Gradient Descent



Mini-batch Gradient Descent



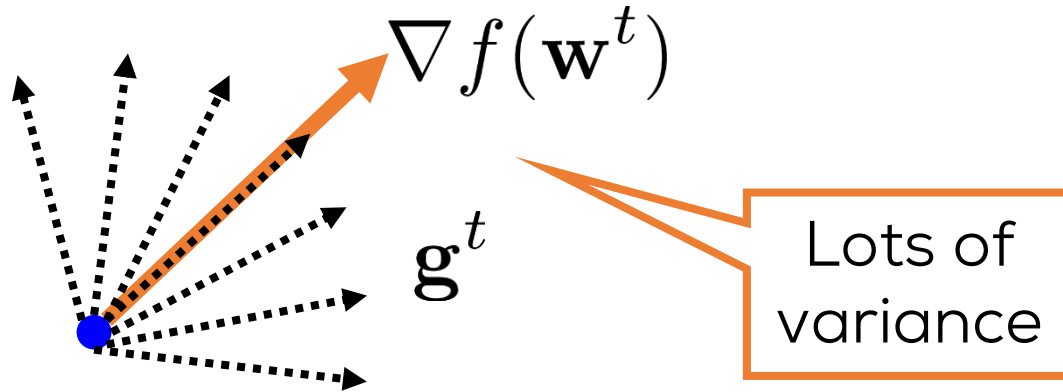
Mini-batch Gradient Descent



MINI-BATCH SGD

1. Initialize \mathbf{w}^0
2. Select B unif. random points $\{I_{t,1}, \dots, I_{t,B}\} \in [n]$
3. Let $\mathbf{g}^t \leftarrow \frac{1}{B} \sum_{i=1}^B \nabla f_{I_{t,i}}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

Mini-batch Gradient Descent



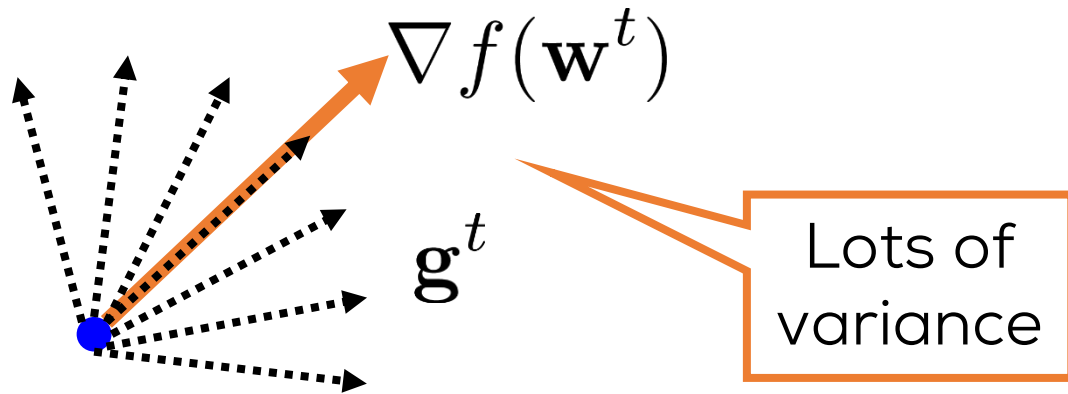
MINI-BATCH SGD

1. Initialize \mathbf{w}^0
2. Select B unif. random points $\{I_{t,1}, \dots, I_{t,B}\} \in [n]$
3. Let $\mathbf{g}^t \leftarrow \frac{1}{B} \sum_{i=1}^B \nabla f_{I_{t,i}}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

Without replacement?



Mini-batch Gradient Descent



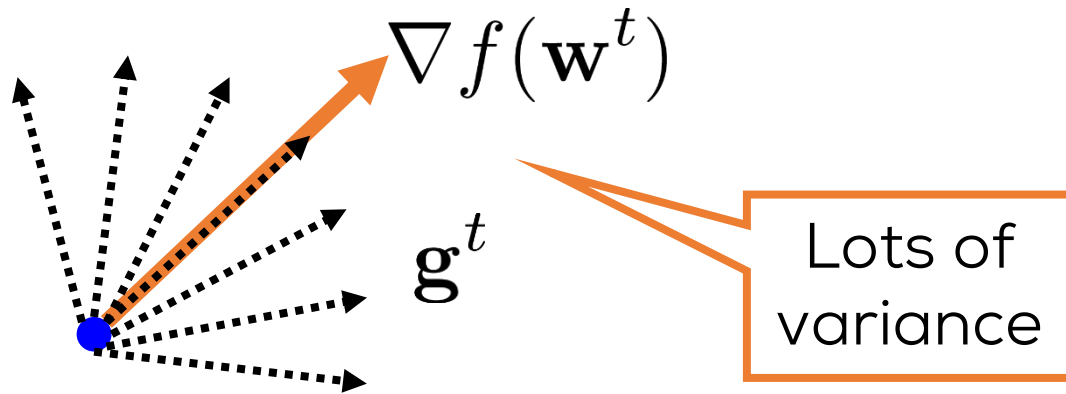
MINI-BATCH SGD

1. Initialize \mathbf{w}^0
2. Select B unif. random points $\{I_{t,1}, \dots, I_{t,B}\} \in [n]$
3. Let $\mathbf{g}^t \leftarrow \frac{1}{B} \sum_{i=1}^B \nabla f_{I_{t,i}}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

Without replacement?

$O(B \cdot d)$ time per iter!

Mini-batch Gradient Descent



MINI-BATCH SGD

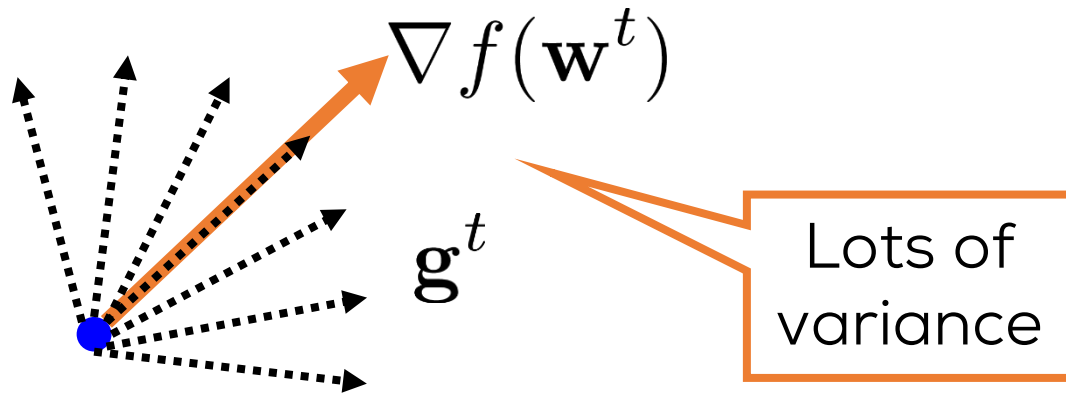
1. Initialize \mathbf{w}^0
2. Select B unif. random points $\{I_{t,1}, \dots, I_{t,B}\} \in [n]$
3. Let $\mathbf{g}^t \leftarrow \frac{1}{B} \sum_{i=1}^B \nabla f_{I_{t,i}}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

$$\mathbb{E}[\mathbf{g}^t | \mathbf{w}^t] = \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

Without replacement?

$O(B \cdot d)$ time per iter!

Mini-batch Gradient Descent



MINI-BATCH SGD

1. Initialize \mathbf{w}^0
2. Select B unif. random points $\{I_{t,1}, \dots, I_{t,B}\} \in [n]$
3. Let $\mathbf{g}^t \leftarrow \frac{1}{B} \sum_{i=1}^B \nabla f_{I_{t,i}}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

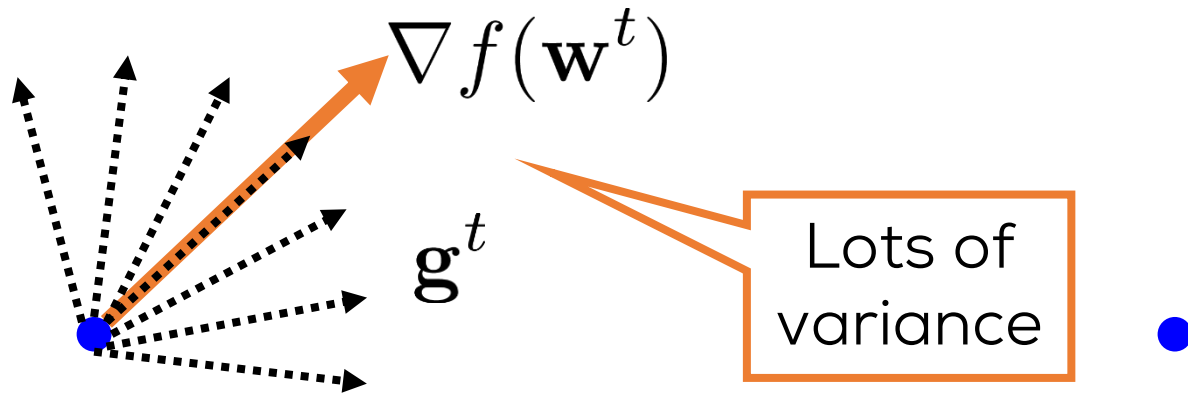
$$\mathbb{E}[\mathbf{g}^t | \mathbf{w}^t] = \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

Without replacement?

Usually $B \ll n$

$O(B \cdot d)$ time per iter!

Mini-batch Gradient Descent



MINI-BATCH SGD

1. Initialize \mathbf{w}^0
2. Select B unif. random points $\{I_{t,1}, \dots, I_{t,B}\} \in [n]$
3. Let $\mathbf{g}^t \leftarrow \frac{1}{B} \sum_{i=1}^B \nabla f_{I_{t,i}}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

$$\mathbb{E}[\mathbf{g}^t | \mathbf{w}^t] = \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

Without
replacement?

Usually $B \ll n$

$O(B \cdot d)$ time
per iter!

Mini-batch Gradient Descent



MINI-BATCH SGD

1. Initialize \mathbf{w}^0
2. Select B unif. random points $\{I_{t,1}, \dots, I_{t,B}\} \in [n]$
3. Let $\mathbf{g}^t \leftarrow \frac{1}{B} \sum_{i=1}^B \nabla f_{I_{t,i}}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

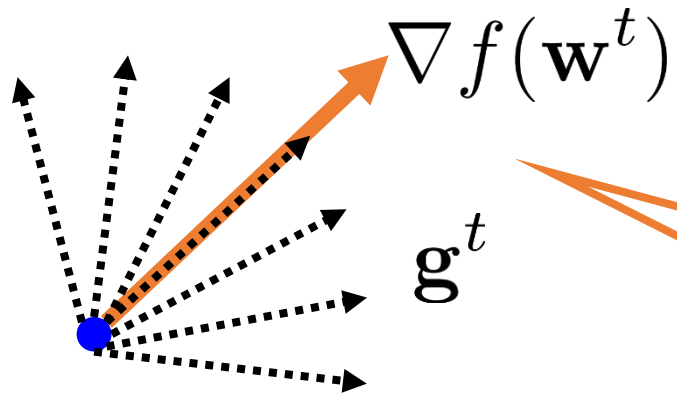
$$\mathbb{E}[\mathbf{g}^t | \mathbf{w}^t] = \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

Without
replacement?

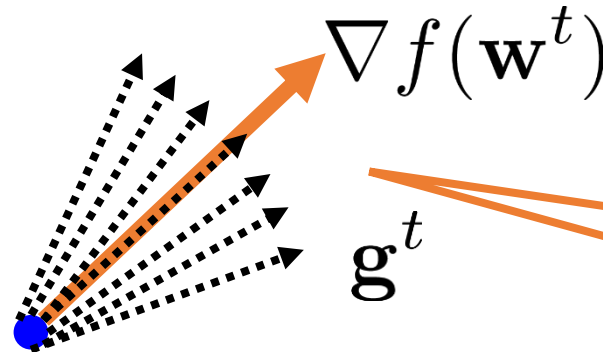
Usually $B \ll n$

$O(B \cdot d)$ time
per iter!

Mini-batch Gradient Descent



Lots of
variance



No variance if
all points chosen
i.e with $B = n$

Less variance
with mini-batch

MINI-BATCH SGD

1. Initialize \mathbf{w}^0
2. Select B unif. random points $\{I_{t,1}, \dots, I_{t,B}\} \in [n]$
3. Let $\mathbf{g}^t \leftarrow \frac{1}{B} \sum_{i=1}^B \nabla f_{I_{t,i}}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat until convergence

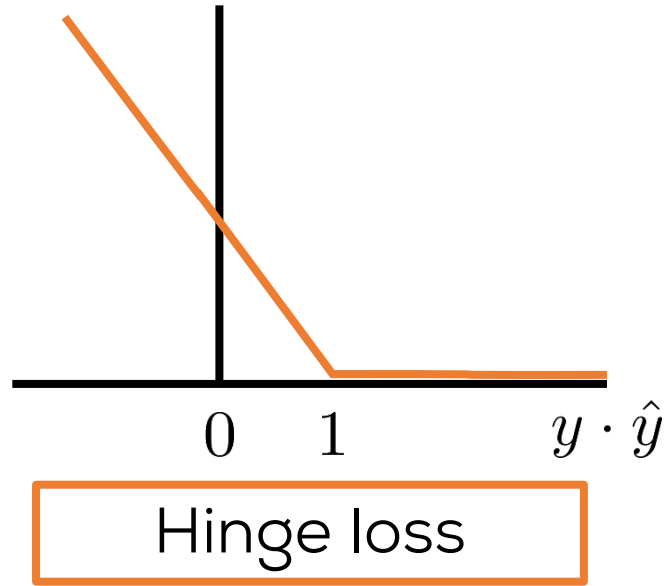
$$\mathbb{E}[\mathbf{g}^t | \mathbf{w}^t] = \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$$

Without
replacement?

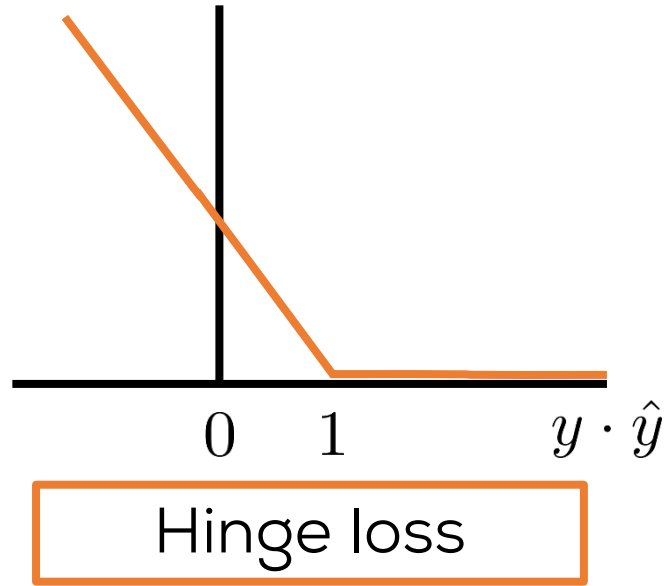
Usually $B \ll n$

$O(B \cdot d)$ time
per iter!

App: The Perceptron Rule via OGD



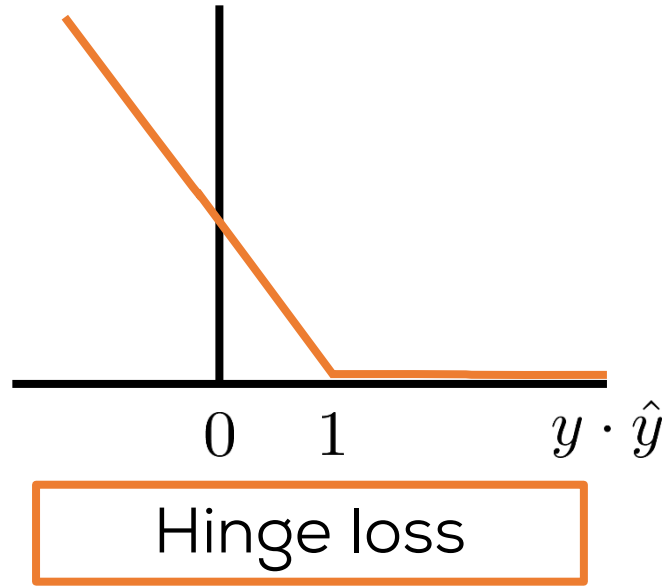
App: The Perceptron Rule via OGD



ONLINE GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. Let $\mathbf{g}^t \leftarrow \nabla f(\mathbf{w}^t, \mathbf{z}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat

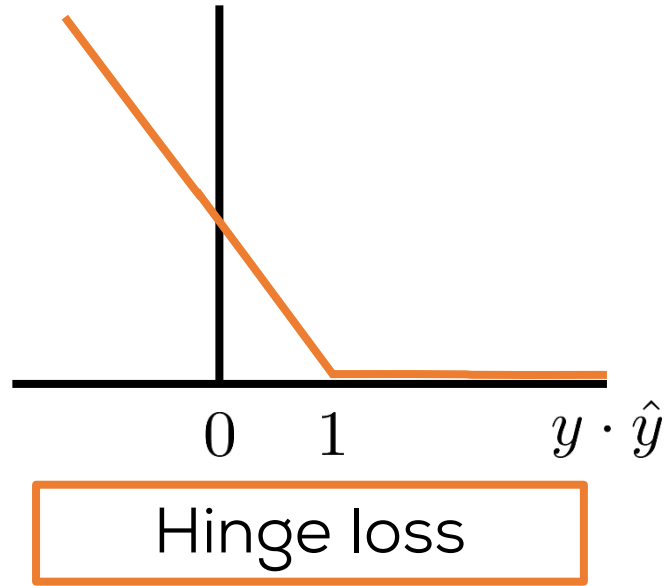
App: The Perceptron Rule via OGD



ONLINE GRADIENT DESCENT $y^t \in \{-1, +1\}$

1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. Let $\mathbf{g}^t \leftarrow \nabla f(\mathbf{w}^t, \mathbf{z}^t) + \nabla r(\mathbf{w}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat

App: The Perceptron Rule via OGD

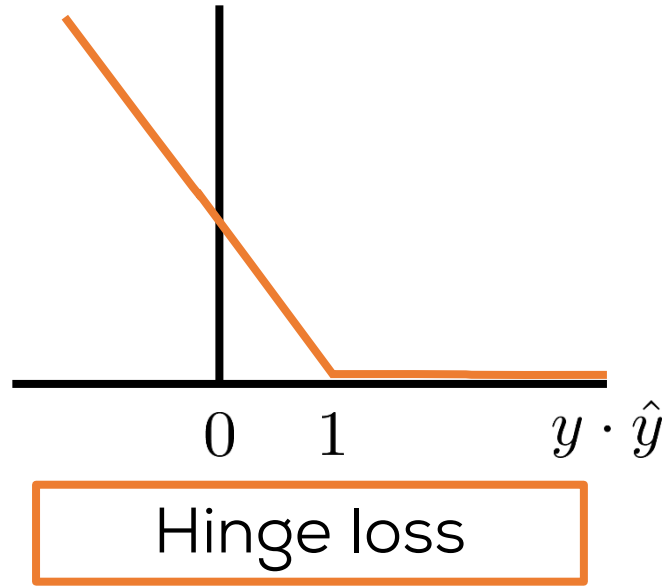


ONLINE GRADIENT DESCENT $y^t \in \{-1, +1\}$

1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. Let $\mathbf{g}^t \leftarrow \nabla f(\mathbf{w}^t, \mathbf{z}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat

Forget reg.
for now

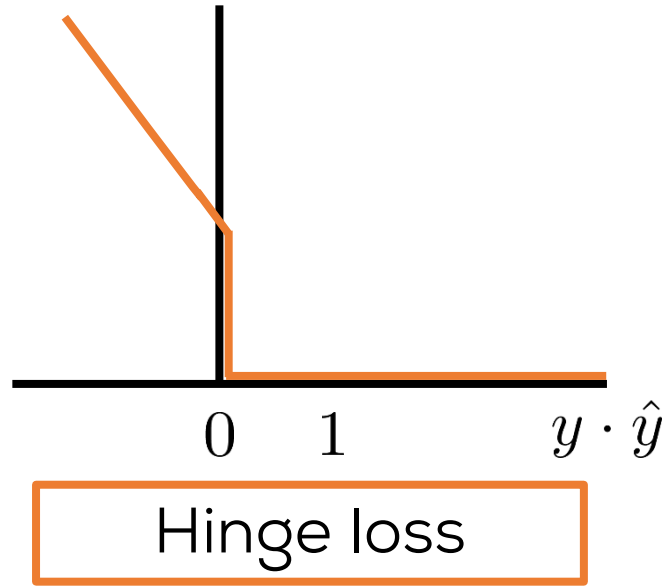
App: The Perceptron Rule via OGD



ONLINE GRADIENT DESCENT $y^t \in \{-1, +1\}$

1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. Let $\mathbf{g}^t \leftarrow \nabla f(\mathbf{w}^t, \mathbf{z}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat

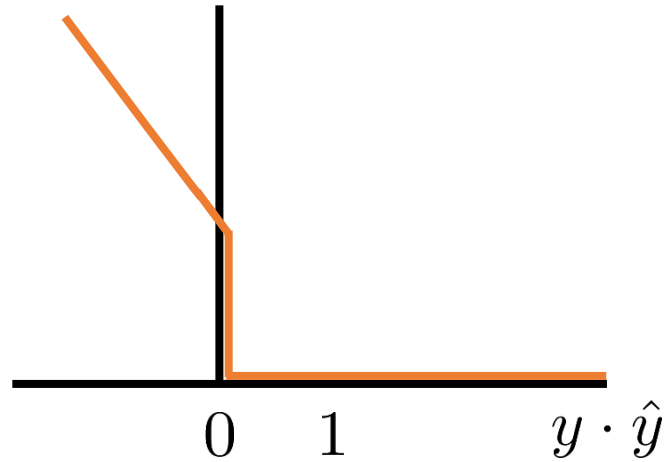
App: The Perceptron Rule via OGD



ONLINE GRADIENT DESCENT $y^t \in \{-1, +1\}$

1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. Let $\mathbf{g}^t \leftarrow \nabla f(\mathbf{w}^t, \mathbf{z}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat

App: The Perceptron Rule via OGD

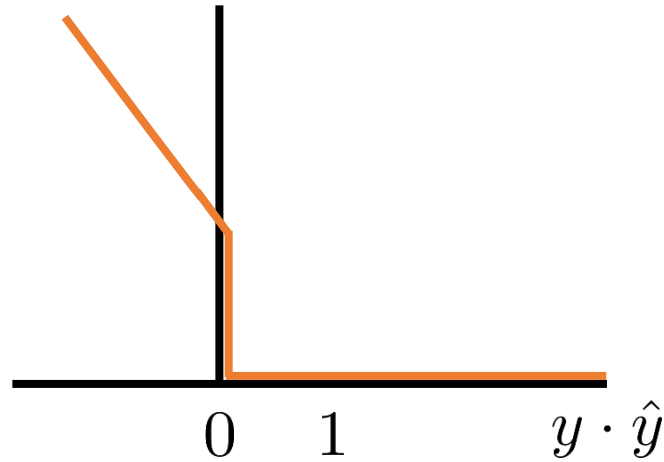


Truncated Hinge loss

ONLINE GRADIENT DESCENT $y^t \in \{-1, +1\}$

1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. Let $\mathbf{g}^t \leftarrow \nabla f(\mathbf{w}^t, \mathbf{z}^t)$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat

App: The Perceptron Rule via OGD

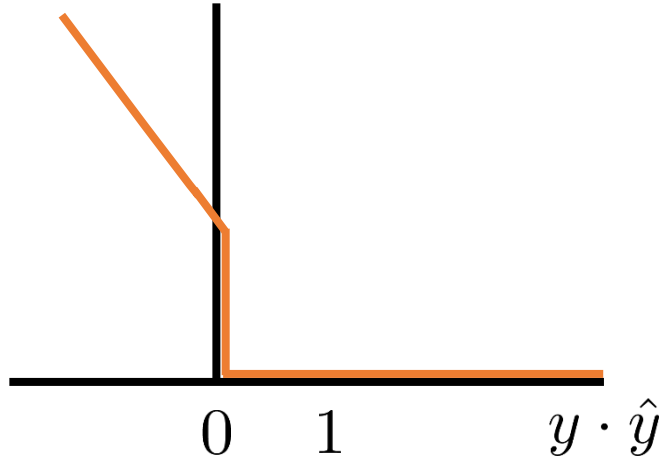


Truncated Hinge loss

ONLINE GRADIENT DESCENT $y^t \in \{-1, +1\}$

1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. If $y^t \langle \mathbf{w}^t, \mathbf{x}^t \rangle < 0$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
5. Repeat

App: The Perceptron Rule via OGD

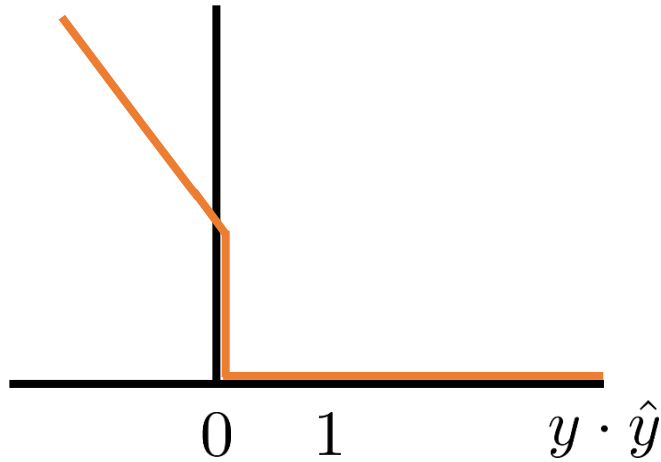


Truncated Hinge loss

ONLINE GRADIENT DESCENT $y^t \in \{-1, +1\}$

1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. If $y^t \langle \mathbf{w}^t, \mathbf{x}^t \rangle < 0$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \eta_t \cdot y^t \cdot \mathbf{x}^t$
5. Repeat

App: The Perceptron Rule via OGD



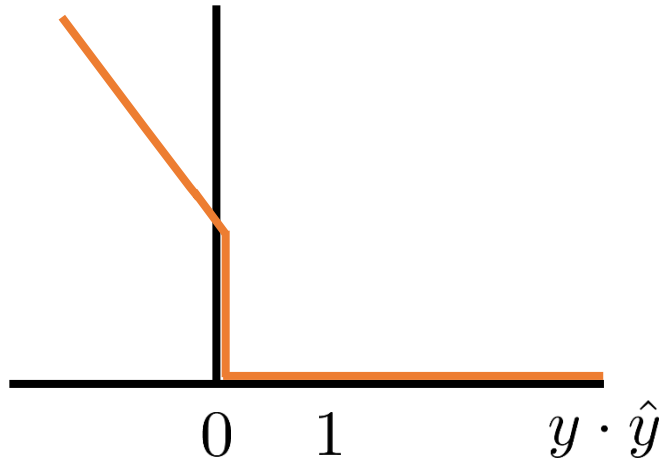
Truncated Hinge loss

Make updates only
when making a mistake

ONLINE GRADIENT DESCENT $y^t \in \{-1, +1\}$

1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. If $y^t \langle \mathbf{w}^t, \mathbf{x}^t \rangle < 0$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \eta_t \cdot y^t \cdot \mathbf{x}^t$
5. Repeat

App: The Perceptron Rule via OGD



Truncated Hinge loss

ONLINE GRADIENT DESCENT $y^t \in \{-1, +1\}$

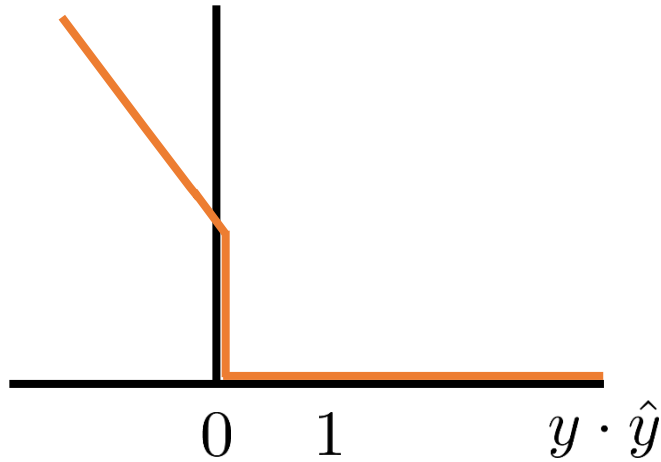
1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. If $y^t \langle \mathbf{w}^t, \mathbf{x}^t \rangle < 0$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \eta_t \cdot y^t \cdot \mathbf{x}^t$
5. Repeat

Make updates only when making a mistake

Push \mathbf{w}^t "away"/"toward" \mathbf{x}^t depending on mistake



App: The Perceptron Rule via OGD



Truncated Hinge loss

PERCEPTRON

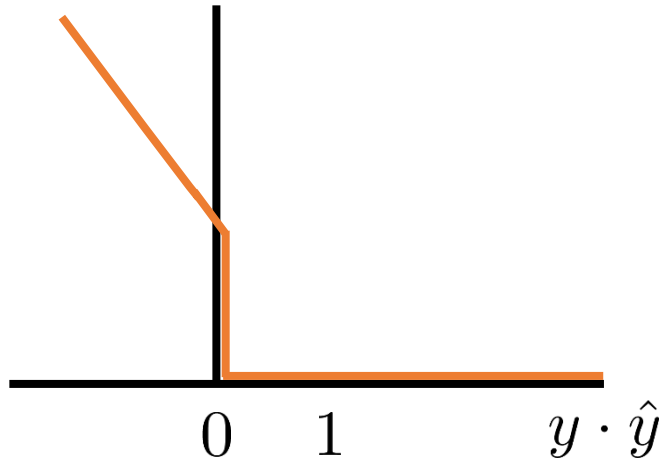
$$y^t \in \{-1, +1\}$$

1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. If $y^t \langle \mathbf{w}^t, \mathbf{x}^t \rangle < 0$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \eta_t \cdot y^t \cdot \mathbf{x}^t$
5. Repeat

Make updates only when making a mistake

Push \mathbf{w}^t "away"/"toward" \mathbf{x}^t depending on mistake

App: The Perceptron Rule via OGD



Truncated Hinge loss

PERCEPTRON

$$y^t \in \{-1, +1\}$$

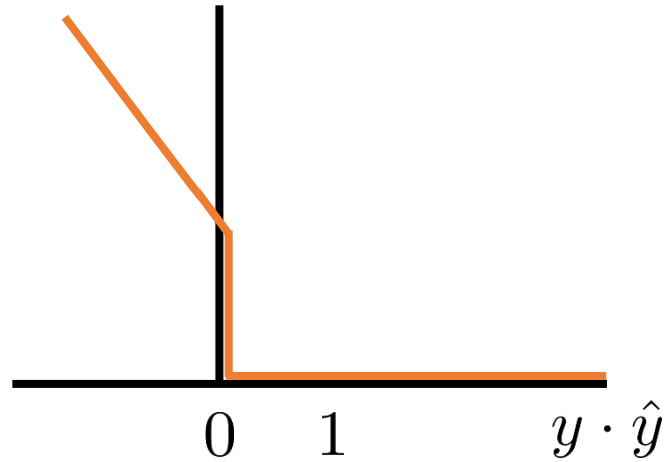
1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. If $y^t \langle \mathbf{w}^t, \mathbf{x}^t \rangle < 0$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \eta_t \cdot y^t \cdot \mathbf{x}^t$
5. Repeat

Make updates only when making a mistake

Push \mathbf{w}^t "away"/"toward" \mathbf{x}^t depending on mistake

Perceptron Logistic Regression??

App: The Perceptron Rule via OGD



Truncated Hinge loss

Make updates only
when making a mistake

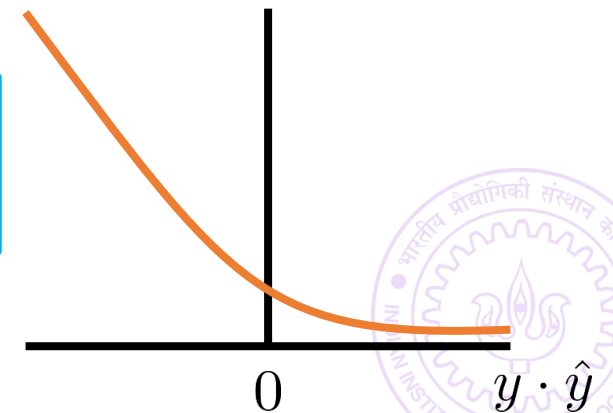
Push \mathbf{w}^t "away"/"toward" \mathbf{x}^t
depending on mistake

PERCEPTRON

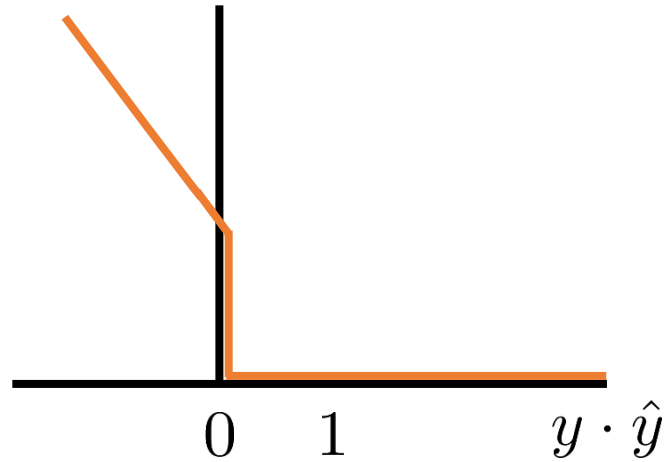
1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. If $y^t \langle \mathbf{w}^t, \mathbf{x}^t \rangle < 0$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \eta_t \cdot y^t \cdot \mathbf{x}^t$
5. Repeat

$$y^t \in \{-1, +1\}$$

Perceptron Logistic
Regression??



App: The Perceptron Rule via OGD



Truncated Hinge loss

Make updates only
when making a mistake

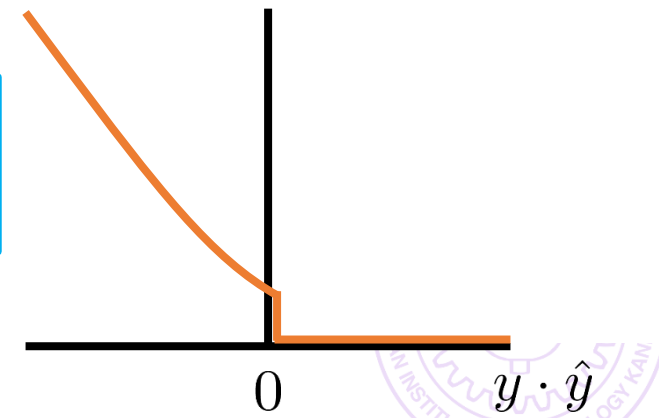
Push \mathbf{w}^t "away"/"toward" \mathbf{x}^t
depending on mistake

PERCEPTRON

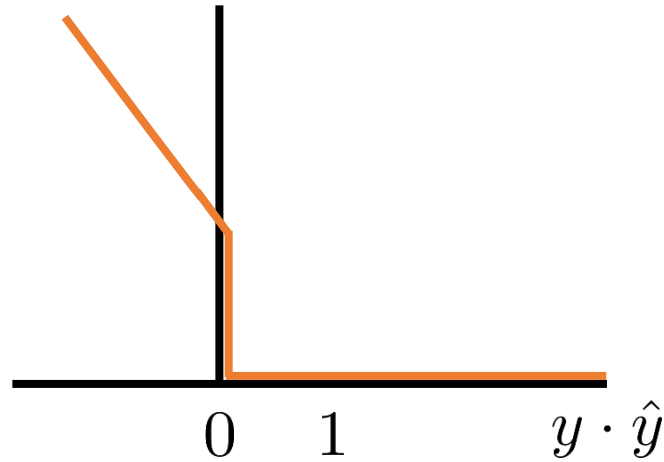
1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. If $y^t \langle \mathbf{w}^t, \mathbf{x}^t \rangle < 0$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \eta_t \cdot y^t \cdot \mathbf{x}^t$
5. Repeat

$$y^t \in \{-1, +1\}$$

Perceptron Logistic
Regression??



App: The Perceptron Rule via OGD



Truncated Hinge loss

Make updates only when making a mistake

Push \mathbf{w}^t "away"/"toward" \mathbf{x}^t depending on mistake

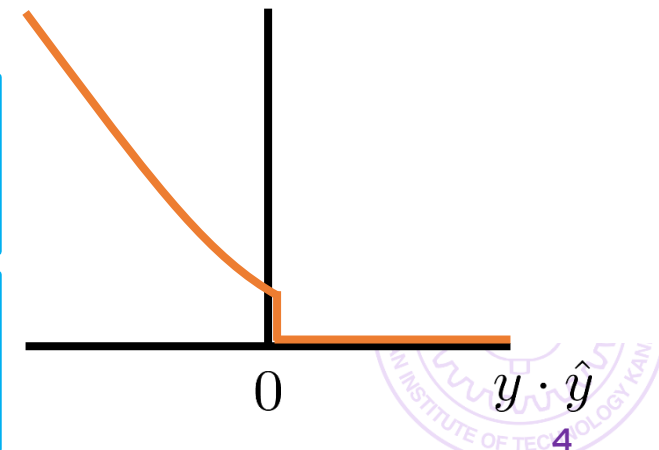
PERCEPTRON

1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. If $y^t \langle \mathbf{w}^t, \mathbf{x}^t \rangle < 0$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \eta_t \cdot y^t \cdot \mathbf{x}^t$
5. Repeat

$$y^t \in \{-1, +1\}$$

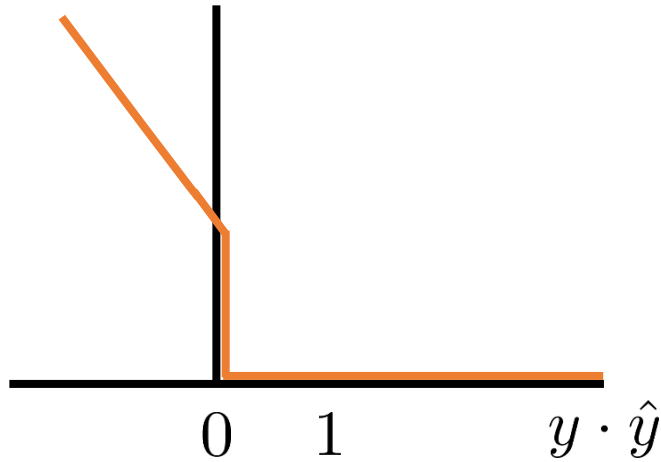
Perceptron Logistic Regression??

Perceptron Crammer Singer for multi-classification?



App: The Perceptron Rule via OGD

Can do stochastic perceptron too



Truncated Hinge loss

PERCEPTRON

$y^t \in \{-1, +1\}$

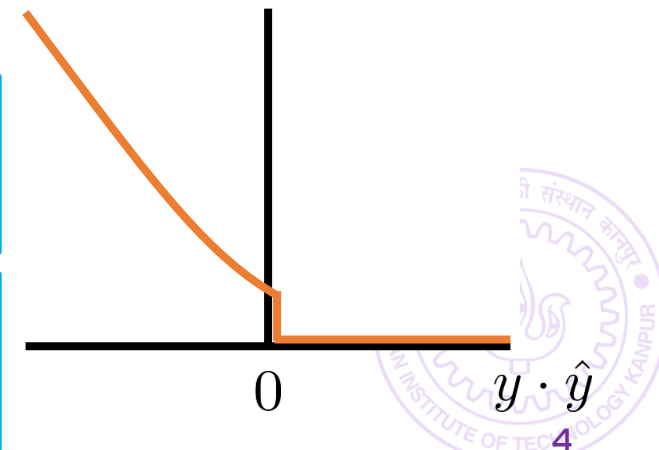
1. Initialize \mathbf{w}^0
2. Receive a data point $\mathbf{z}^t = (y^t, \mathbf{x}^t)$
3. If $y^t \langle \mathbf{w}^t, \mathbf{x}^t \rangle < 0$
4. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \eta_t \cdot y^t \cdot \mathbf{x}^t$
5. Repeat

Make updates only when making a mistake

Push \mathbf{w}^t "away"/"toward" \mathbf{x}^t depending on mistake

Perceptron Logistic Regression??

Perceptron Crammer Singer for multi-classification?



Coordinate Descent

September 1, 2017



Coordinate Descent

- Similar to gradient descent except update one variable at a time
- Notation: $\nabla_i f(\mathbf{w}) = [\nabla f(\mathbf{w})]_i$ i.e. the i -th directional derivative

Coordinate Descent

- Similar to gradient descent except update one variable at a time
- Notation: $\nabla_i f(\mathbf{w}) = [\nabla f(\mathbf{w})]_i$ i.e. the i -th directional derivative

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$

Coordinate Descent

- Similar to gradient descent except update one variable at a time
- Notation: $\nabla_i f(\mathbf{w}) = [\nabla f(\mathbf{w})]_i$ i.e. the i -th directional derivative

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$

COORDINATE DESCENT

1. Initialize $\mathbf{w}^0 \in \mathbb{R}^d$
2. Select a coordinate $j_t \in [d]$
3. Let $\mathbf{g}_{j_t}^t \leftarrow \nabla_{j_t} f(\mathbf{w}^t) + \nabla_{j_t} r(\mathbf{w}^t) \in \mathbb{R}$
4. Update $\mathbf{w}_{j_t}^{t+1} \leftarrow \mathbf{w}_{j_t}^t - \eta_t \cdot \mathbf{g}_{j_t}^t$
5. Preserve other coord. $\mathbf{w}_k^{t+1} \leftarrow \mathbf{w}_k^t, k \neq j_t$
6. Repeat until convergence

Coordinate Descent

- Similar to gradient descent except update one variable at a time
- Notation: $\nabla_i f(\mathbf{w}) = [\nabla f(\mathbf{w})]_i$ i.e. the i -th directional derivative

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$

COORDINATE DESCENT

1. Initialize $\mathbf{w}^0 \in \mathbb{R}^d$
2. Select a coordinate $j_t \in [d]$
3. Let $\mathbf{g}_{j_t}^t \leftarrow \nabla_{j_t} f(\mathbf{w}^t) + \nabla_{j_t} r(\mathbf{w}^t) \in \mathbb{R}$
4. Update $\mathbf{w}_{j_t}^{t+1} \leftarrow \mathbf{w}_{j_t}^t - \eta_t \cdot \mathbf{g}_{j_t}^t$
5. Preserve other coord. $\mathbf{w}_k^{t+1} \leftarrow \mathbf{w}_k^t, k \neq j_t$
6. Repeat until convergence

Random (Stochastic CD), Cyclic
 $1, 2, \dots, d, 1, 2, 3 \dots d, \dots$

Coordinate Descent

- Similar to gradient descent except update one variable at a time
- Notation: $\nabla_i f(\mathbf{w}) = [\nabla f(\mathbf{w})]_i$ i.e. the i -th directional derivative

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$

COORDINATE DESCENT

1. Initialize $\mathbf{w}^0 \in \mathbb{R}^d$
2. Select a coordinate $j_t \in [d]$
3. Let $\mathbf{g}_{j_t}^t \leftarrow \nabla_{j_t} f(\mathbf{w}^t) + \nabla_{j_t} r(\mathbf{w}^t) \in \mathbb{R}$
4. Update $\mathbf{w}_{j_t}^{t+1} \leftarrow \mathbf{w}_{j_t}^t - \eta_t \cdot \mathbf{g}_{j_t}^t$
5. Preserve other coord. $\mathbf{w}_k^{t+1} \leftarrow \mathbf{w}_k^t, k \neq j_t$
6. Repeat until convergence

Random (Stochastic CD), Cyclic
1,2, ..., d, 1,2,3 ... d, ...

Only $O(n)$ time
per iter!

Coordinate Descent

- Similar to gradient descent except update one variable at a time
- Notation: $\nabla_i f(\mathbf{w}) = [\nabla f(\mathbf{w})]_i$ i.e. the i -th directional derivative

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$

Useful when
 $d \gg n$

Random (Stochastic
CD), Cyclic
 $1, 2, \dots, d, 1, 2, 3 \dots d, \dots$

COORDINATE DESCENT

1. Initialize $\mathbf{w}^0 \in \mathbb{R}^d$
2. Select a coordinate $j_t \in [d]$
3. Let $\mathbf{g}_{j_t}^t \leftarrow \nabla_{j_t} f(\mathbf{w}^t) + \nabla_{j_t} r(\mathbf{w}^t) \in \mathbb{R}$
4. Update $\mathbf{w}_{j_t}^{t+1} \leftarrow \mathbf{w}_{j_t}^t - \eta_t \cdot \mathbf{g}_{j_t}^t$
5. Preserve other coord. $\mathbf{w}_k^{t+1} \leftarrow \mathbf{w}_k^t, k \neq j_t$
6. Repeat until convergence

Only $O(n)$ time
per iter!

Coordinate Descent

- Similar to gradient descent except update one variable at a time
- Notation: $\nabla_i f(\mathbf{w}) = [\nabla f(\mathbf{w})]_i$ i.e. the i -th directional derivative

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$

Useful when
 $d \gg n$

Random (Stochastic
CD), Cyclic
 $1, 2, \dots, d, 1, 2, 3 \dots d, \dots$

COORDINATE DESCENT

1. Initialize $\mathbf{w}^0 \in \mathbb{R}^d$
2. Select a coordinate $j_t \in [d]$
3. Let $\mathbf{g}_{j_t}^t \leftarrow \nabla_{j_t} f(\mathbf{w}^t) + \nabla_{j_t} r(\mathbf{w}^t) \in \mathbb{R}$
4. Update $\mathbf{w}_{j_t}^{t+1} \leftarrow \mathbf{w}_{j_t}^t - \eta_t \cdot \mathbf{g}_{j_t}^t$

Only $O(n)$ time
per iter!

Can work with a block
of coordinates instead

Preserve other coord. $\mathbf{w}_k^{t+1} \leftarrow \mathbf{w}_k^t, k \neq j_t$

Repeat until convergence



Coordinate Descent

- Similar to gradient descent except update one variable at a time
- Notation: $\nabla_i f(\mathbf{w}) = [\nabla f(\mathbf{w})]_i$ i.e. the i -th directional derivative

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$

Useful when
 $d \gg n$

Random (Stochastic
CD), Cyclic
 $1, 2, \dots, d, 1, 2, 3 \dots d, \dots$

COORDINATE DESCENT

1. Initialize $\mathbf{w}^0 \in \mathbb{R}^d$

2. Select a coordinate $j_t \in [d]$

3. Let $\mathbf{g}_{j_t}^t \leftarrow \nabla_{j_t} f(\mathbf{w}^t) + \nabla_{j_t} r(\mathbf{w}^t) \in \mathbb{R}$

4. Update $\mathbf{w}_{j_t}^{t+1} \leftarrow \mathbf{w}_{j_t}^t - \eta_t \cdot \mathbf{g}_{j_t}^t$

5. Preserve other coord. $\mathbf{w}_k^{t+1} \leftarrow \mathbf{w}_k^t, k \neq j_t$

6. Repeat until convergence

Only $O(n)$ time
per iter!

Block Coord.
Descent

Can work with a block
of coordinates instead



Coordinate Descent

- Similar to gradient descent except update one variable at a time
- Notation: $\nabla_i f(\mathbf{w}) = [\nabla f(\mathbf{w})]_i$ i.e. the i -th directional derivative

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$

Useful when
 $d \gg n$

Somewhat
like mini-
batch SGD

Block Coord.
Descent

Can work with a block
of coordinates instead

COORDINATE DESCENT

1. Initialize $\mathbf{w}^0 \in \mathbb{R}^d$
2. Select a coordinate $j_t \in [d]$
3. Let $\mathbf{g}_{j_t}^t \leftarrow \nabla_{j_t} f(\mathbf{w}^t) + \nabla_{j_t} r(\mathbf{w}^t) \in \mathbb{R}$
4. Update $\mathbf{w}_{j_t}^{t+1} \leftarrow \mathbf{w}_{j_t}^t - \eta_t \cdot \mathbf{g}_{j_t}^t$
5. Preserve other coord. $\mathbf{w}_k^{t+1} \leftarrow \mathbf{w}_k^t, k \neq j_t$
6. Repeat until convergence

Random (Stochastic
CD), Cyclic
 $1, 2, \dots, d, 1, 2, 3 \dots d, \dots$

Only $O(n)$ time
per iter!



App: Linear Regression via CD

September 1, 2017



App: Linear Regression via CD

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

App: Linear Regression via CD

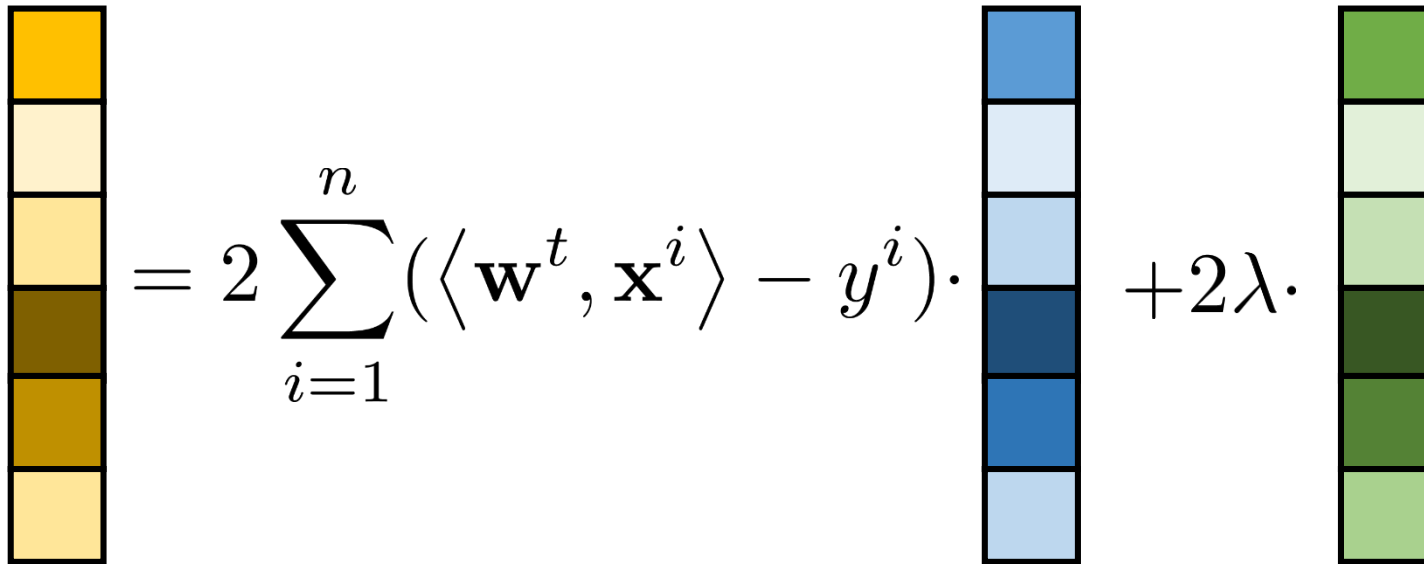
$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}^t = 2 \sum_{i=1}^n \left(\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i \right) \cdot \mathbf{x}^i + 2\lambda \cdot \mathbf{w}^t$$

App: Linear Regression via CD

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}^t = 2 \sum_{i=1}^n ((\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}^i + 2\lambda \cdot \mathbf{w}^t$$

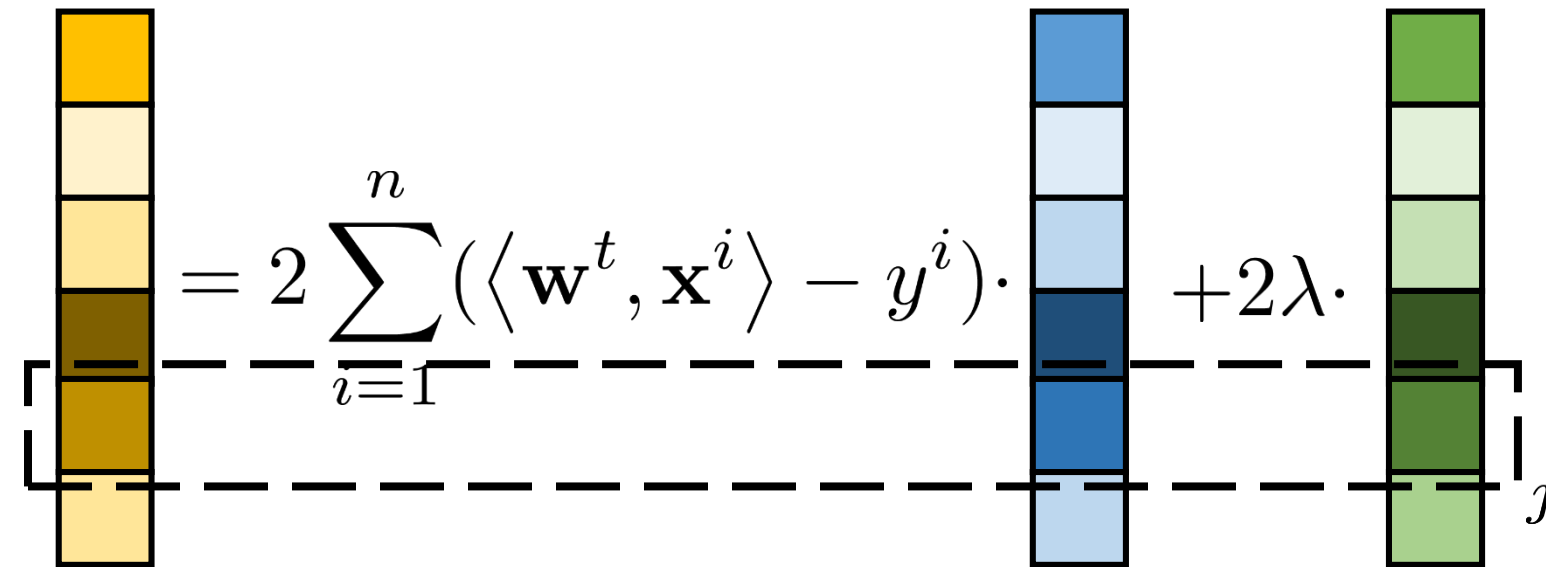


$$= 2 \sum_{i=1}^n ((\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}^i + 2\lambda \cdot \mathbf{w}^t$$

App: Linear Regression via CD

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

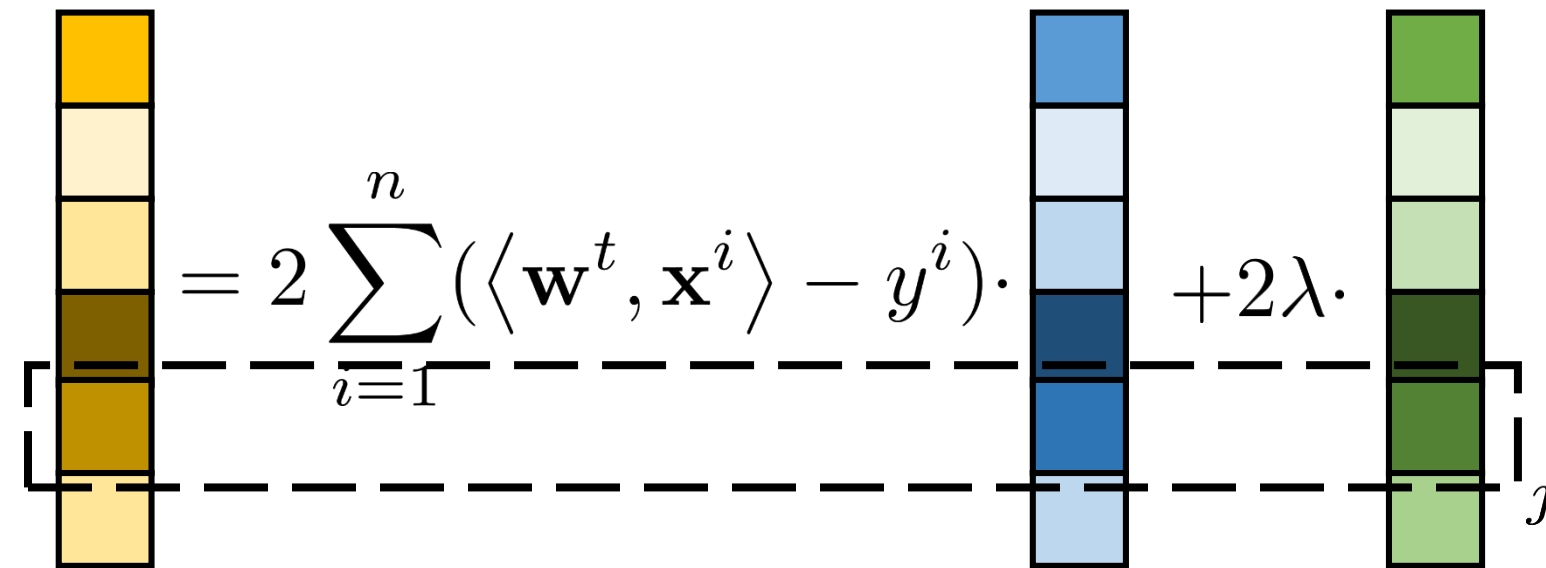
$$\mathbf{g}^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}^i + 2\lambda \cdot \mathbf{w}^t$$



App: Linear Regression via CD

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

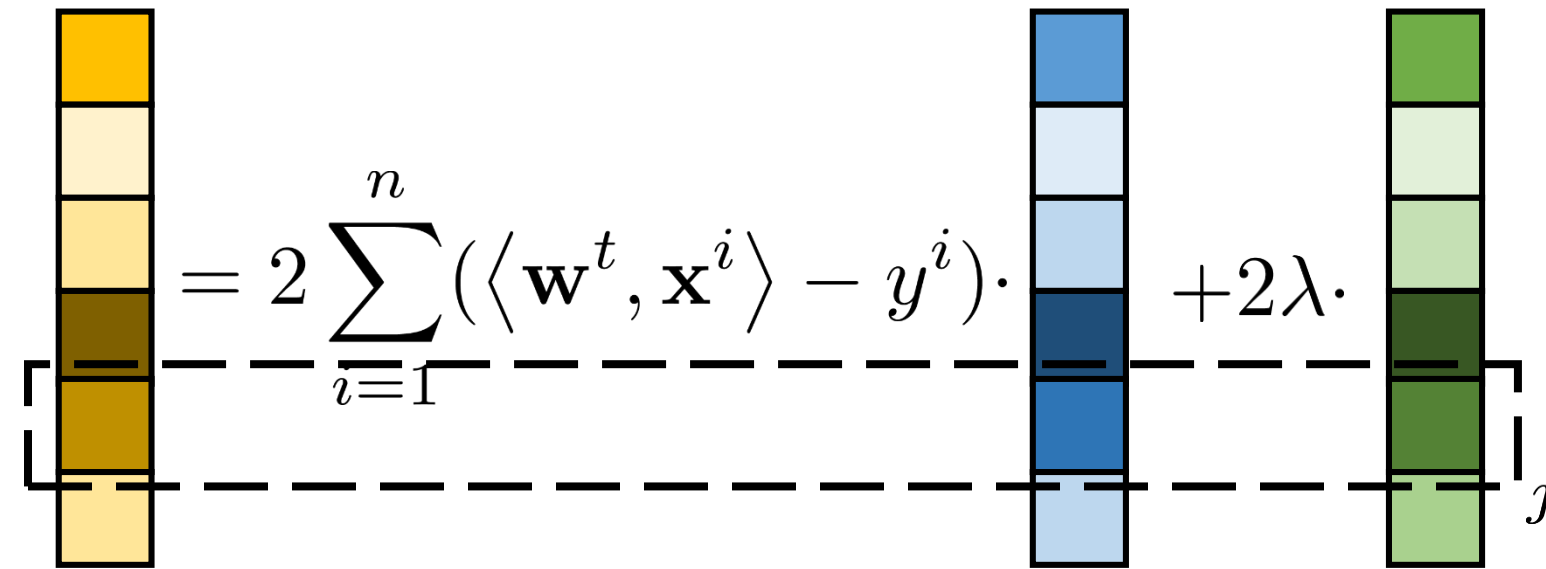
$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}^i + 2\lambda \cdot \mathbf{w}^t$$



App: Linear Regression via CD

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

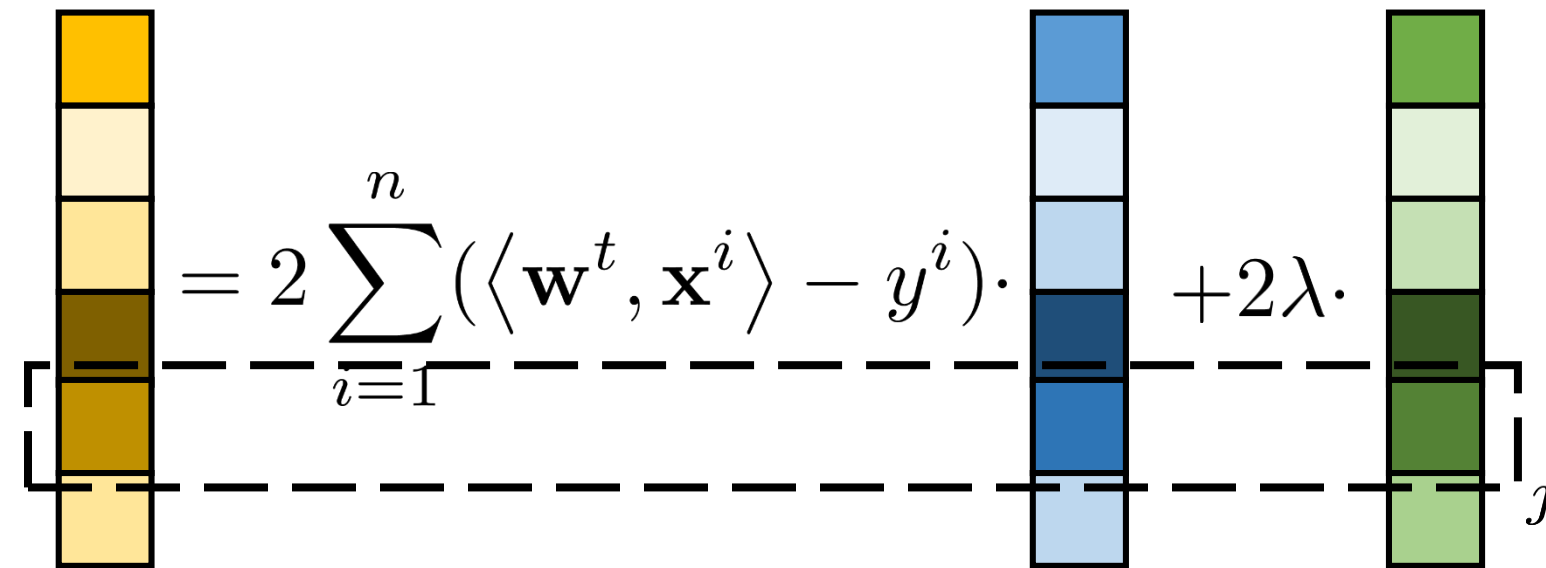
$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}^t$$



App: Linear Regression via CD

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n ((\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i) + 2\lambda \cdot \mathbf{w}_j^t$$

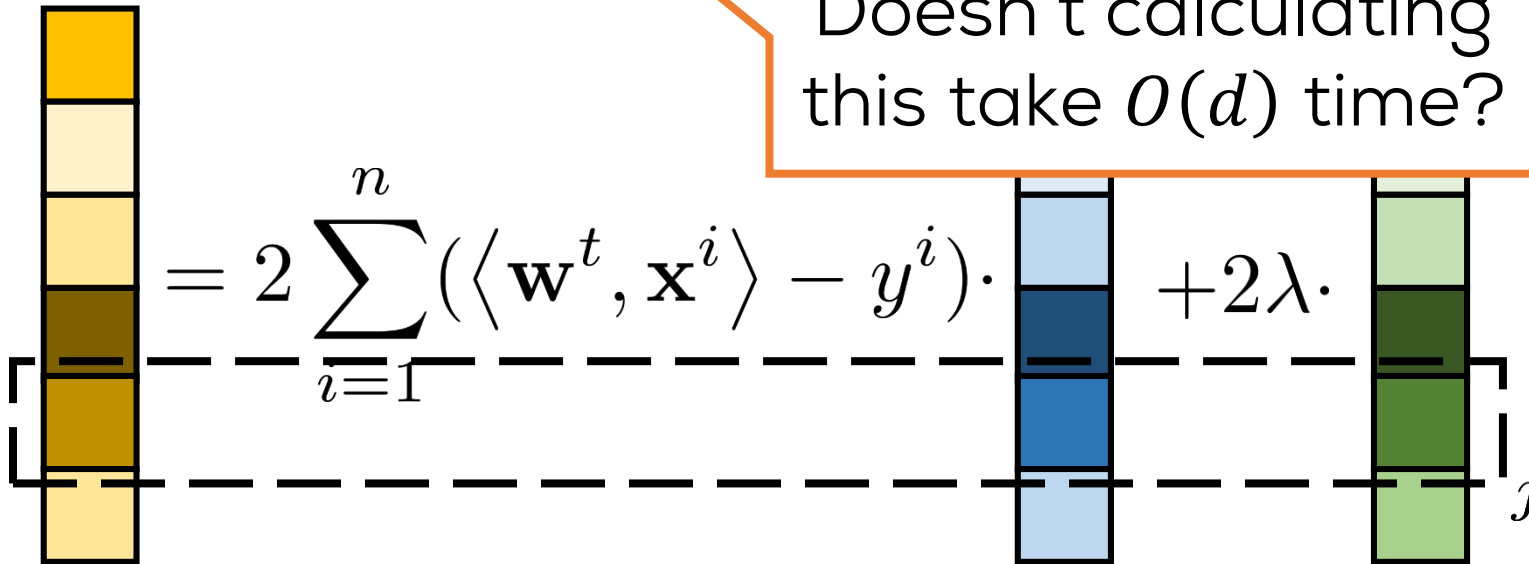


App: Linear Regression via CD

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

Doesn't calculating this take $O(d)$ time?



App: Linear Regression via CD

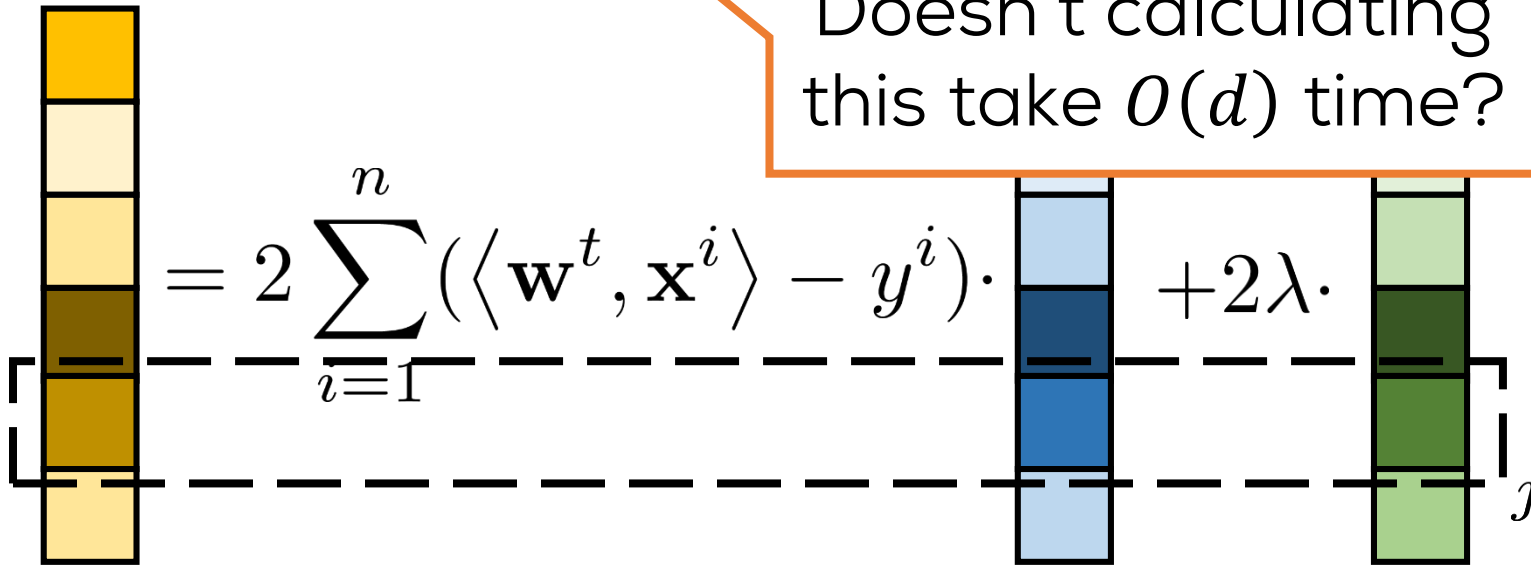


$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(nd)$ time
in total

Doesn't calculating
this take $O(d)$ time?



App: Linear Regression via CD



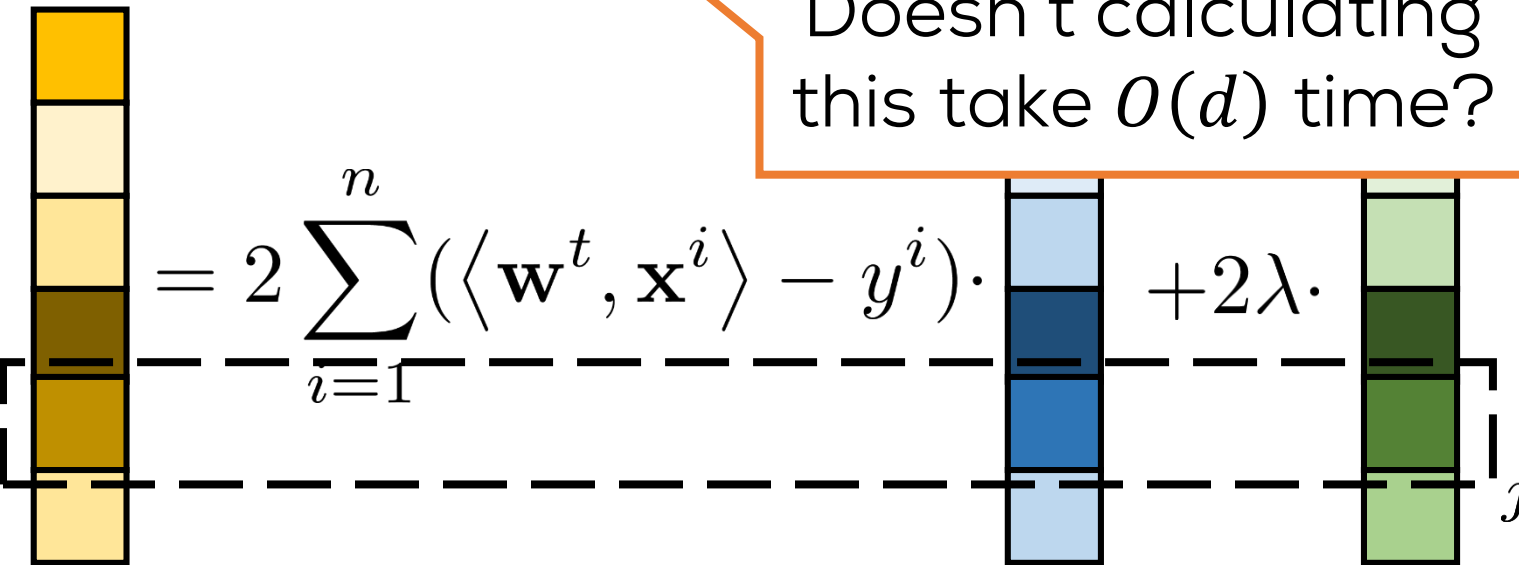
$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(nd)$ time
in total

Doesn't calculating
this take $O(d)$ time?

Only if you are
not careful



App: Linear Regression via CD

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(nd)$ time
in total



App: Linear Regression via CD

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$$\mathbf{w}^0 = \mathbf{0}$$

$O(nd)$ time
in total



App: Linear Regression via CD

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(nd)$ time
in total

$$\mathbf{w}^0 = \mathbf{0}$$

$$\langle \mathbf{w}^0, \mathbf{x}^i \rangle = 0$$



App: Linear Regression via CD

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(nd)$ time
in total

$$\mathbf{w}^0 = \mathbf{0}$$

$$\langle \mathbf{w}^0, \mathbf{x}^i \rangle = 0$$

$O(n)$ time



App: Linear Regression via CD

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(nd)$ time
in total

$$\mathbf{w}^0 = \mathbf{0} \quad \mathbf{w}^t$$

$$\langle \mathbf{w}^0, \mathbf{x}^i \rangle = 0$$

$O(n)$ time



App: Linear Regression via CD

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(nd)$ time
in total

$$\mathbf{w}^0 = \mathbf{0} \quad \mathbf{w}^t$$

$$\langle \mathbf{w}^0, \mathbf{x}^i \rangle = 0 \quad \langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

$O(n)$ time



App: Linear Regression via CD

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(nd)$ time
in total

$$\mathbf{w}^0 = \mathbf{0} \quad \mathbf{w}^t$$

$$\langle \mathbf{w}^0, \mathbf{x}^i \rangle = 0 \quad \langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

$O(n)$ time

Available



App: Linear Regression via CD



$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(nd)$ time
in total

$$\mathbf{w}^0 = \mathbf{0} \quad \mathbf{w}^t$$

$$\langle \mathbf{w}^0, \mathbf{x}^i \rangle = 0 \quad \langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

$O(n)$ time

Available

Assumed

App: Linear Regression via CD



$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(nd)$ time
in total

$$\mathbf{w}^0 = \mathbf{0} \quad \mathbf{w}^t \quad \mathbf{w}_{j_t}^{t+1} = \mathbf{w}_{j_t}^t - \eta_t \cdot \mathbf{g}_{j_t}^t$$

$$\langle \mathbf{w}^0, \mathbf{x}^i \rangle = 0 \quad \langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

$O(n)$ time

Available

Assumed

App: Linear Regression via CD



$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(nd)$ time
in total

$$\mathbf{w}^0 = \mathbf{0}$$

$$\mathbf{w}^t$$

$$\mathbf{w}_{j_t}^{t+1} = \mathbf{w}_{j_t}^t - \eta_t \cdot \mathbf{g}_{j_t}^t$$

$$\langle \mathbf{w}^0, \mathbf{x}^i \rangle = 0$$

$$\langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

$$\langle \mathbf{w}^{t+1}, \mathbf{x}^i \rangle = \langle \mathbf{w}^t, \mathbf{x}^i \rangle - \eta_t \cdot \mathbf{g}_{j_t}^t \cdot \mathbf{x}_{j_t}^i$$

$O(n)$ time

Available

Assumed

App: Linear Regression via CD



$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(nd)$ time
in total

$$\mathbf{w}^0 = \mathbf{0}$$

$$\mathbf{w}^t$$

$$\mathbf{w}_{j_t}^{t+1} = \mathbf{w}_{j_t}^t - \eta_t \cdot \mathbf{g}_{j_t}^t$$

$$\langle \mathbf{w}^0, \mathbf{x}^i \rangle = 0$$

$$\langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

$$\langle \mathbf{w}^{t+1}, \mathbf{x}^i \rangle = \langle \mathbf{w}^t, \mathbf{x}^i \rangle - \eta_t \cdot \mathbf{g}_{j_t}^t \cdot \mathbf{x}_{j_t}^i$$

$O(n)$ time

Available

Already
available

Assumed

App: Linear Regression via CD



$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(nd)$ time
in total

$$\mathbf{w}^0 = \mathbf{0}$$

$$\mathbf{w}^t$$

$$\mathbf{w}_{j_t}^{t+1} = \mathbf{w}_{j_t}^t - \eta_t \cdot \mathbf{g}_{j_t}^t$$

$$\langle \mathbf{w}^0, \mathbf{x}^i \rangle = 0$$

$$\langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

$$\langle \mathbf{w}^{t+1}, \mathbf{x}^i \rangle = \langle \mathbf{w}^t, \mathbf{x}^i \rangle - \eta_t \cdot \mathbf{g}_{j_t}^t \cdot \mathbf{x}_{j_t}^i$$

$O(n)$ time

Available

Assumed

Already
available

$O(1)$ time per
data point!

App: Linear Regression via CD



$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(nd)$ time
in total

$$\mathbf{w}^0 = \mathbf{0}$$

$$\mathbf{w}^t$$

$$\mathbf{w}_{j_t}^{t+1} = \mathbf{w}_{j_t}^t - \eta_t \cdot \mathbf{g}_{j_t}^t$$

$$\langle \mathbf{w}^0, \mathbf{x}^i \rangle = 0$$

$$\langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

$$\langle \mathbf{w}^{t+1}, \mathbf{x}^i \rangle = \langle \mathbf{w}^t, \mathbf{x}^i \rangle - \eta_t \cdot \mathbf{g}_{j_t}^t \cdot \mathbf{x}_{j_t}^i$$

$O(n)$ time

Available

Assumed

Already
available

$O(1)$ time per
data point!

$O(n)$ time
in total

App: Linear Regression via CD



$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(nd)$ time
in total

$$\mathbf{w}^0 = \mathbf{0}$$

$$\mathbf{w}^t$$

$$\mathbf{w}_{j_t}^{t+1} = \mathbf{w}_{j_t}^t - \eta_t \cdot \mathbf{g}_{j_t}^t$$

$$\langle \mathbf{w}^0, \mathbf{x}^i \rangle = 0$$

$$\langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

$$\langle \mathbf{w}^{t+1}, \mathbf{x}^i \rangle = \langle \mathbf{w}^t, \mathbf{x}^i \rangle - \eta_t \cdot \mathbf{g}_{j_t}^t \cdot \mathbf{x}_{j_t}^i$$

$O(n)$ time

Available

Assumed

Already
available

Induction
complete!

$O(1)$ time per
data point!

$O(n)$ time
in total

App: Linear Regression via CD



$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(n)$ time
in total

$$\mathbf{w}^0 = \mathbf{0}$$

$$\mathbf{w}^t$$

$$\mathbf{w}_{j_t}^{t+1} = \mathbf{w}_{j_t}^t - \eta_t \cdot \mathbf{g}_{j_t}^t$$

$$\langle \mathbf{w}^0, \mathbf{x}^i \rangle = 0$$

$$\langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

$$\langle \mathbf{w}^{t+1}, \mathbf{x}^i \rangle = \langle \mathbf{w}^t, \mathbf{x}^i \rangle - \eta_t \cdot \mathbf{g}_{j_t}^t \cdot \mathbf{x}_{j_t}^i$$

$O(n)$ time

Available

Assumed

Already
available

Induction
complete!

$O(1)$ time per
data point!

$O(n)$ time
in total

App: Linear Regression via CD



$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{g}_j^t = 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}_j^i + 2\lambda \cdot \mathbf{w}_j^t$$

$O(n)$ time
in total

$$\mathbf{w}^0 = \mathbf{0}$$

$$\mathbf{w}^t$$

$$\mathbf{w}_{j_t}^{t+1} = \mathbf{w}_{j_t}^t - \eta_t \cdot \mathbf{g}_{j_t}^t$$

$$\langle \mathbf{w}^0, \mathbf{x}^i \rangle = 0$$

$$\langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

$$\langle \mathbf{w}^{t+1}, \mathbf{x}^i \rangle = \langle \mathbf{w}^t, \mathbf{x}^i \rangle - \eta_t \cdot \mathbf{g}_{j_t}^t \cdot \mathbf{x}_{j_t}^i$$

$O(n)$ time

Available

Assumed

Already
available

Induction
complete!

$O(1)$ time per
data point!

$O(n)$ time
in total

Handling Constraints

September 1, 2017



The Tale of Two Techniques

September 1, 2017



The Tale of Two Techniques

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^r (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2$$

$$\text{s.t. } \|\mathbf{w}\|_2 \leq r$$

The Tale of Two Techniques

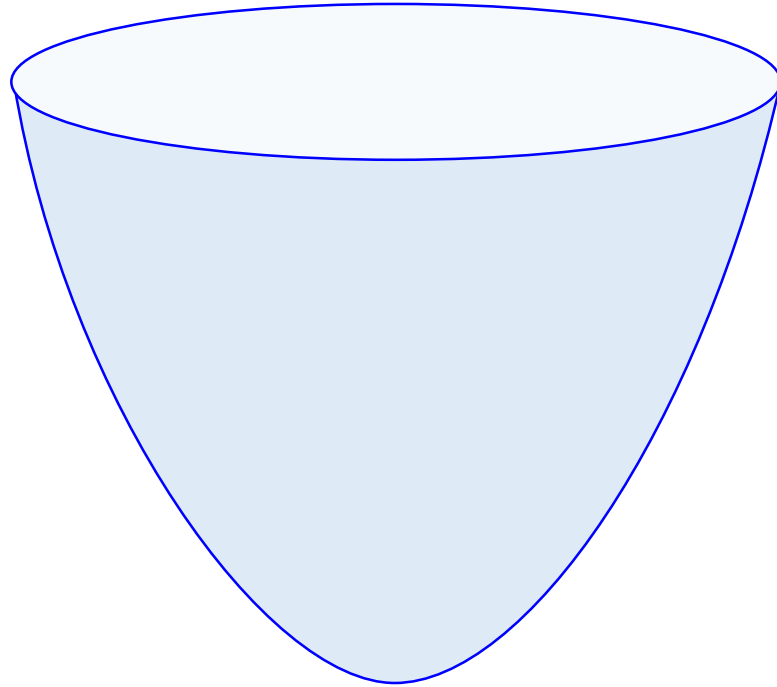
$$\hat{\mathbf{w}} = \arg \min f(\mathbf{w})$$

s.t. $\mathbf{w} \in \mathcal{C}$

The Tale of Two Techniques

$$\hat{\mathbf{w}} = \arg \min f(\mathbf{w})$$

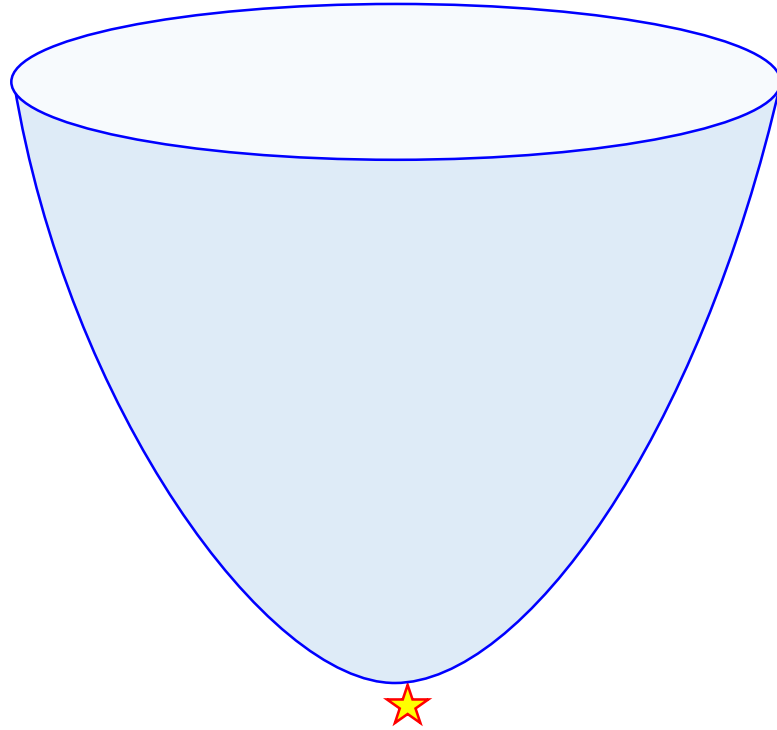
s.t. $\mathbf{w} \in \mathcal{C}$



The Tale of Two Techniques

$$\hat{\mathbf{w}} = \arg \min f(\mathbf{w})$$

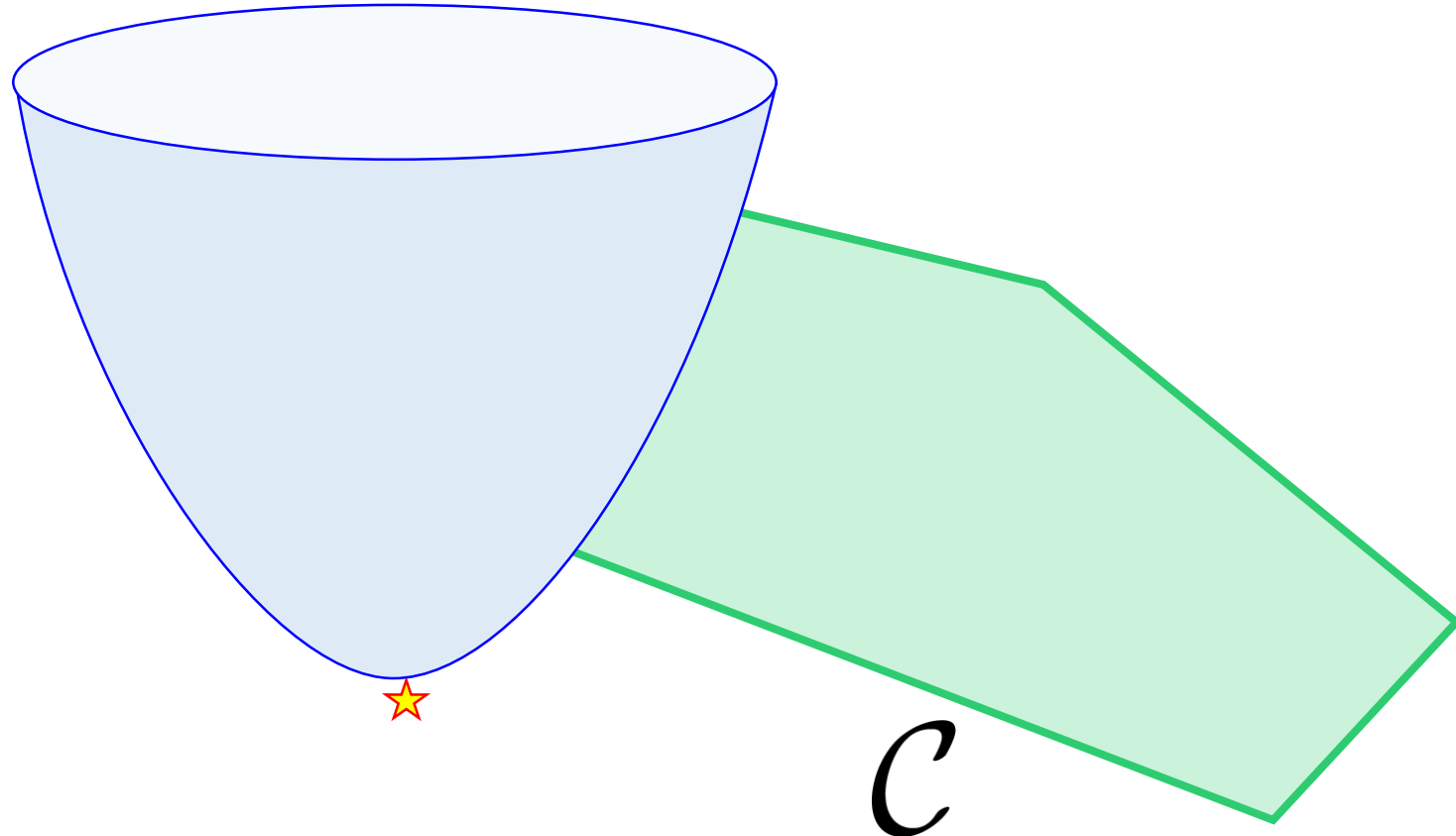
s.t. $\mathbf{w} \in \mathcal{C}$



The Tale of Two Techniques

$$\hat{\mathbf{w}} = \arg \min f(\mathbf{w})$$

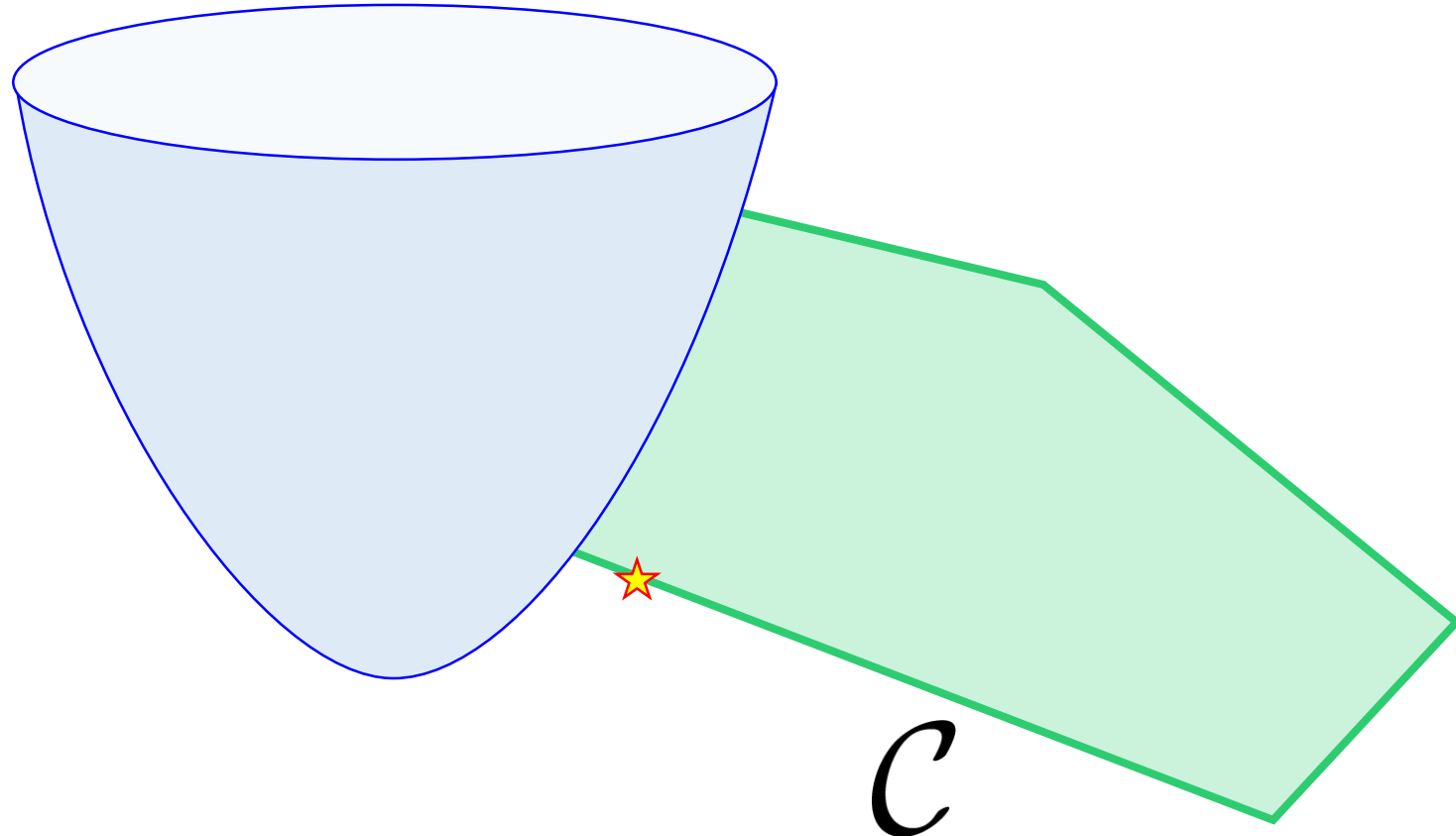
s.t. $\mathbf{w} \in \mathcal{C}$



The Tale of Two Techniques

$$\hat{\mathbf{w}} = \arg \min f(\mathbf{w})$$

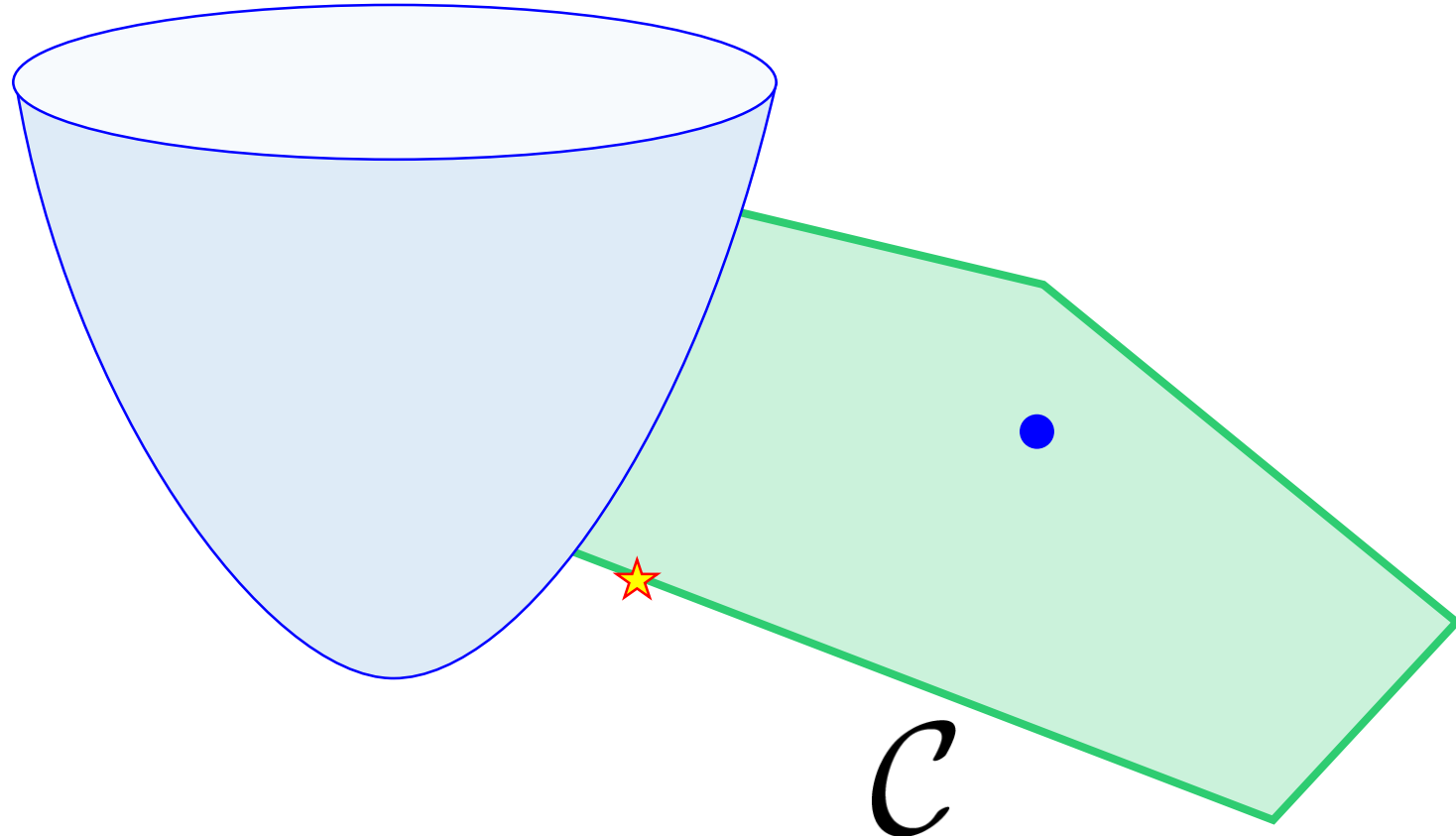
s.t. $\mathbf{w} \in \mathcal{C}$



The Tale of Two Techniques

$$\hat{\mathbf{w}} = \arg \min f(\mathbf{w})$$

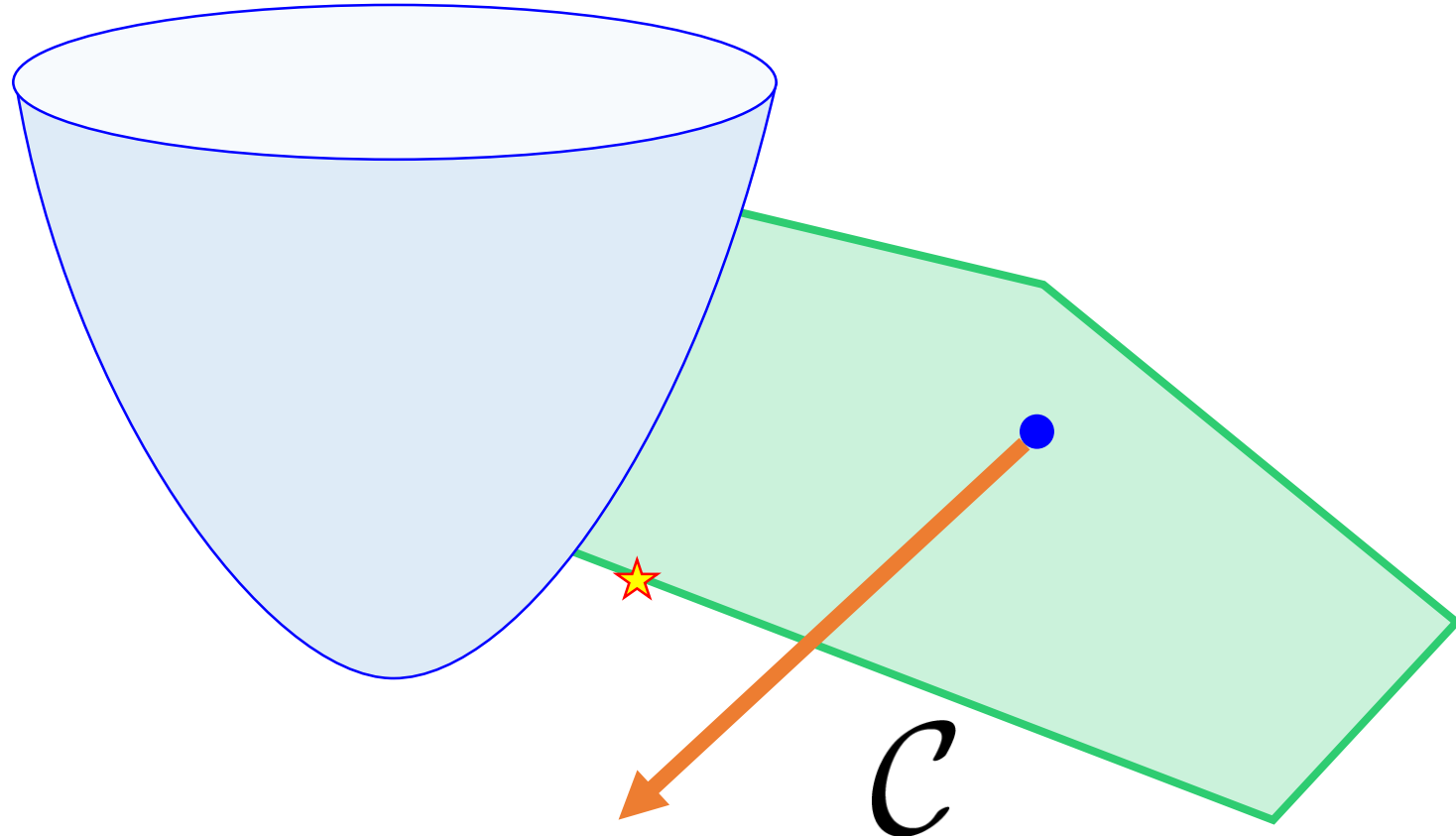
s.t. $\mathbf{w} \in \mathcal{C}$



The Tale of Two Techniques

$$\hat{\mathbf{w}} = \arg \min f(\mathbf{w})$$

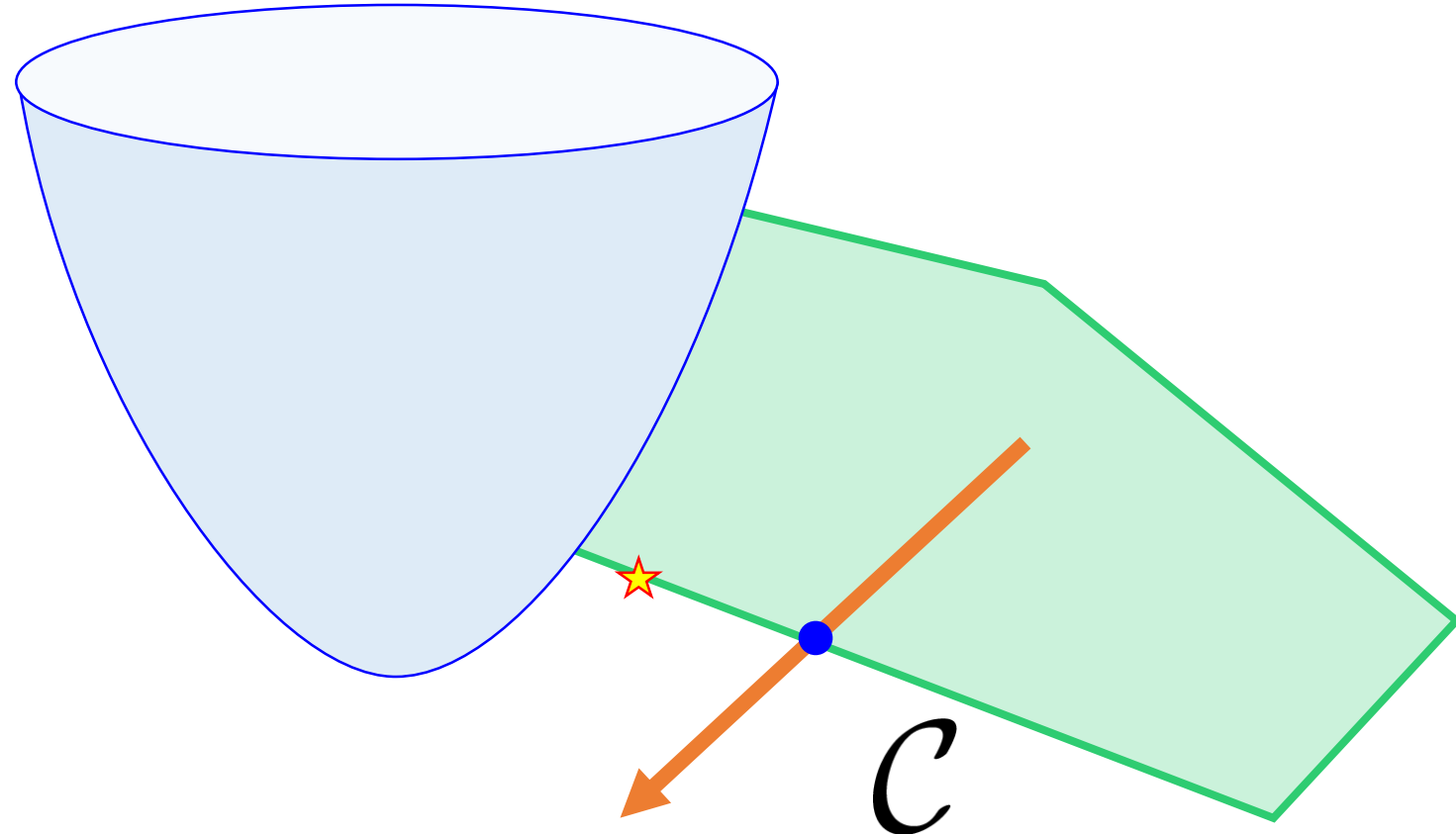
s.t. $\mathbf{w} \in \mathcal{C}$



The Tale of Two Techniques

$$\hat{\mathbf{w}} = \arg \min f(\mathbf{w})$$
$$\text{s.t. } \mathbf{w} \in \mathcal{C}$$

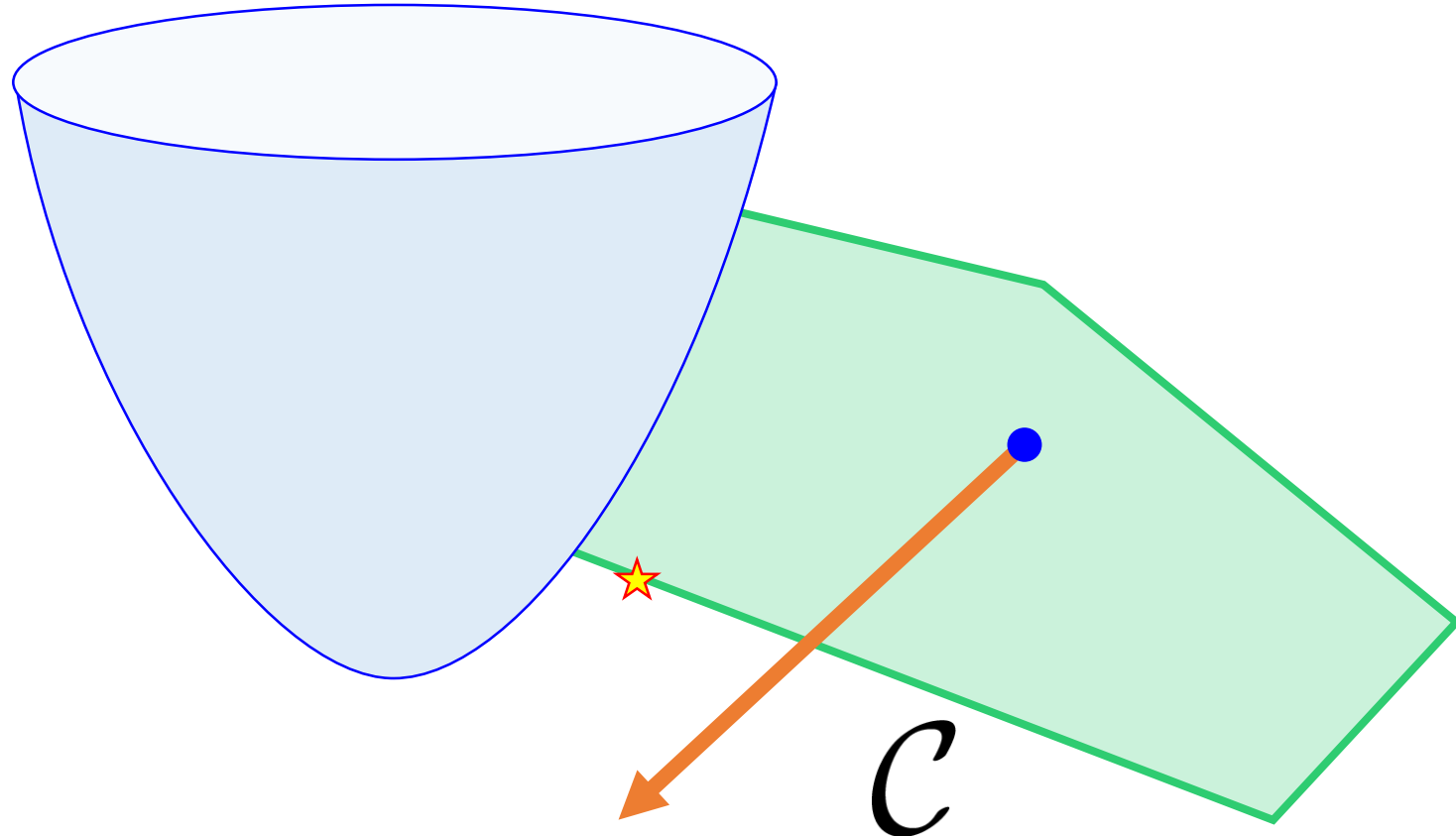
Interior Point
Method



The Tale of Two Techniques

$$\hat{\mathbf{w}} = \arg \min f(\mathbf{w})$$
$$\text{s.t. } \mathbf{w} \in \mathcal{C}$$

Interior Point
Method

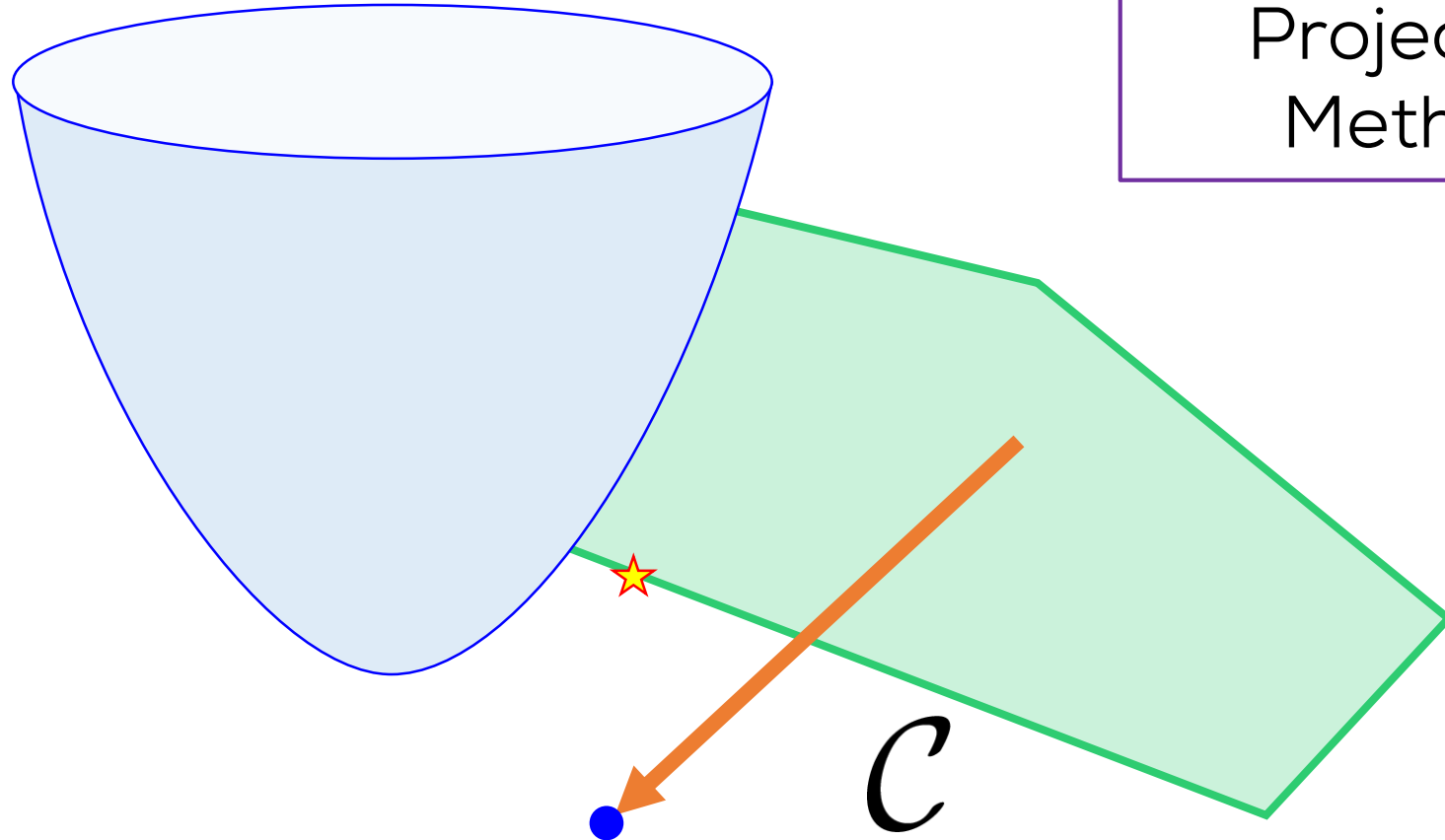


The Tale of Two Techniques

$$\hat{\mathbf{w}} = \arg \min f(\mathbf{w})$$
$$\text{s.t. } \mathbf{w} \in \mathcal{C}$$

Interior Point
Method

Projected
Method

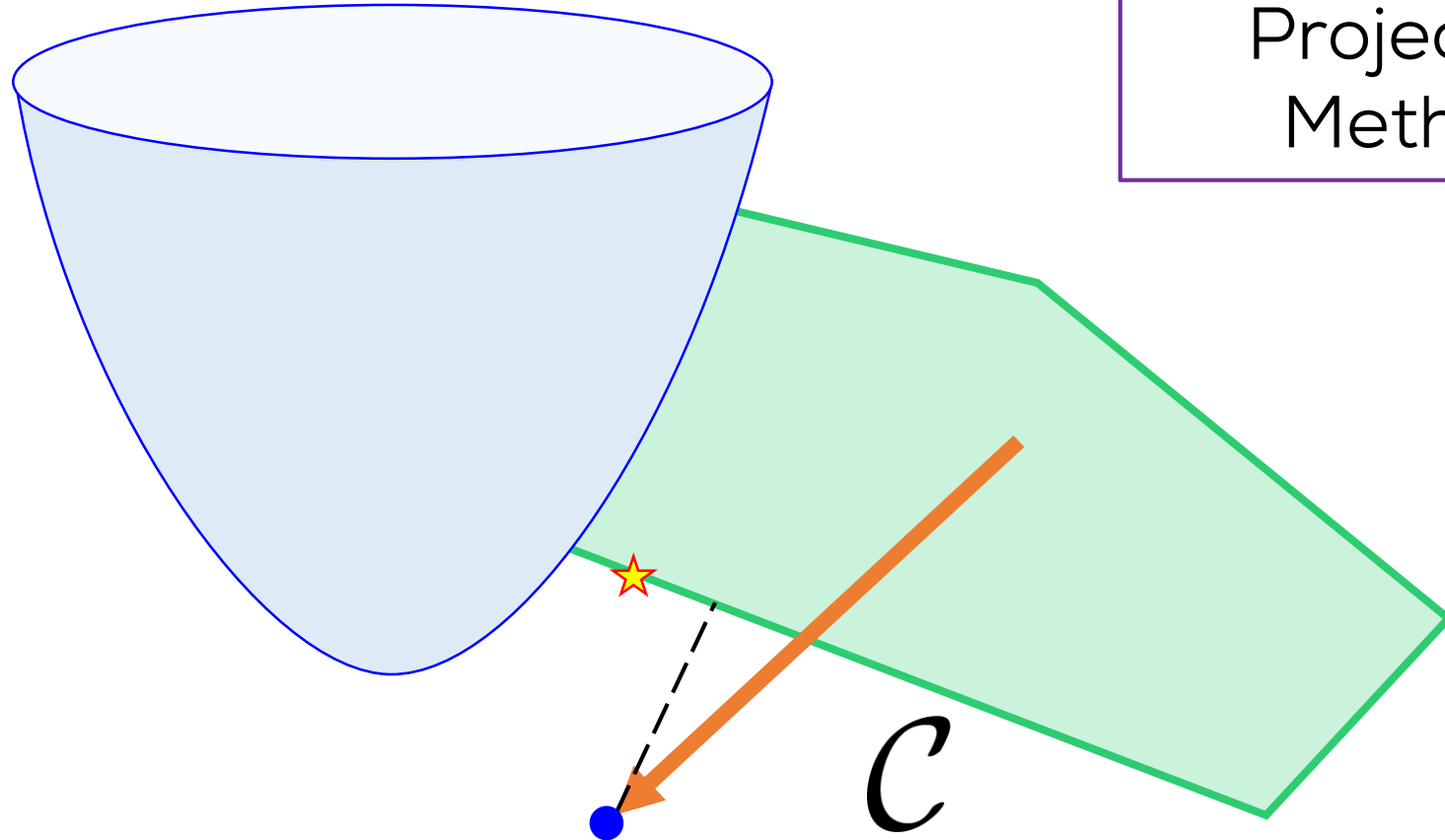


The Tale of Two Techniques

$$\hat{\mathbf{w}} = \arg \min f(\mathbf{w})$$
$$\text{s.t. } \mathbf{w} \in \mathcal{C}$$

Interior Point
Method

Projected
Method

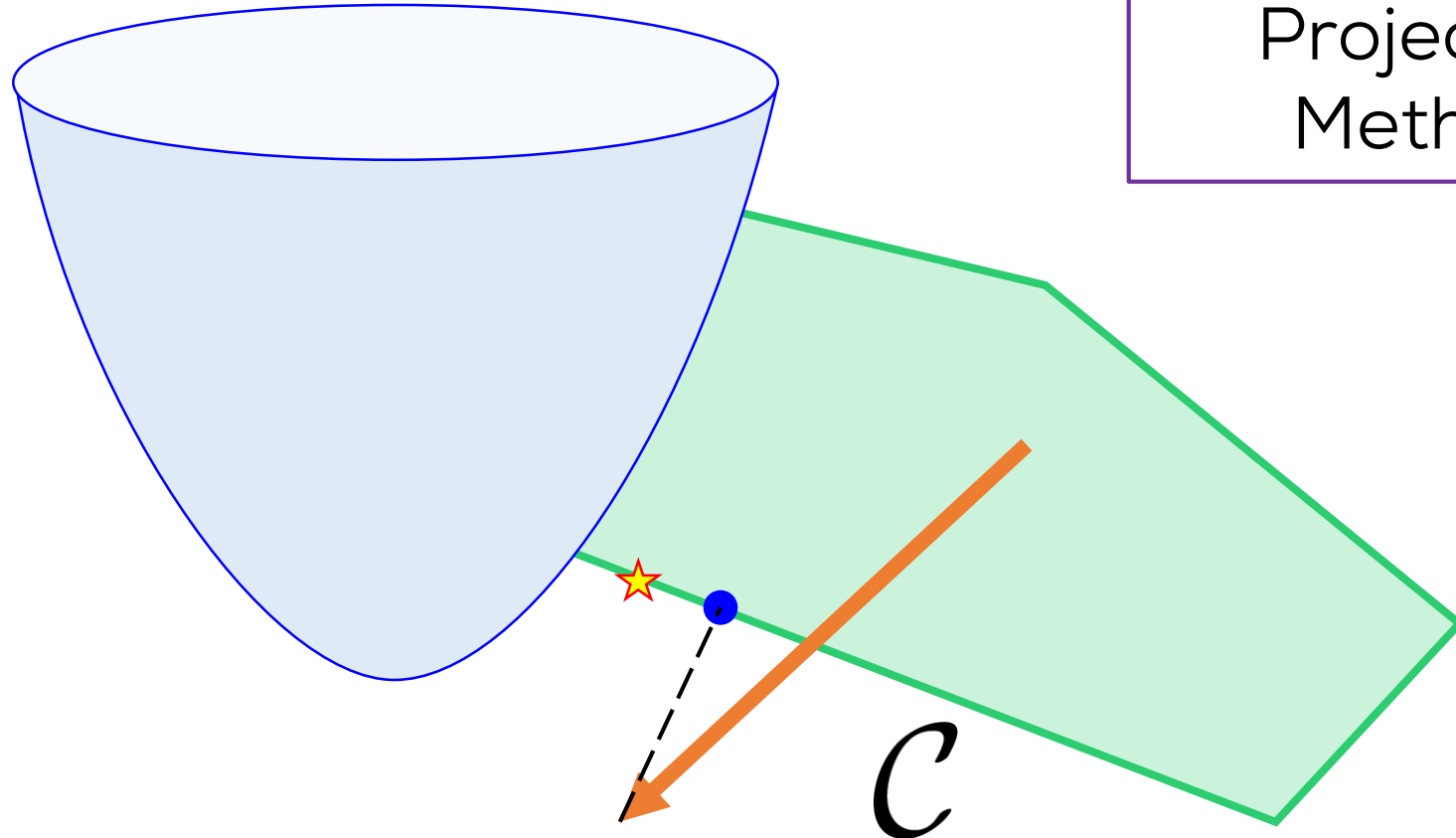


The Tale of Two Techniques

$$\hat{\mathbf{w}} = \arg \min f(\mathbf{w})$$
$$\text{s.t. } \mathbf{w} \in \mathcal{C}$$

Interior Point
Method

Projected
Method



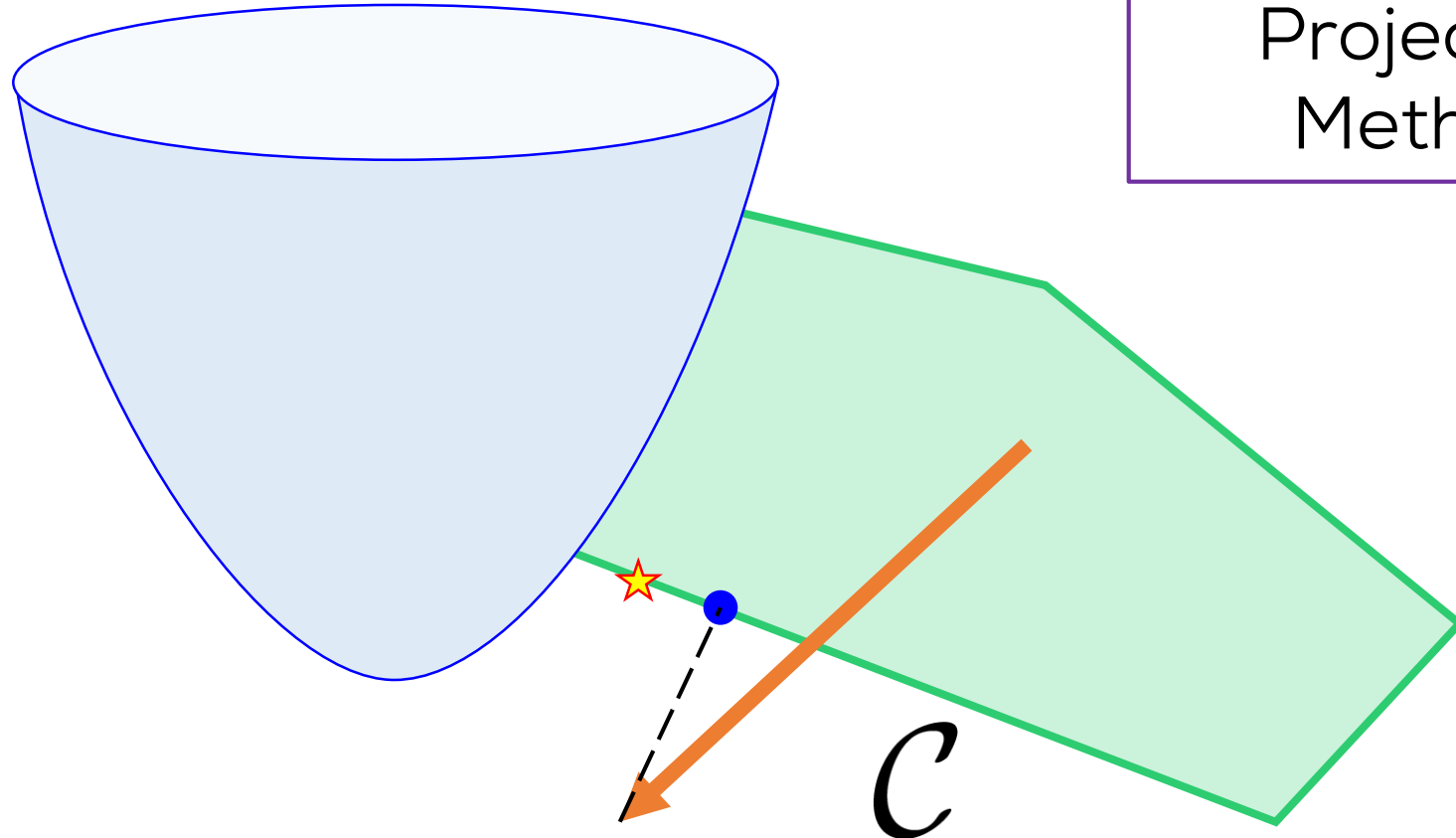
The Tale of Two Techniques

$$\hat{\mathbf{w}} = \arg \min f(\mathbf{w})$$
$$\text{s.t. } \mathbf{w} \in \mathcal{C}$$

Interior Point
Method

Powerful but can
be slow, difficult
to implement

Projected
Method

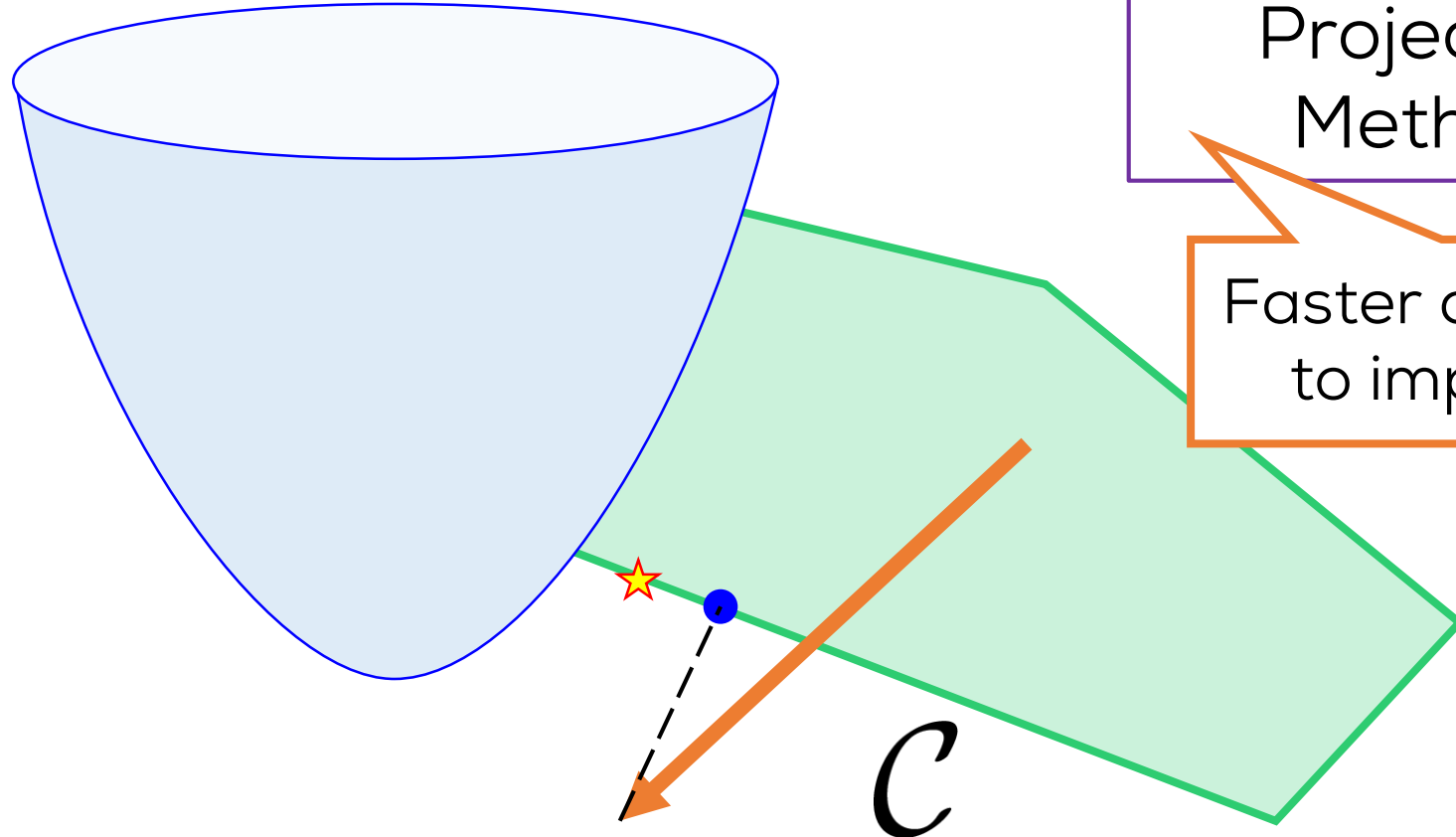


The Tale of Two Techniques

$$\hat{\mathbf{w}} = \arg \min f(\mathbf{w})$$
$$\text{s.t. } \mathbf{w} \in \mathcal{C}$$

Interior Point
Method

Powerful but can
be slow, difficult
to implement



Projected
Method

Faster and easier
to implement

The Tale of Two Techniques

$$\hat{\mathbf{w}} = \arg \min f(\mathbf{w})$$

s.t. $\mathbf{w} \in \mathcal{C}$

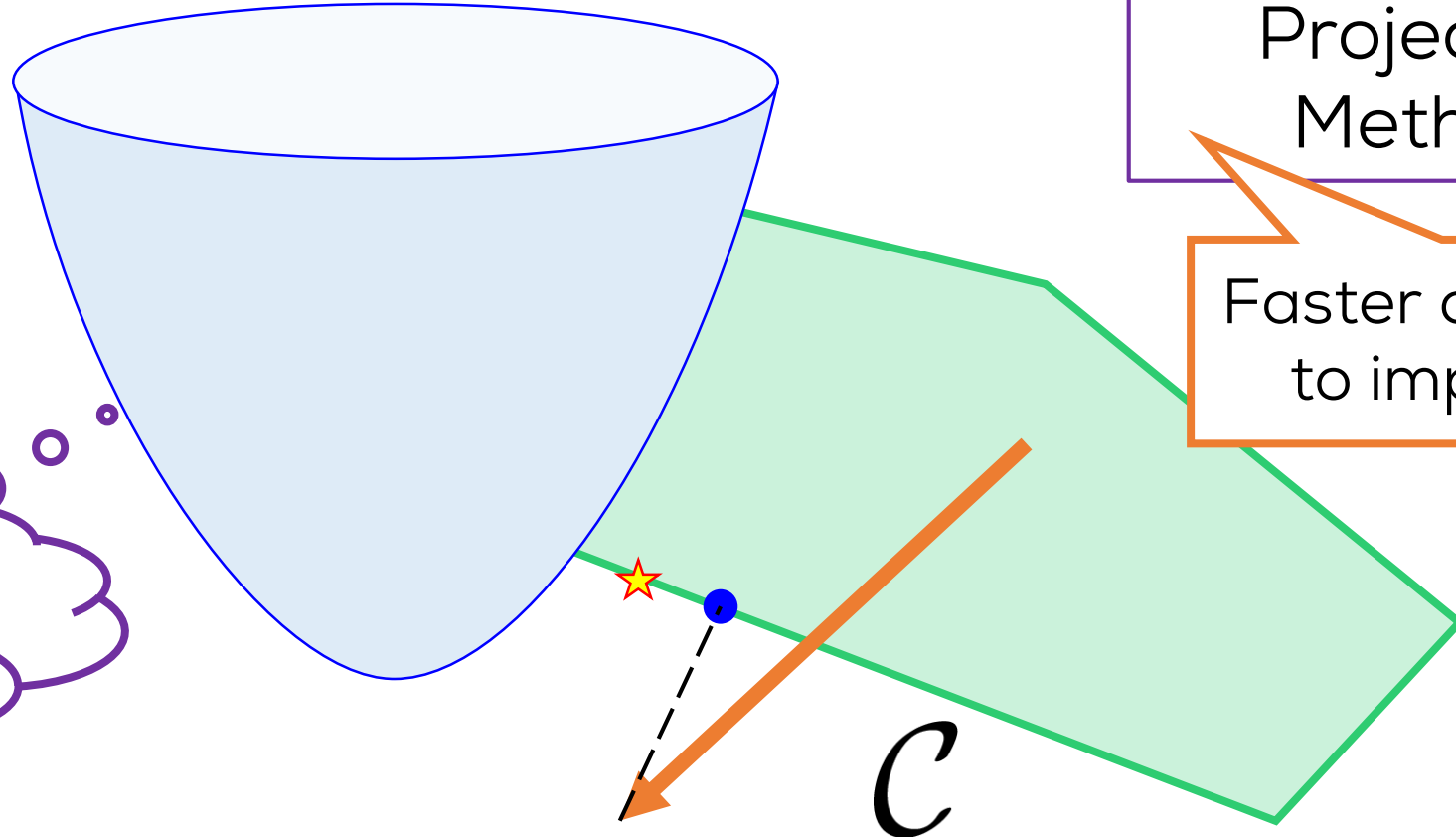
Interior Point
Method

Powerful but can
be slow, difficult
to implement

What if the set
is non-convex?

Projected
Method

Faster and easier
to implement



The Tale of Two Techniques

$$\hat{\mathbf{w}} = \arg \min f(\mathbf{w})$$

s.t. $\mathbf{w} \in \mathcal{C}$

Interior Point
Method

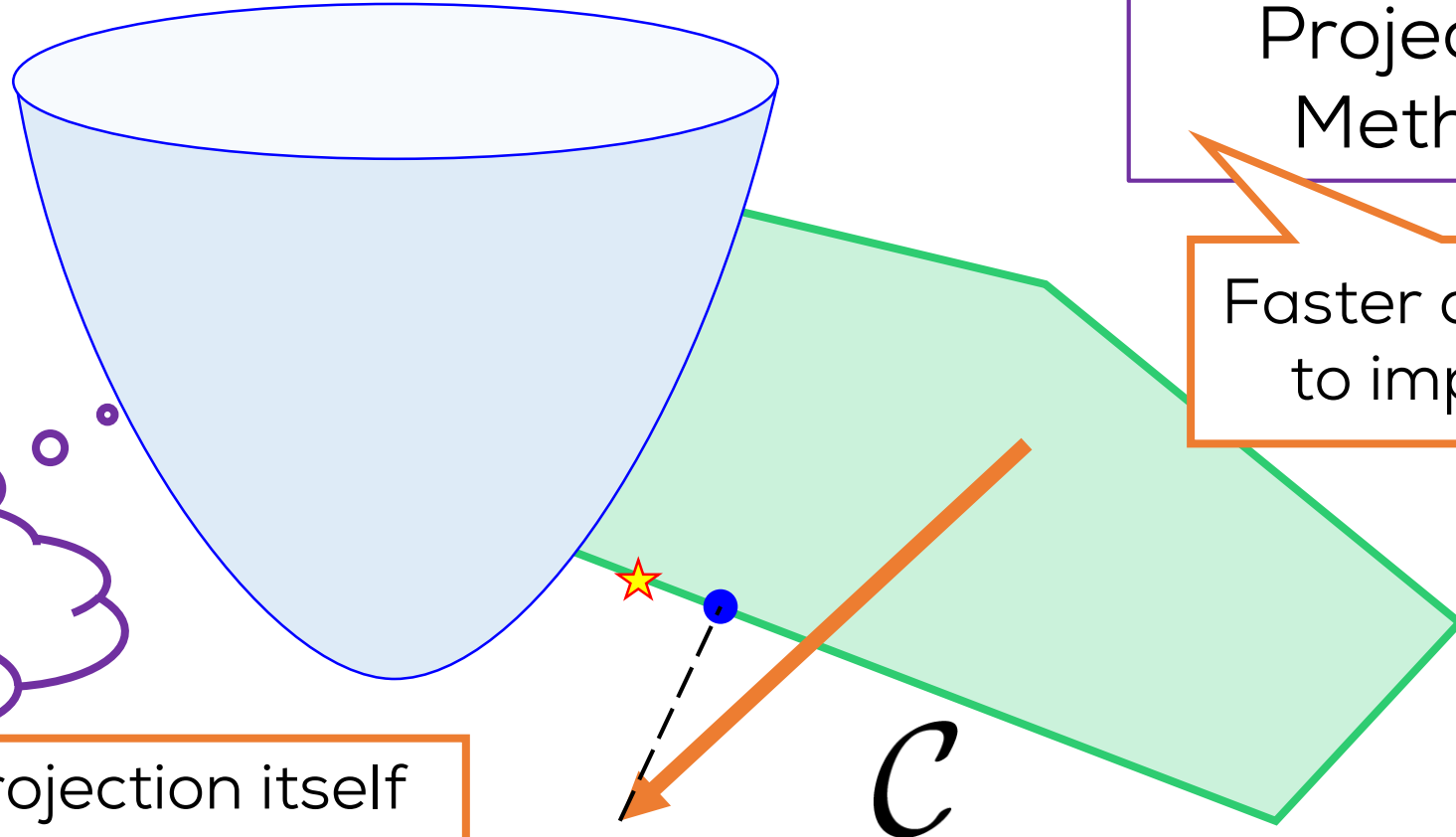
Powerful but can
be slow, difficult
to implement

What if the set
is non-convex?

Projection itself
can be NP-hard!

Projected
Method

Faster and easier
to implement



Projected Gradient Descent

September 1, 2017



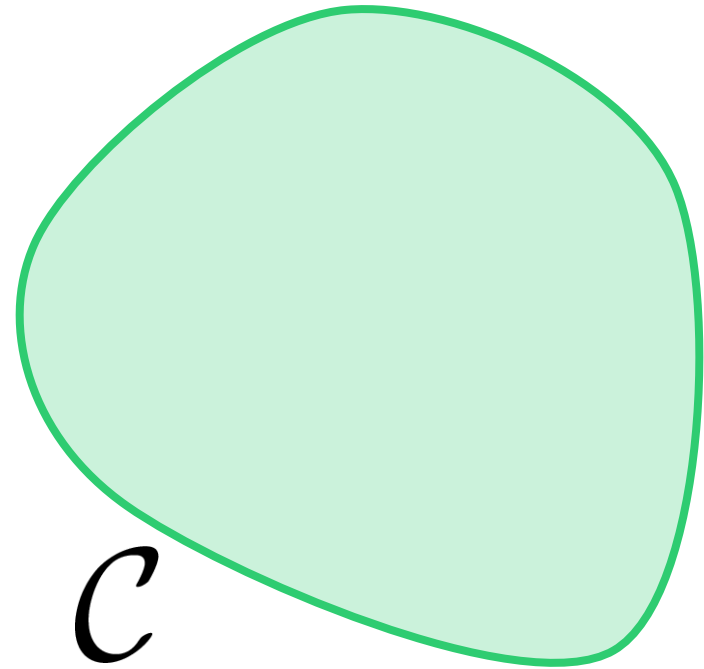
Projected Gradient Descent

$$\begin{aligned} \arg \min \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & \mathbf{w} \in \mathcal{C} \end{aligned}$$

September 1, 2017



Projected Gradient Descent



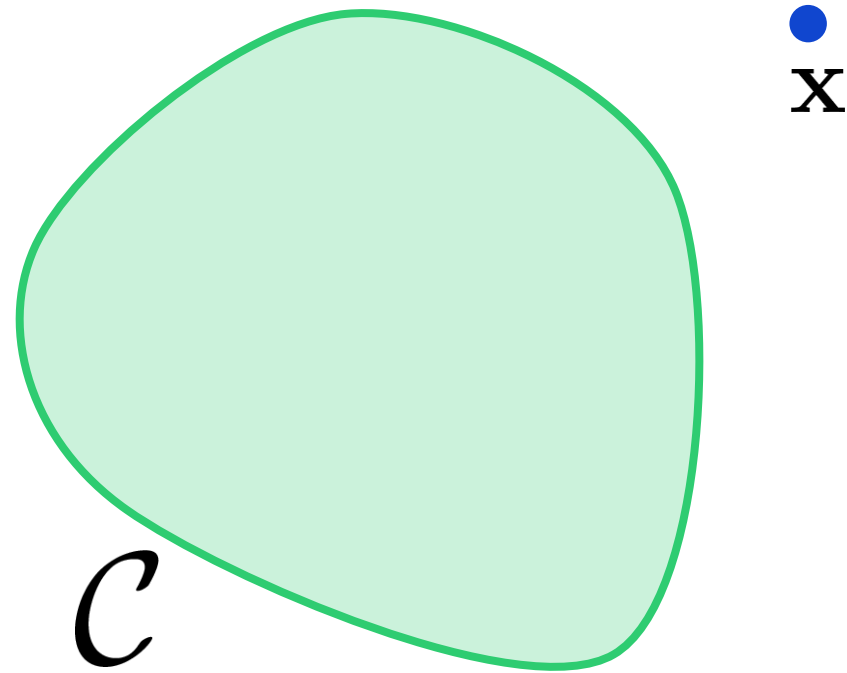
\mathcal{C}

$$\begin{aligned} \arg \min \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & \mathbf{w} \in \mathcal{C} \end{aligned}$$

September 1, 2017



Projected Gradient Descent

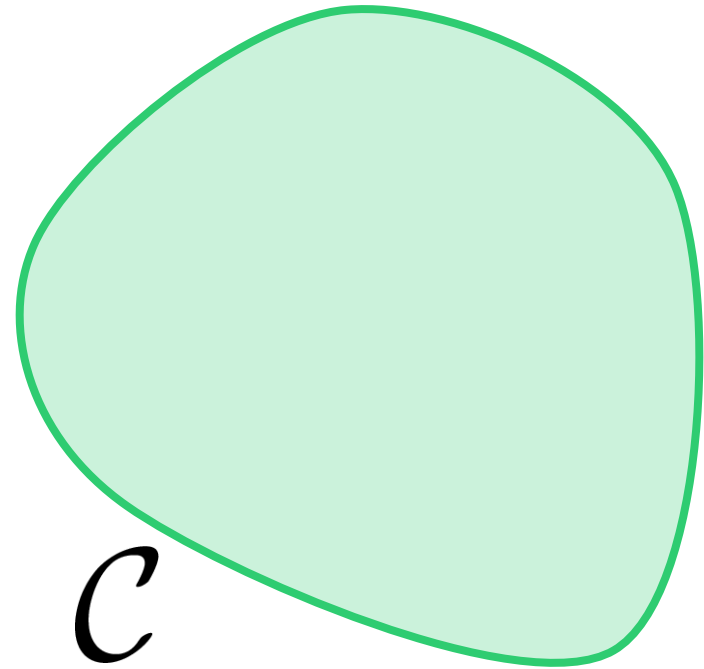


$$\begin{aligned} \arg \min f(\mathbf{w}) \\ \text{s.t. } \mathbf{w} \in \mathcal{C} \end{aligned}$$

September 1, 2017



Projected Gradient Descent

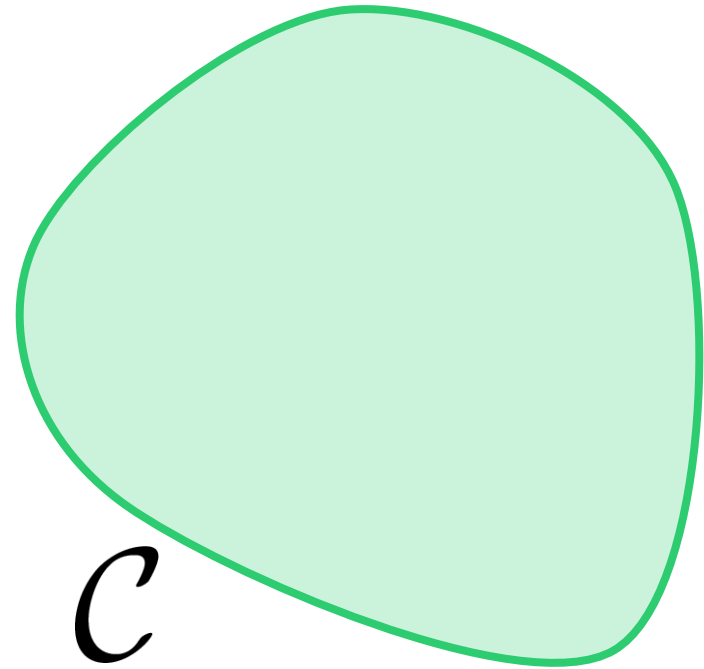


•
 \mathbf{x}

$$\Pi_{\mathcal{C}}(\mathbf{x}) = \arg \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{z} - \mathbf{x}\|$$

$$\begin{aligned} \arg \min f(\mathbf{w}) \\ \text{s.t. } \mathbf{w} \in \mathcal{C} \end{aligned}$$

Projected Gradient Descent



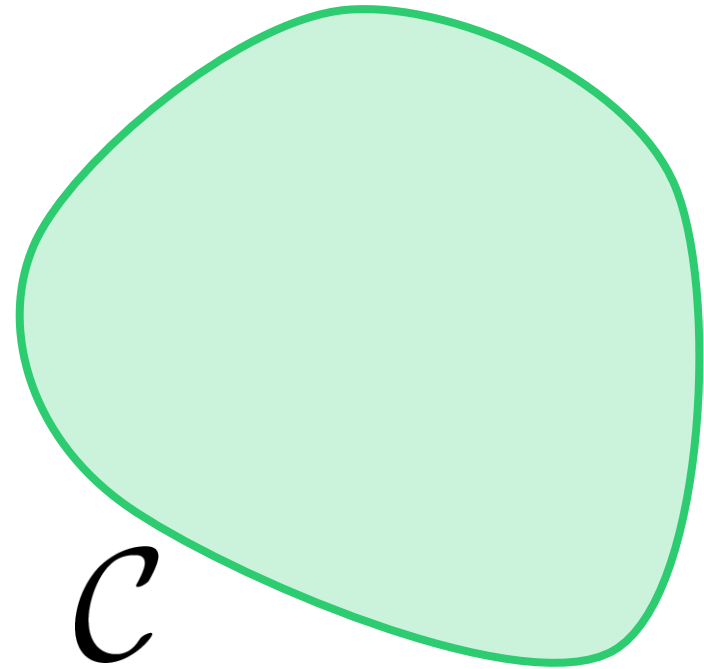
•
 \mathbf{x}

$$\Pi_{\mathcal{C}}(\mathbf{x}) = \arg \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{z} - \mathbf{x}\|_2$$

Can use other
norms too!

$$\begin{aligned} \arg \min \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & \mathbf{w} \in \mathcal{C} \end{aligned}$$

Projected Gradient Descent

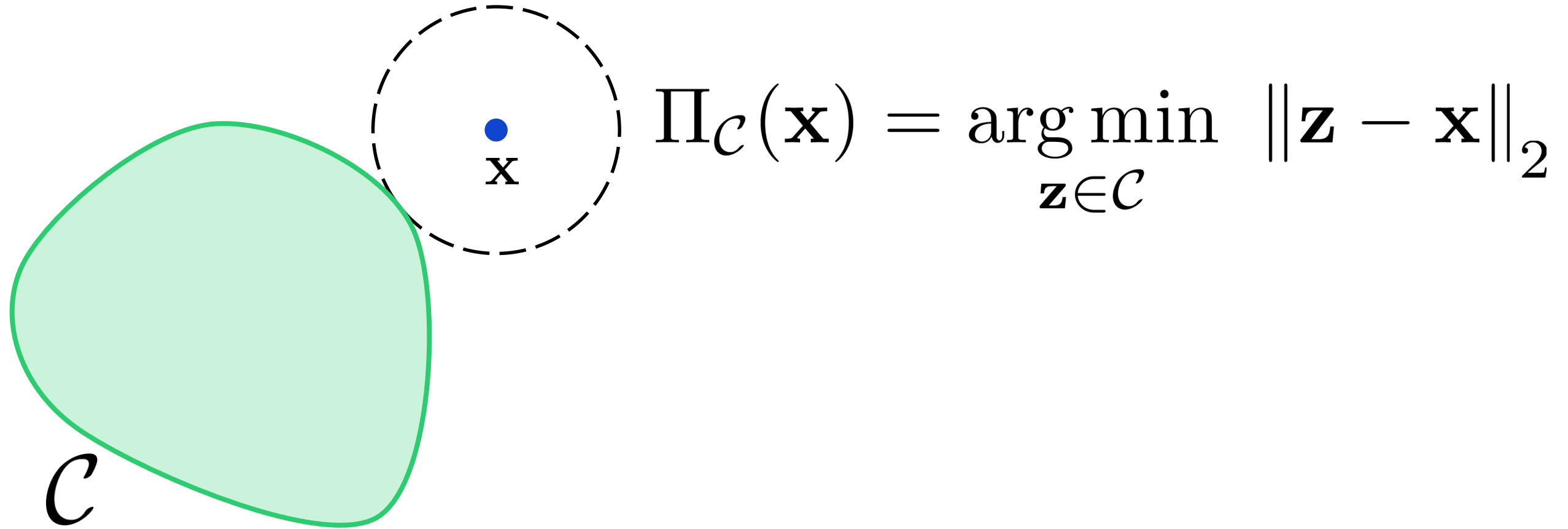


•
 \mathbf{x}

$$\Pi_{\mathcal{C}}(\mathbf{x}) = \arg \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{z} - \mathbf{x}\|_2$$

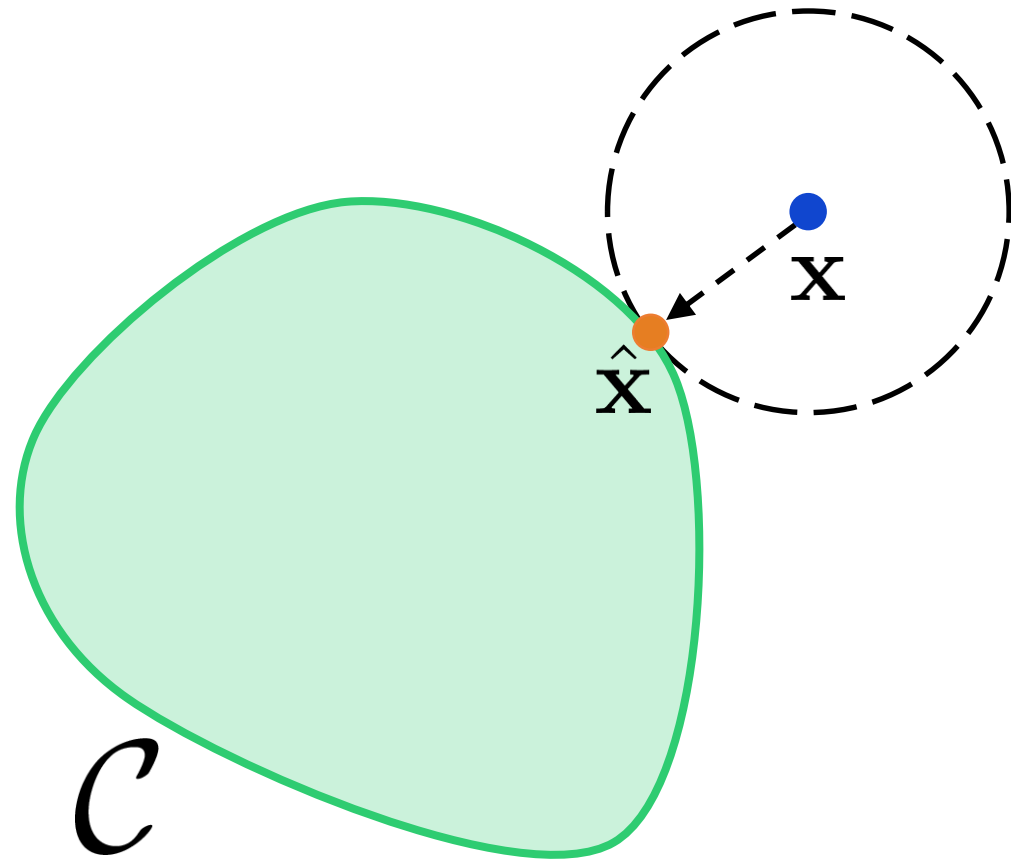
$$\begin{aligned} \arg \min \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & \mathbf{w} \in \mathcal{C} \end{aligned}$$

Projected Gradient Descent



$$\begin{aligned} \arg \min f(\mathbf{w}) \\ \text{s.t. } \mathbf{w} \in \mathcal{C} \end{aligned}$$

Projected Gradient Descent

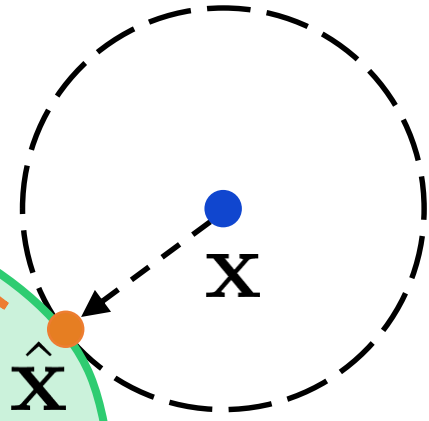


$$\Pi_{\mathcal{C}}(\mathbf{x}) = \arg \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{z} - \mathbf{x}\|_2$$

$$\begin{aligned} \arg \min & f(\mathbf{w}) \\ \text{s.t. } & \mathbf{w} \in \mathcal{C} \end{aligned}$$

Projected Gradient Descent

Projection



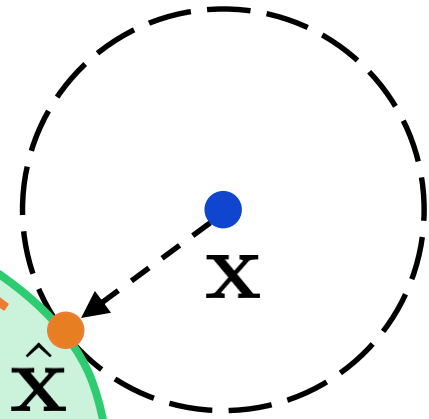
$$\Pi_{\mathcal{C}}(\mathbf{x}) = \arg \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{z} - \mathbf{x}\|_2$$

$$\begin{aligned} \arg \min \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & \mathbf{w} \in \mathcal{C} \end{aligned}$$

Projected Gradient Descent

Doesn't change
the problem

Projection



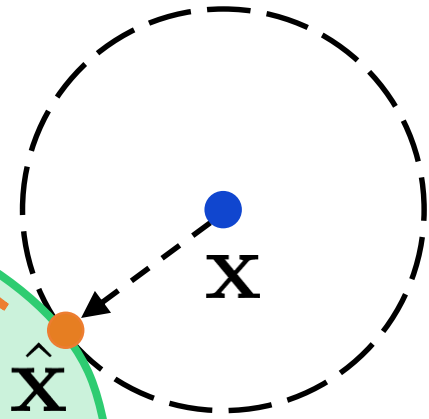
$$\Pi_{\mathcal{C}}(\mathbf{x}) = \arg \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{z} - \mathbf{x}\|_2^2$$

\mathcal{C}

$$\begin{aligned} \arg \min \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & \mathbf{w} \in \mathcal{C} \end{aligned}$$

Projected Gradient Descent

Projection



$$\Pi_{\mathcal{C}}(\mathbf{x}) = \arg \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{z} - \mathbf{x}\|_2^2$$

\mathcal{C}

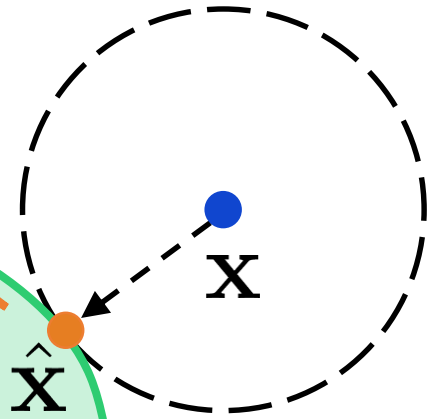
$$\begin{aligned} \arg \min \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & \mathbf{w} \in \mathcal{C} \end{aligned}$$

Projected Gradient Descent

Projection

Projection is just another optimization problem!

$$\Pi_{\mathcal{C}}(\mathbf{x}) = \arg \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{z} - \mathbf{x}\|_2^2$$



\mathcal{C}

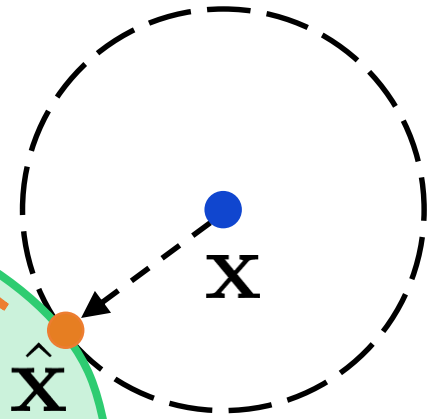
$$\begin{aligned} \arg \min \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & \mathbf{w} \in \mathcal{C} \end{aligned}$$

Projected Gradient Descent

Projection

Projection is just another optimization problem!

$$\Pi_{\mathcal{C}}(\mathbf{x}) = \arg \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{z} - \mathbf{x}\|_2^2$$



$\hat{\mathbf{x}}$

\mathcal{C}

$$\begin{aligned} \arg \min \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & \mathbf{w} \in \mathcal{C} \end{aligned}$$

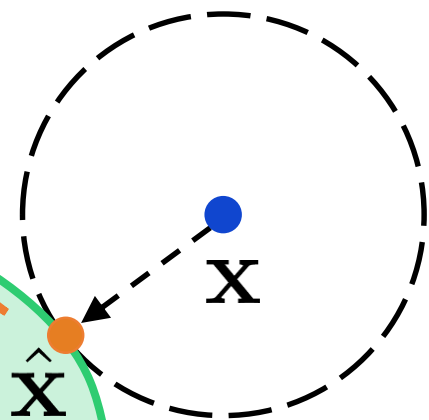
PROJECTED GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Take gradient step $\mathbf{u}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
3. Satisfy constraint $\mathbf{w}^{t+1} \leftarrow \Pi_{\mathcal{C}}(\mathbf{u}^{t+1})$
4. Repeat until convergence

Projected Gradient Descent

Projection is just another optimization problem!

Projection



$$\Pi_C(\mathbf{x}) = \arg \min_{\mathbf{z} \in C} \|\mathbf{z} - \mathbf{x}\|_2^2$$

Projected Stochastic Gradient Descent??

PROJECTED GRADIENT DESCENT

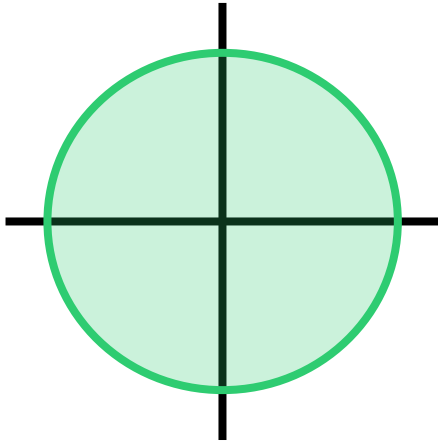
1. Initialize \mathbf{w}^0
2. Take gradient step $\mathbf{u}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
3. Satisfy constraint $\mathbf{w}^{t+1} \leftarrow \Pi_C(\mathbf{u}^{t+1})$
4. Repeat until convergence

$$\begin{aligned} \arg \min & f(\mathbf{w}) \\ \text{s.t. } & \mathbf{w} \in C \end{aligned}$$

A few useful projections

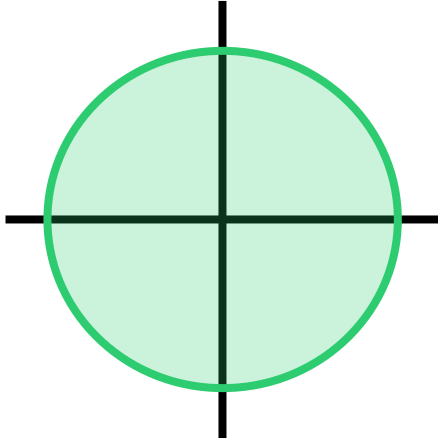
A few useful projections

$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq 1\}$$



A few useful projections

$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq 1\}$$

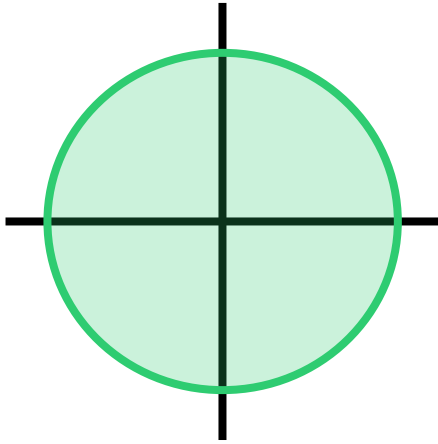


$$\mathcal{C} = \{\mathbf{x} : \mathbf{x}_i \geq 0\}$$

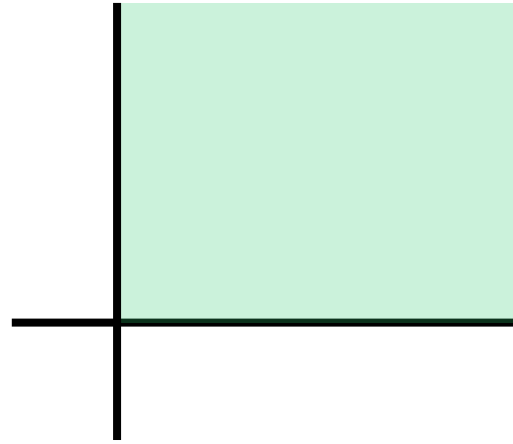


A few useful projections

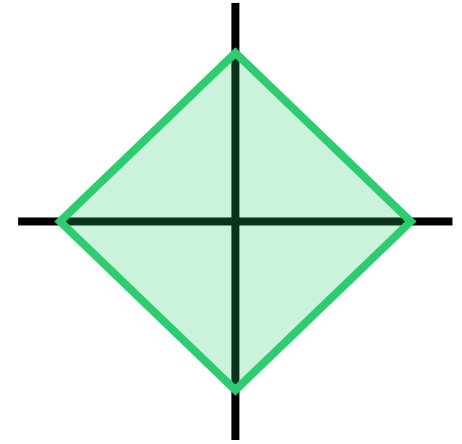
$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq 1\}$$



$$\mathcal{C} = \{\mathbf{x} : \mathbf{x}_i \geq 0\}$$

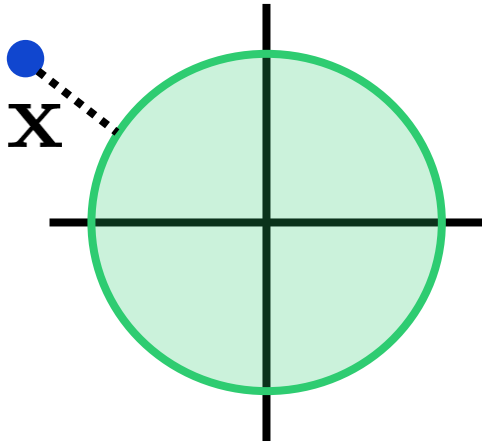


$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_1 \leq 1\}$$

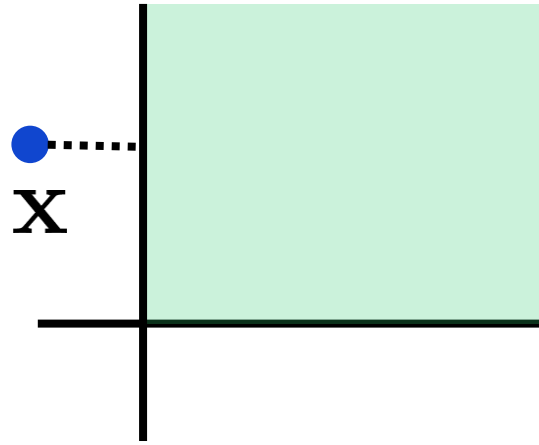


A few useful projections

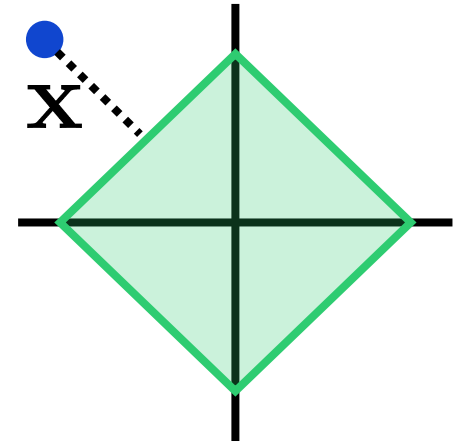
$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq 1\}$$



$$\mathcal{C} = \{\mathbf{x} : \mathbf{x}_i \geq 0\}$$



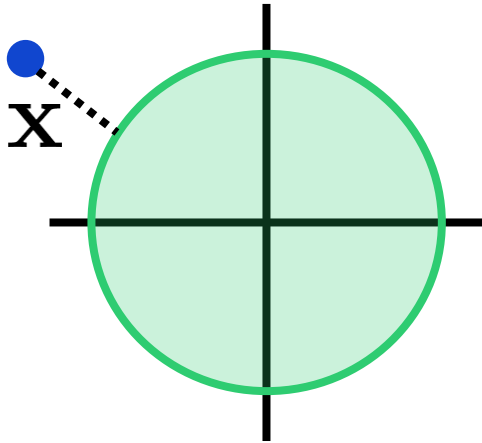
$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_1 \leq 1\}$$



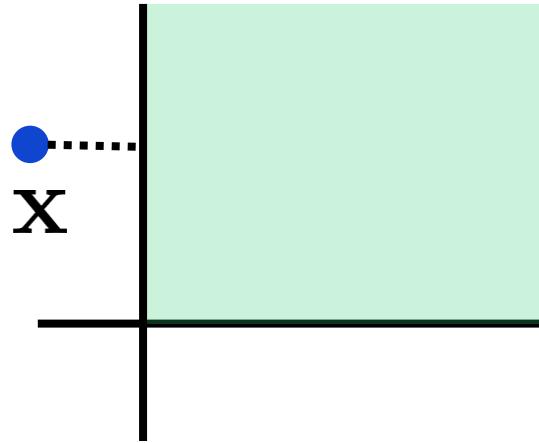
A few useful projections

$$\hat{\mathbf{x}} = \Pi_{\mathcal{C}}(\mathbf{x})$$

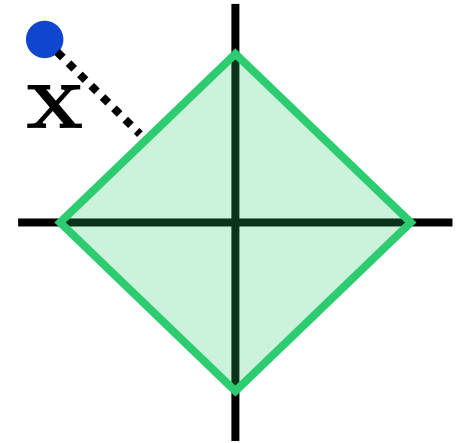
$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq 1\}$$



$$\mathcal{C} = \{\mathbf{x} : \mathbf{x}_i \geq 0\}$$



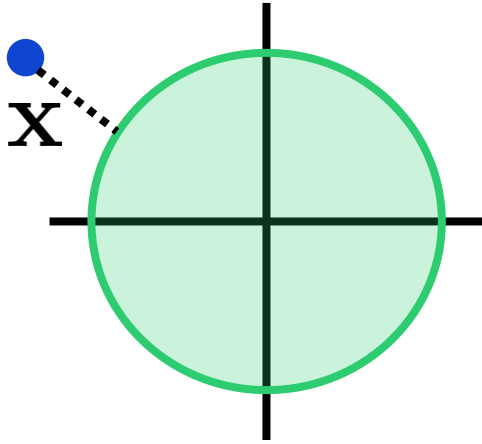
$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_1 \leq 1\}$$



A few useful projections

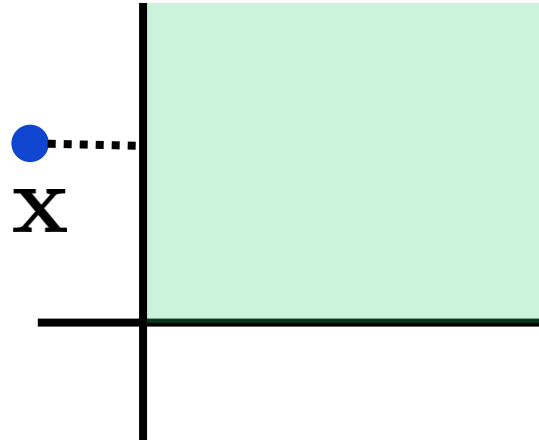
$$\hat{\mathbf{x}} = \Pi_{\mathcal{C}}(\mathbf{x})$$

$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq 1\}$$

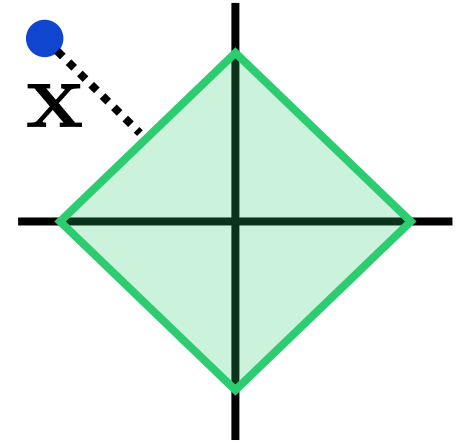


$$\hat{\mathbf{x}} = \begin{cases} \mathbf{x} & \text{if } \|\mathbf{x}\|_2 \leq 1 \\ \frac{\mathbf{x}}{\|\mathbf{x}\|_2} & \text{if } \|\mathbf{x}\|_2 > 1 \end{cases}$$

$$\mathcal{C} = \{\mathbf{x} : \mathbf{x}_i \geq 0\}$$



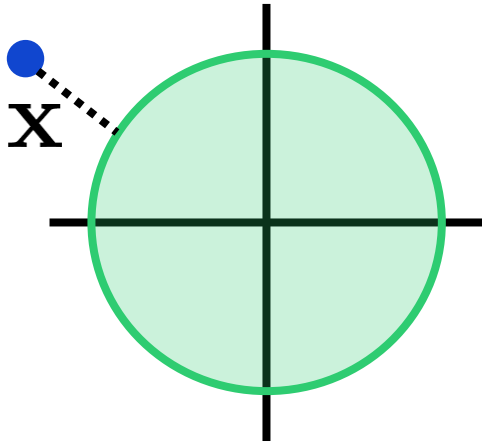
$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_1 \leq 1\}$$



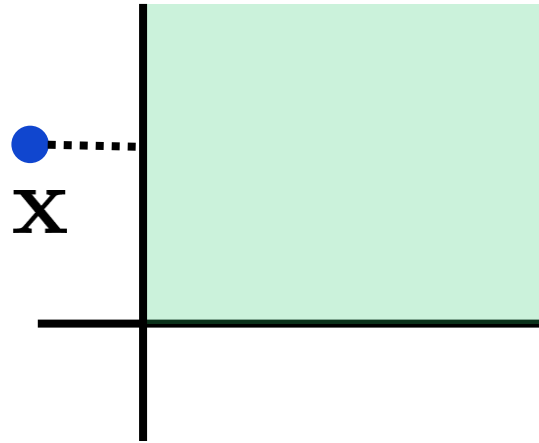
A few useful projections

$$\hat{\mathbf{x}} = \Pi_{\mathcal{C}}(\mathbf{x})$$

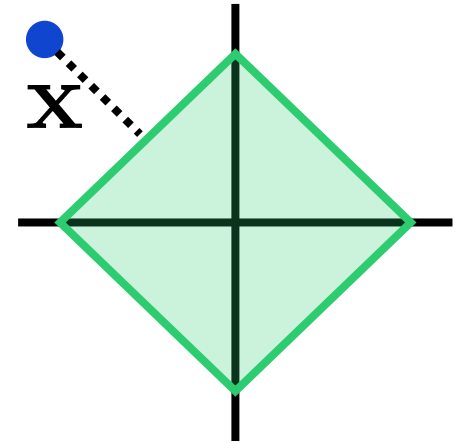
$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq 1\}$$



$$\mathcal{C} = \{\mathbf{x} : \mathbf{x}_i \geq 0\}$$



$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_1 \leq 1\}$$



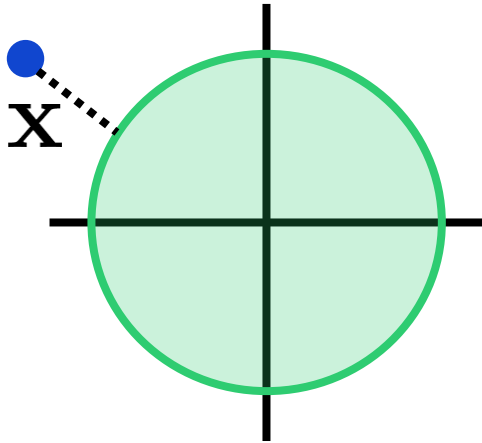
$$\hat{\mathbf{x}} = \begin{cases} \mathbf{x} & \text{if } \|\mathbf{x}\|_2 \leq 1 \\ \frac{\mathbf{x}}{\|\mathbf{x}\|_2} & \text{if } \|\mathbf{x}\|_2 > 1 \end{cases}$$

What if $\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq r\}$?

A few useful projections

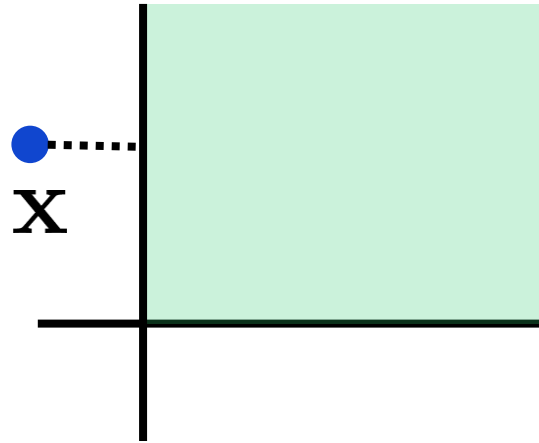
$$\hat{\mathbf{x}} = \Pi_{\mathcal{C}}(\mathbf{x})$$

$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq 1\}$$



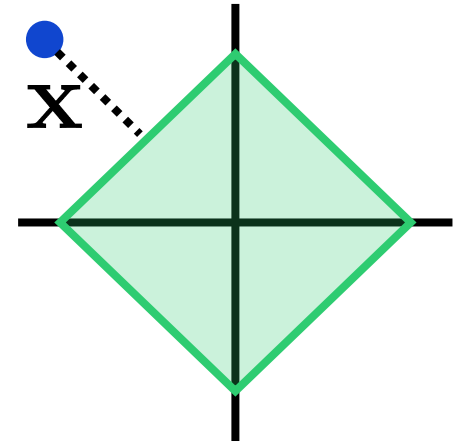
$$\hat{\mathbf{x}} = \begin{cases} \mathbf{x} & \text{if } \|\mathbf{x}\|_2 \leq 1 \\ \frac{\mathbf{x}}{\|\mathbf{x}\|_2} & \text{if } \|\mathbf{x}\|_2 > 1 \end{cases}$$

$$\mathcal{C} = \{\mathbf{x} : \mathbf{x}_i \geq 0\}$$



$$\hat{\mathbf{x}}_i = \begin{cases} \mathbf{x}_i & \text{if } \mathbf{x}_i \geq 0 \\ 0 & \text{if } \mathbf{x}_i < 0 \end{cases}$$

$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_1 \leq 1\}$$

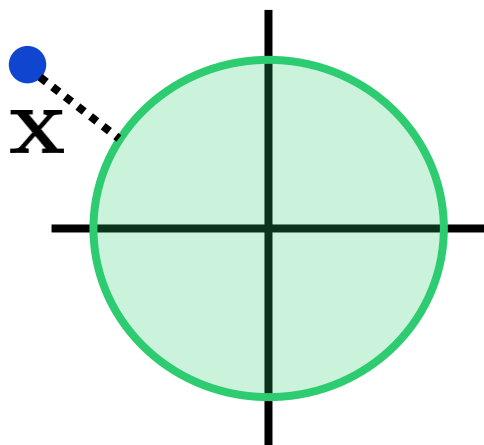


What if $\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq r\}$?

A few useful projections

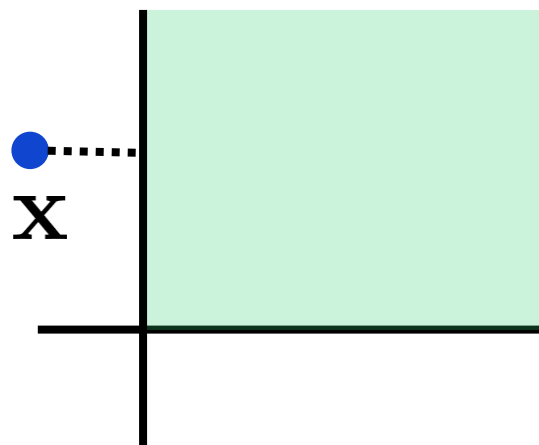
$$\hat{\mathbf{x}} = \Pi_{\mathcal{C}}(\mathbf{x})$$

$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq 1\}$$



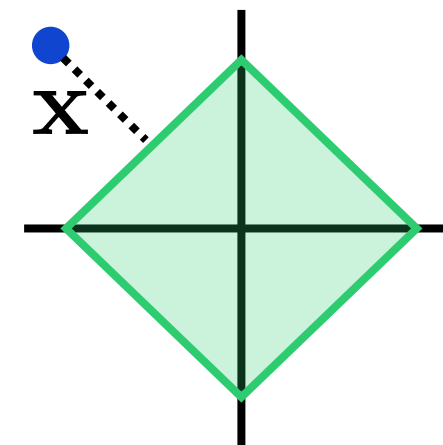
$$\hat{\mathbf{x}} = \begin{cases} \mathbf{x} & \text{if } \|\mathbf{x}\|_2 \leq 1 \\ \frac{\mathbf{x}}{\|\mathbf{x}\|_2} & \text{if } \|\mathbf{x}\|_2 > 1 \end{cases}$$

$$\mathcal{C} = \{\mathbf{x} : \mathbf{x}_i \geq 0\}$$



$$\hat{\mathbf{x}}_i = \begin{cases} \mathbf{x}_i & \text{if } \mathbf{x}_i \geq 0 \\ 0 & \text{if } \mathbf{x}_i < 0 \end{cases}$$

$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_1 \leq 1\}$$



Find out!

What if $\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq r\}$?

Some disclaimers

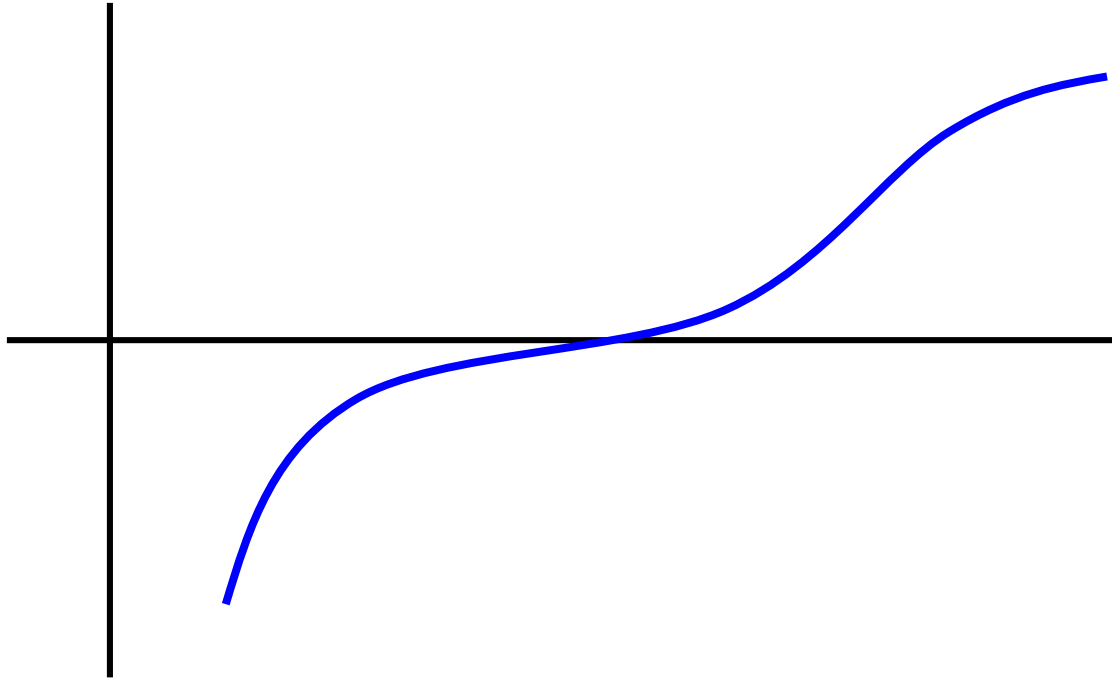
- GD can be slow in convergence
 - Step size tuning can affect performance quite a bit – Armijo, Adagrad
 - Momentum-based methods $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t - \gamma \cdot (\mathbf{w}^t - \mathbf{w}^{t-1})$
 - Second order methods – expensive but super fast!
- Exist superior methods for specific problems
 - RR (L_2 regularized ls regression) – conjugate gradient method
 - Lasso (L_1 regularized ls regression) – proximal gradient method
 - Logistic regression – trust region Newton method
- Details, convergence analysis beyond scope of this course
- CS774: Optimization Techniques
- Next: Newton method, duality, coordinate descent

Momentum weight

Newton Method

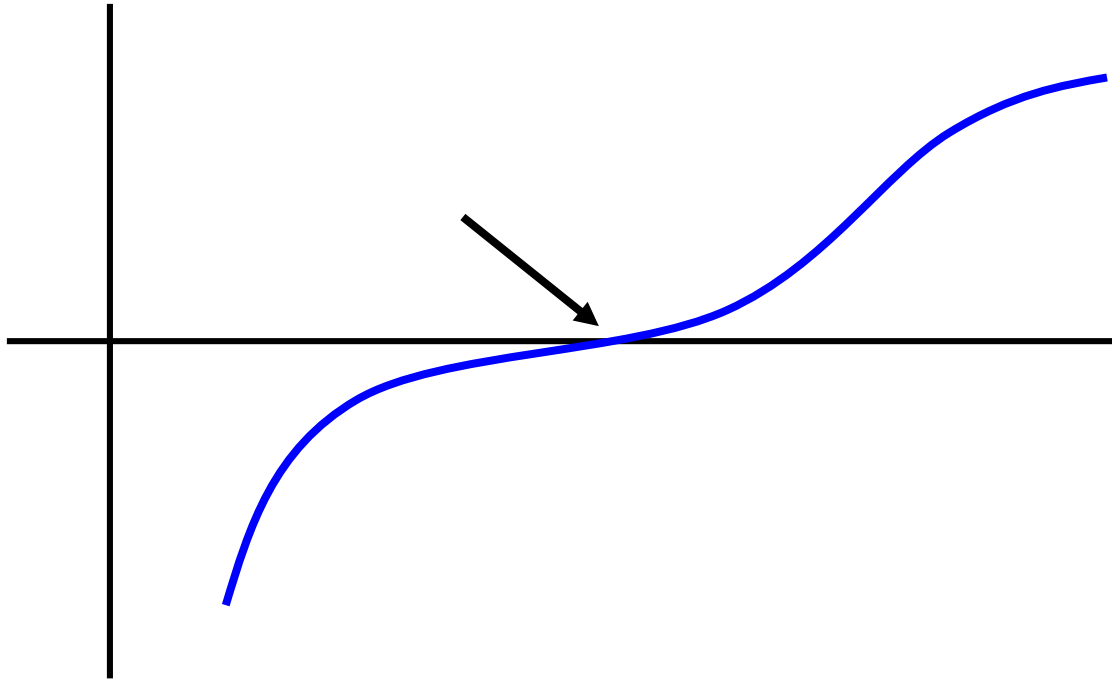
Descent with second order derivatives

Newton-Raphson Method for Root Finding

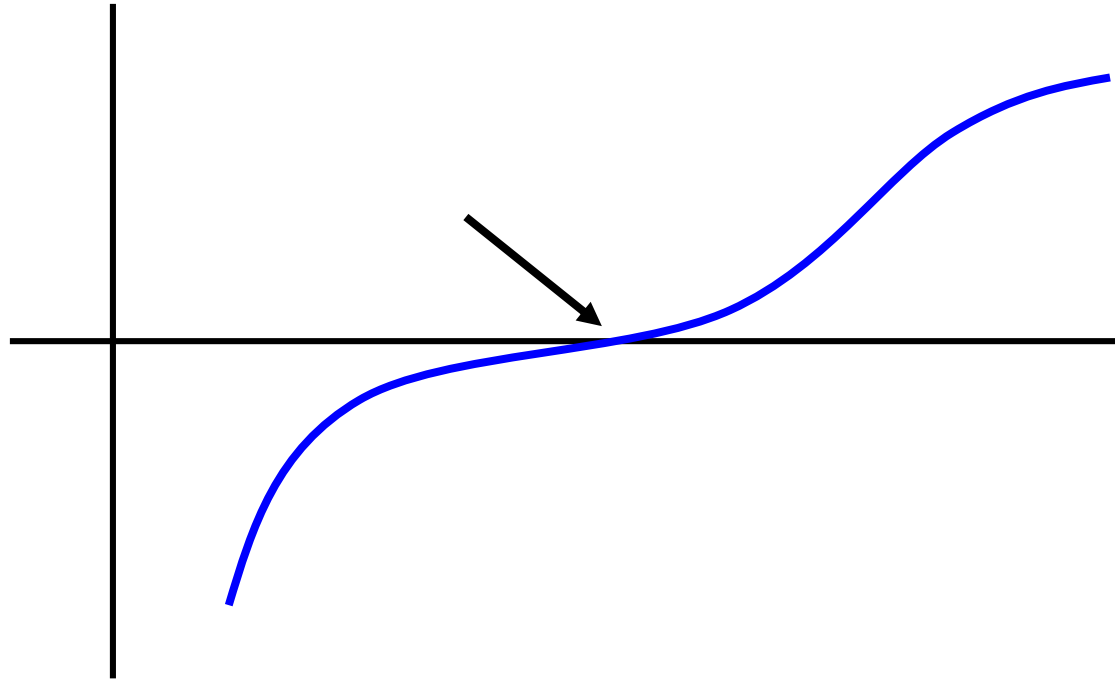


Newton-Raphson Method for Root Finding

$$x : f(x) = 0$$



Newton-Raphson Method for Root Finding

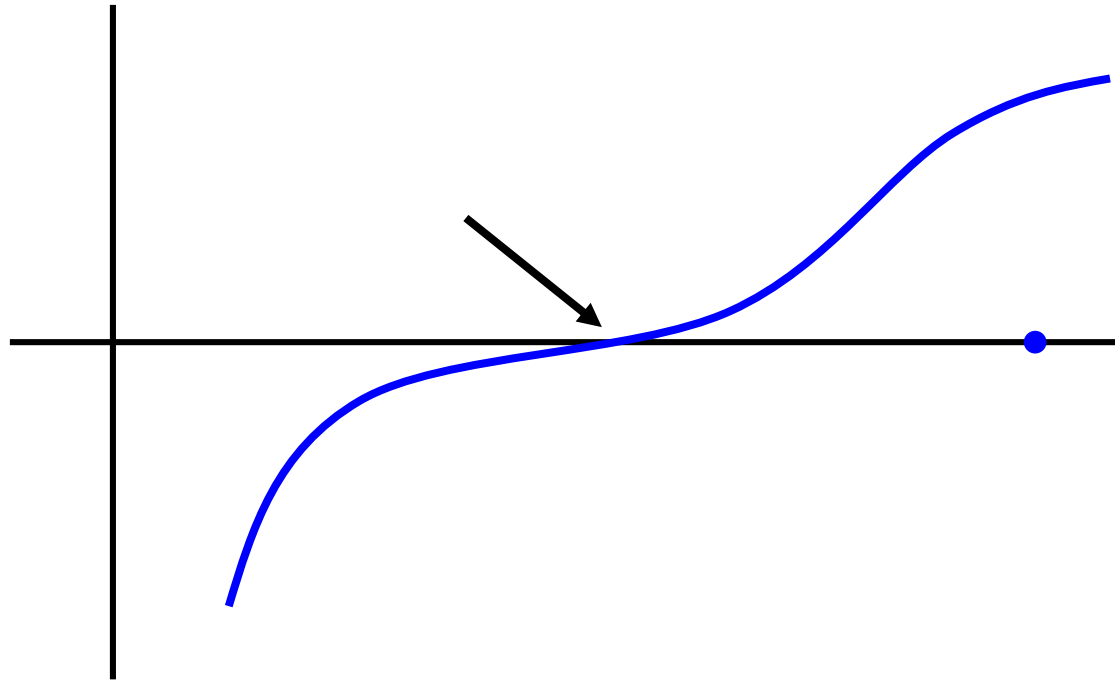


$$x : f(x) = 0$$

Finding roots of linear functions is easy

$$f(x) = ax + b, x_0 = \frac{-b}{a}$$

Newton-Raphson Method for Root Finding

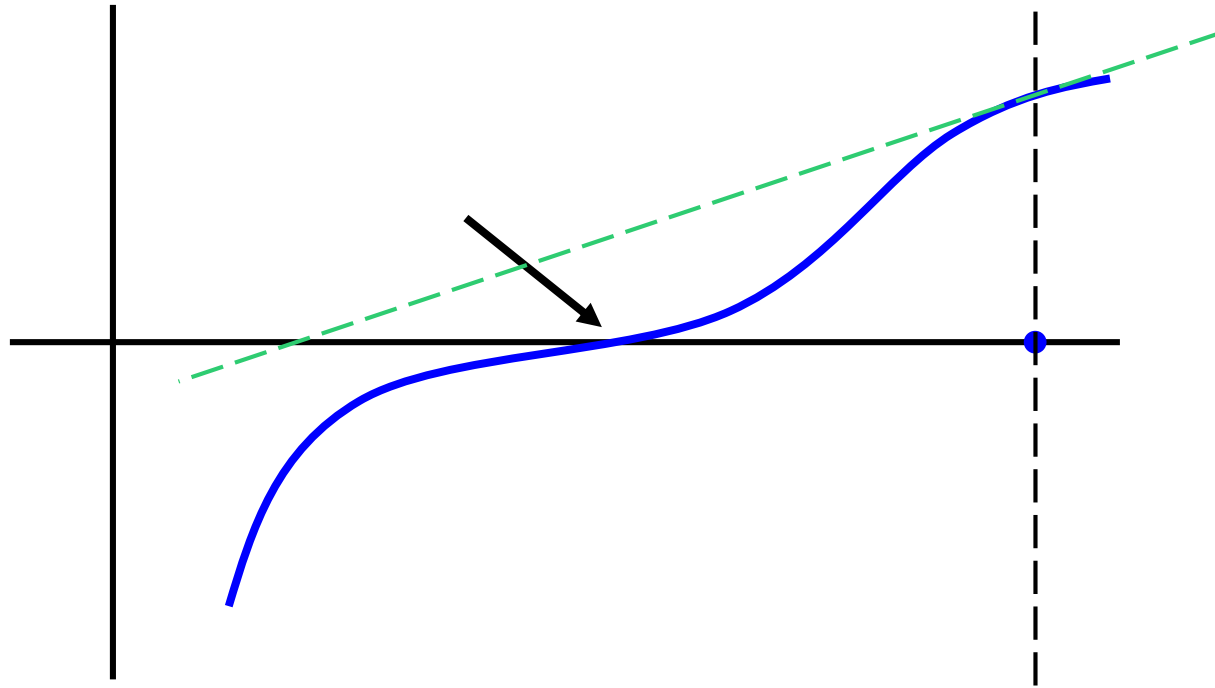


$$x : f(x) = 0$$

Finding roots of linear functions is easy

$$f(x) = ax + b, x_0 = \frac{-b}{a}$$

Newton-Raphson Method for Root Finding

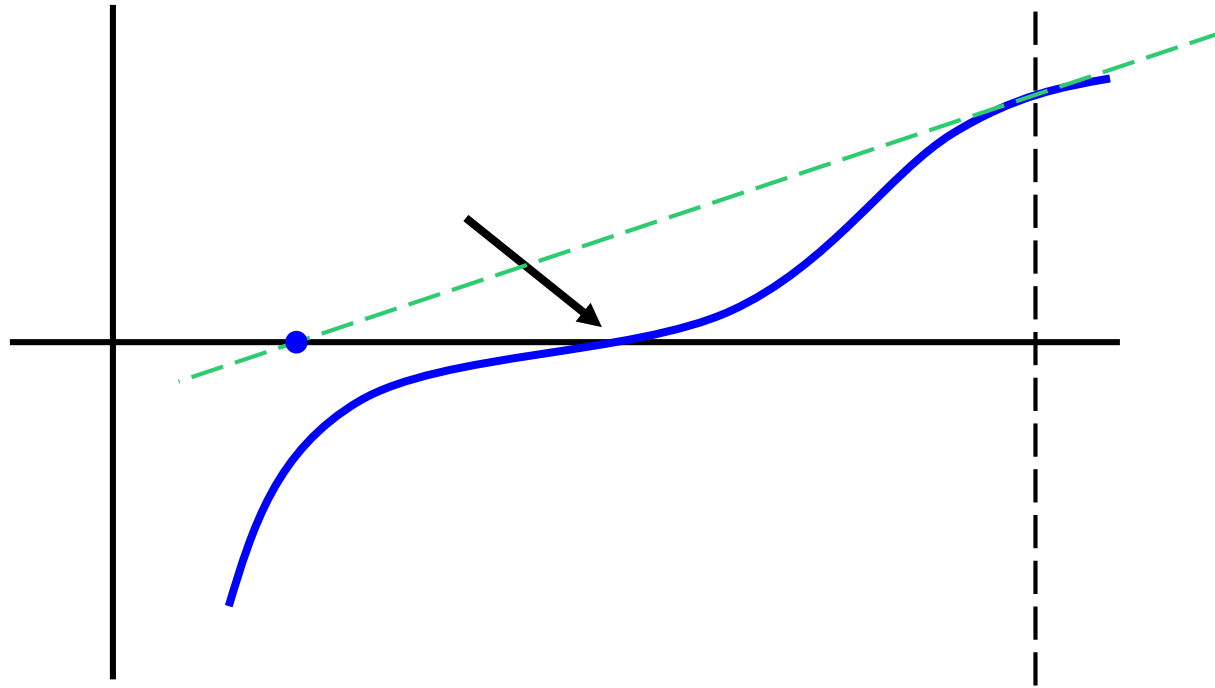


$$x : f(x) = 0$$

Finding roots of linear functions is easy

$$f(x) = ax + b, x_0 = \frac{-b}{a}$$

Newton-Raphson Method for Root Finding

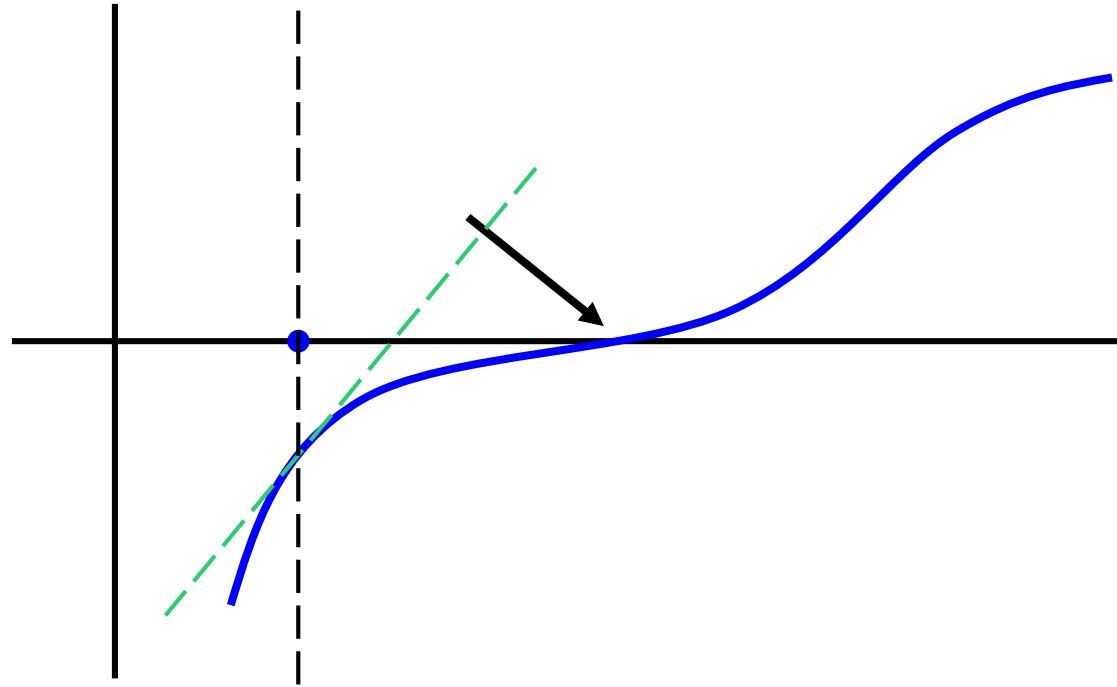


$$x : f(x) = 0$$

Finding roots of linear functions is easy

$$f(x) = ax + b, x_0 = \frac{-b}{a}$$

Newton-Raphson Method for Root Finding

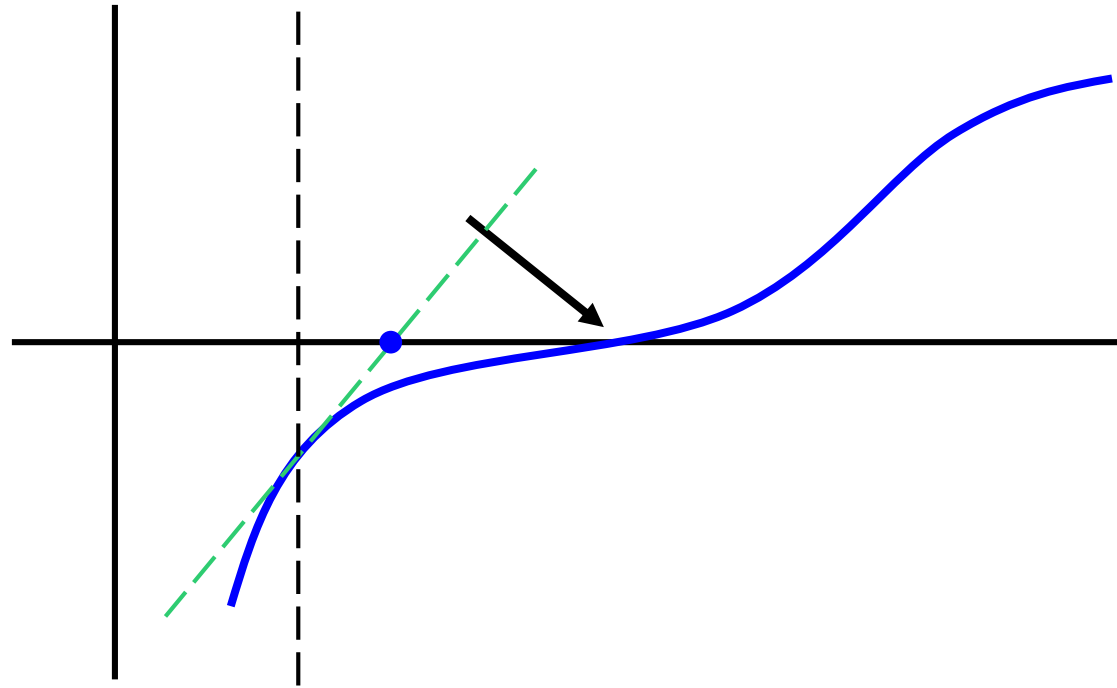


$$x : f(x) = 0$$

Finding roots of linear functions is easy

$$f(x) = ax + b, x_0 = \frac{-b}{a}$$

Newton-Raphson Method for Root Finding

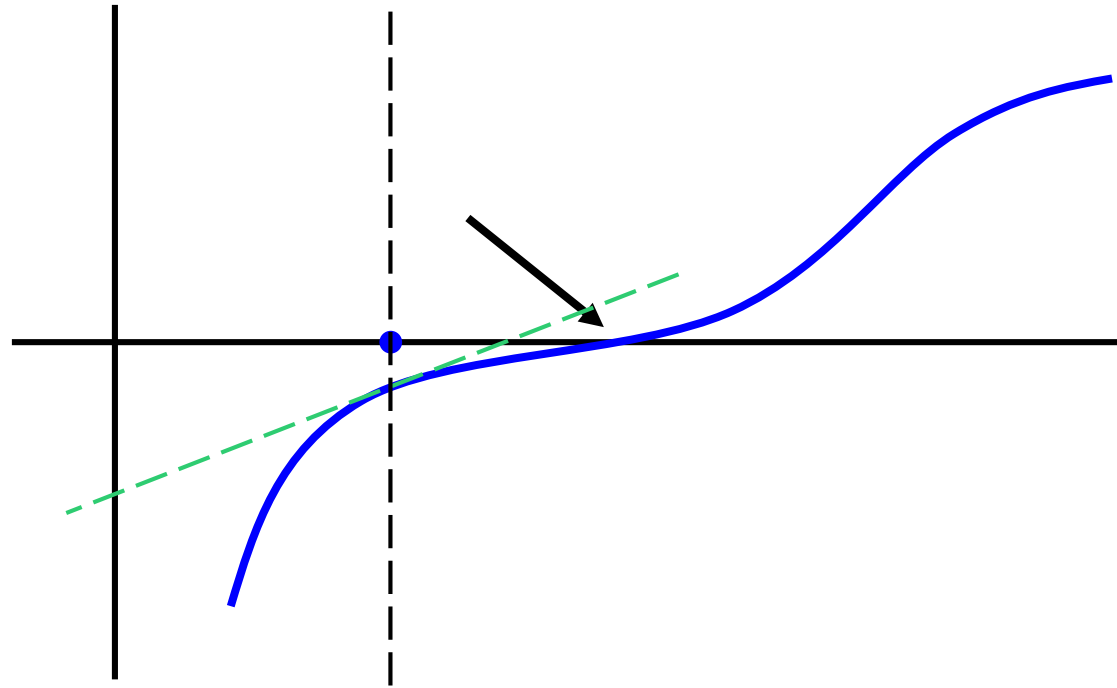


$$x : f(x) = 0$$

Finding roots of linear functions is easy

$$f(x) = ax + b, x_0 = \frac{-b}{a}$$

Newton-Raphson Method for Root Finding

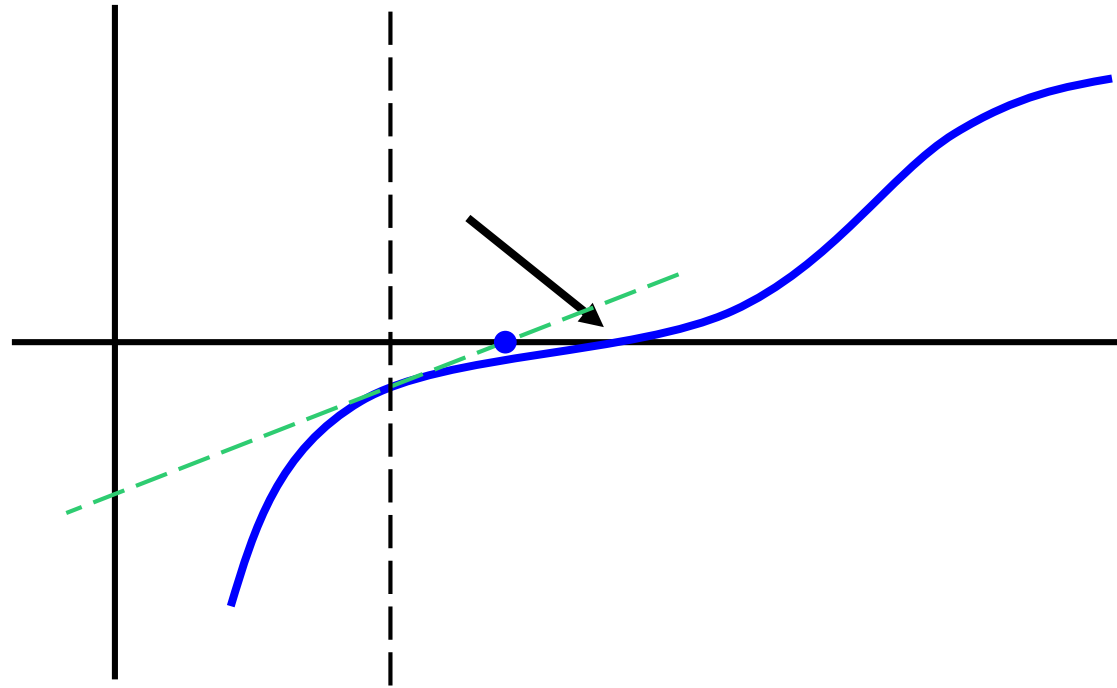


$$x : f(x) = 0$$

Finding roots of linear functions is easy

$$f(x) = ax + b, x_0 = \frac{-b}{a}$$

Newton-Raphson Method for Root Finding

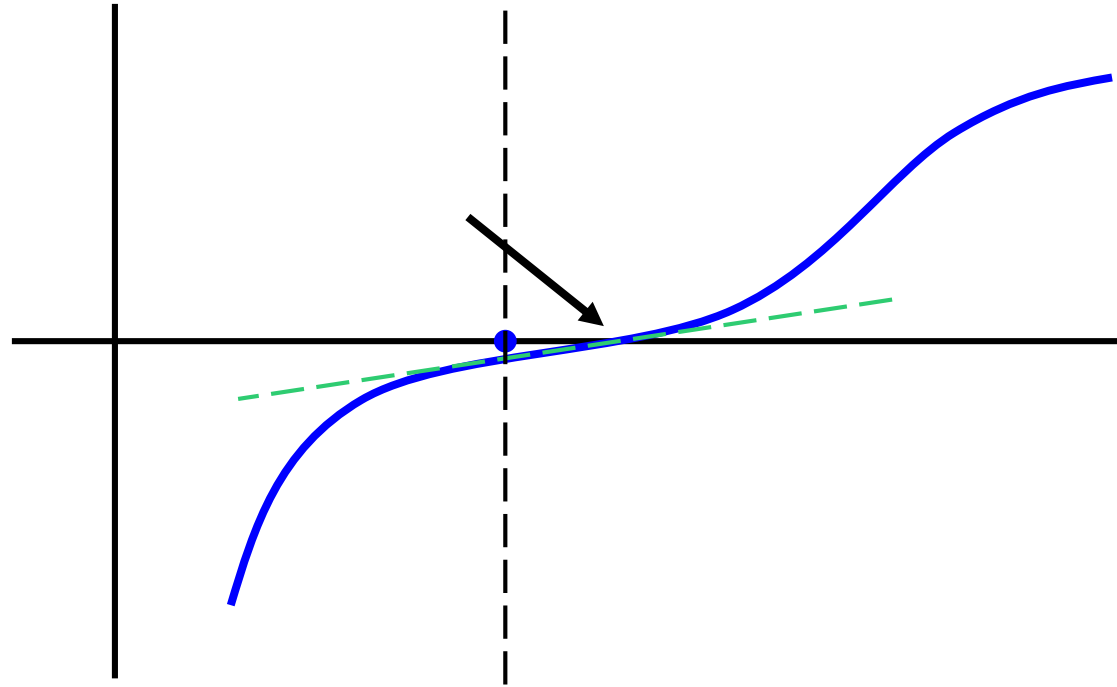


$$x : f(x) = 0$$

Finding roots of linear functions is easy

$$f(x) = ax + b, x_0 = \frac{-b}{a}$$

Newton-Raphson Method for Root Finding

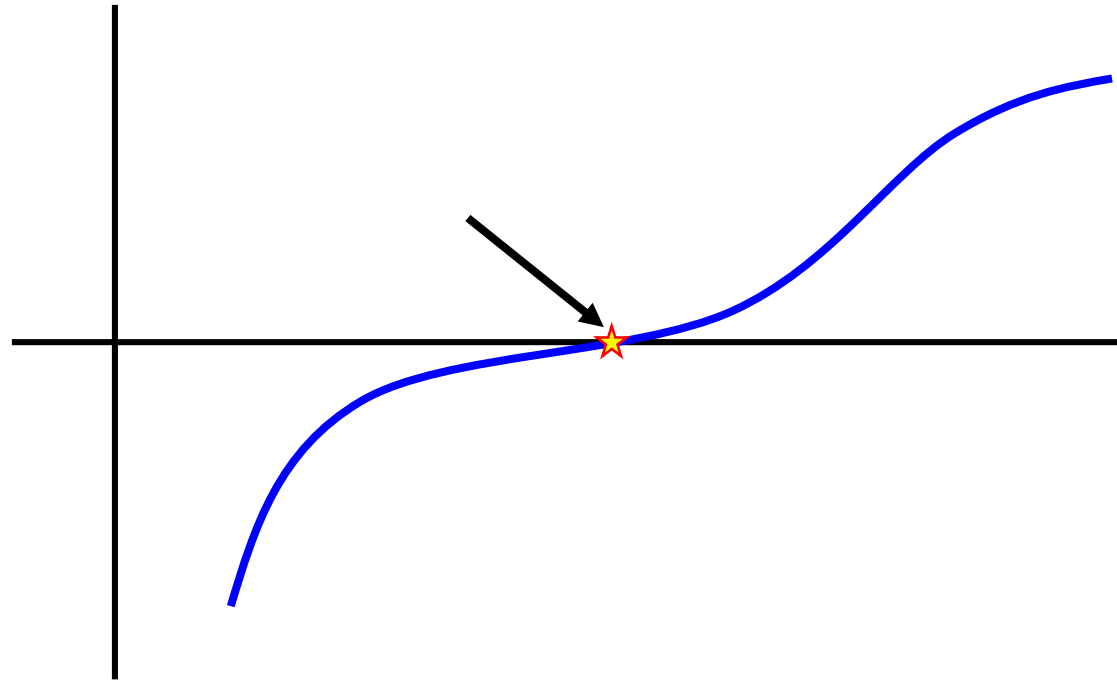


$$x : f(x) = 0$$

Finding roots of linear functions is easy

$$f(x) = ax + b, x_0 = \frac{-b}{a}$$

Newton-Raphson Method for Root Finding

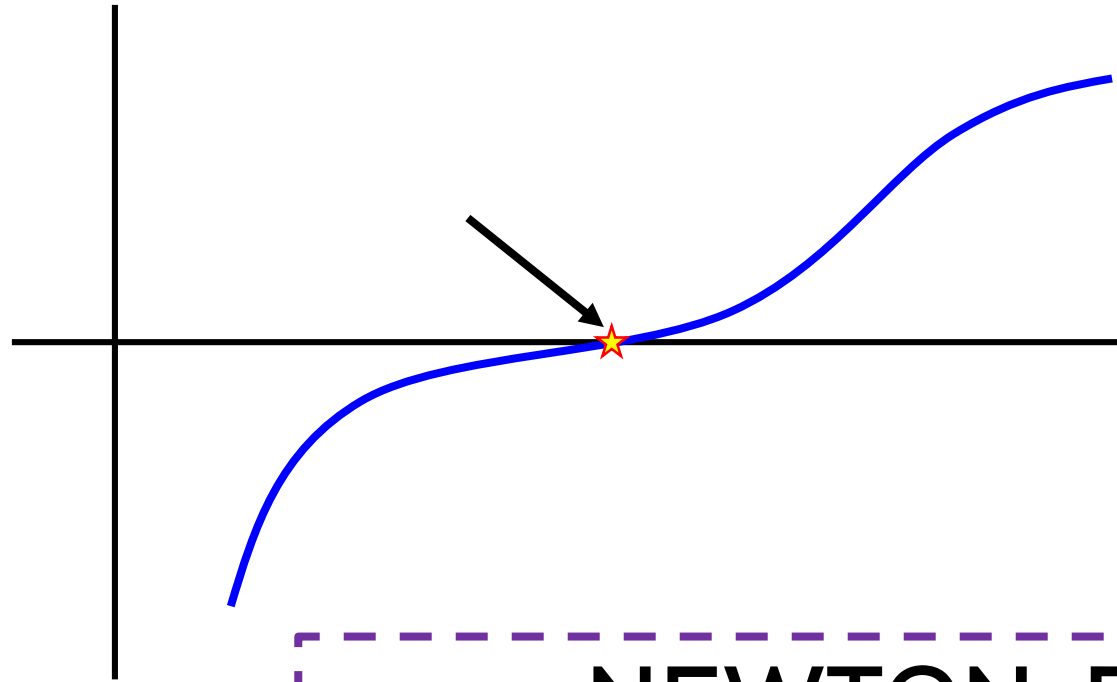


$$x : f(x) = 0$$

Finding roots of linear functions is easy

$$f(x) = ax + b, x_0 = \frac{-b}{a}$$

Newton-Raphson Method for Root Finding



$$x : f(x) = 0$$

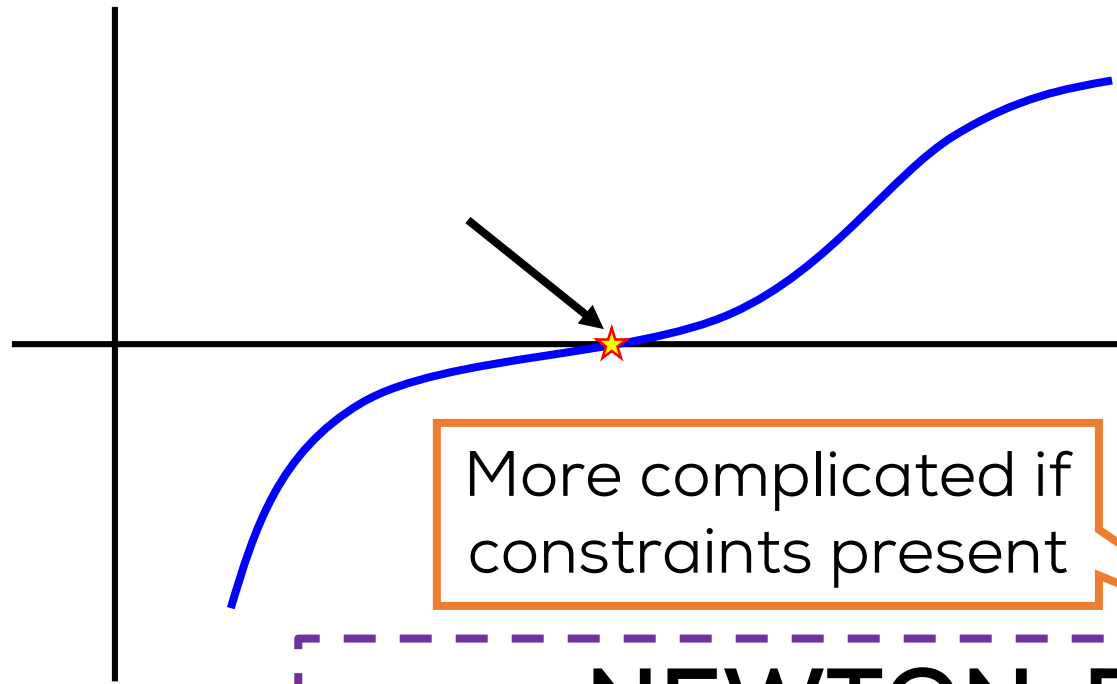
Finding roots of linear functions is easy

$$f(x) = ax + b, x_0 = \frac{-b}{a}$$

NEWTON-RAPHSON METHOD

1. Initialize x^0
2. Approximate $f(\cdot)$ by $\tilde{f}_t(x) = f(x^t) + \langle f'(x^t), x - x^t \rangle$
3. Update $x^{t+1} \leftarrow \text{ROOT}(\tilde{f}_t) = x^t - \frac{f(x^t)}{f'(x^t)}$
4. Repeat until convergence

Newton-Raphson Method for Root Finding



$$x : f(x) = 0$$

Finding roots of linear functions is easy

$$f(x) = ax + b, x_0 = \frac{-b}{a}$$

NEWTON-RAPHSON METHOD

1. Initialize x^0
2. Approximate $f(\cdot)$ by $\tilde{f}_t(x) = f(x^t) + \langle f'(x^t), x - x^t \rangle$
3. Update $x^{t+1} \leftarrow \text{ROOT}(\tilde{f}_t) = x^t - \frac{f(x^t)}{f'(x^t)}$
4. Repeat until convergence

Newton-Raphson Method for Root Finding

$$x : f(x) = 0$$

Finding roots of linear functions is easy

$$f(x) = ax + b, x_0 = \frac{-b}{a}$$

More complicated if constraints present

NEWTON-RAPHSON METHOD

No guarantees in general

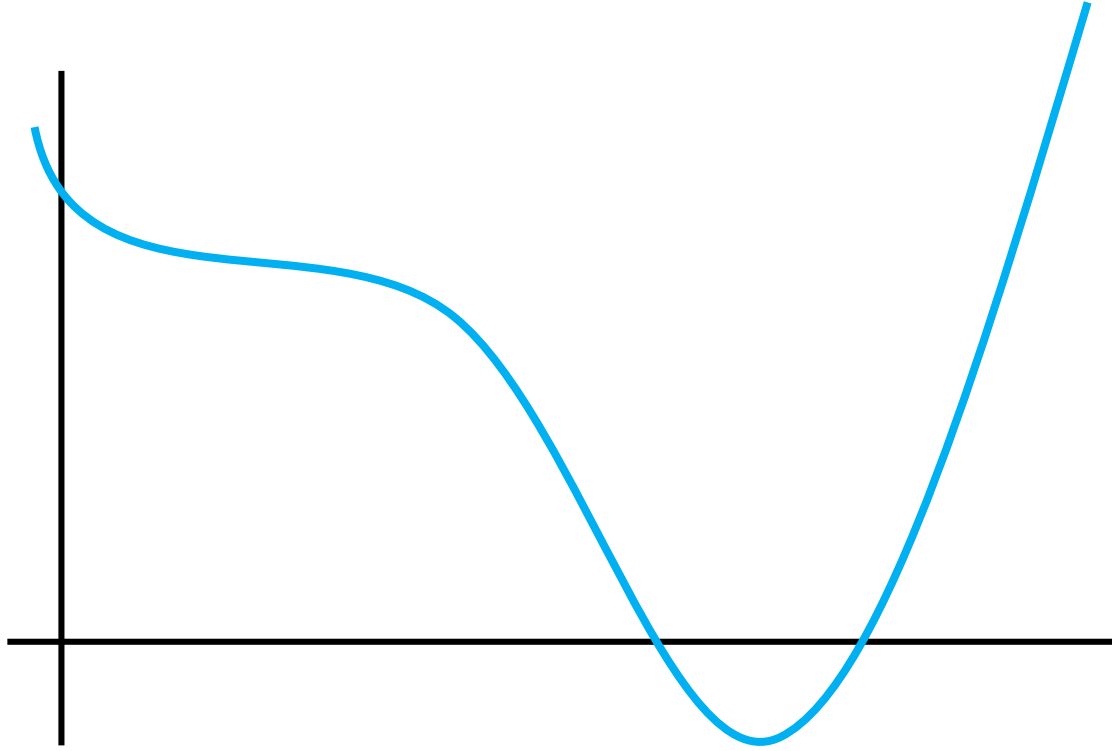
1. Initialize x^0
2. Approximate $f(\cdot)$ by $\tilde{f}_t(x) = f(x^t) + \langle f'(x^t), x - x^t \rangle$
3. Update $x^{t+1} \leftarrow \text{ROOT}(\tilde{f}_t) = x^t - \frac{f(x^t)}{f'(x^t)}$
4. Repeat until convergence

Newton Method

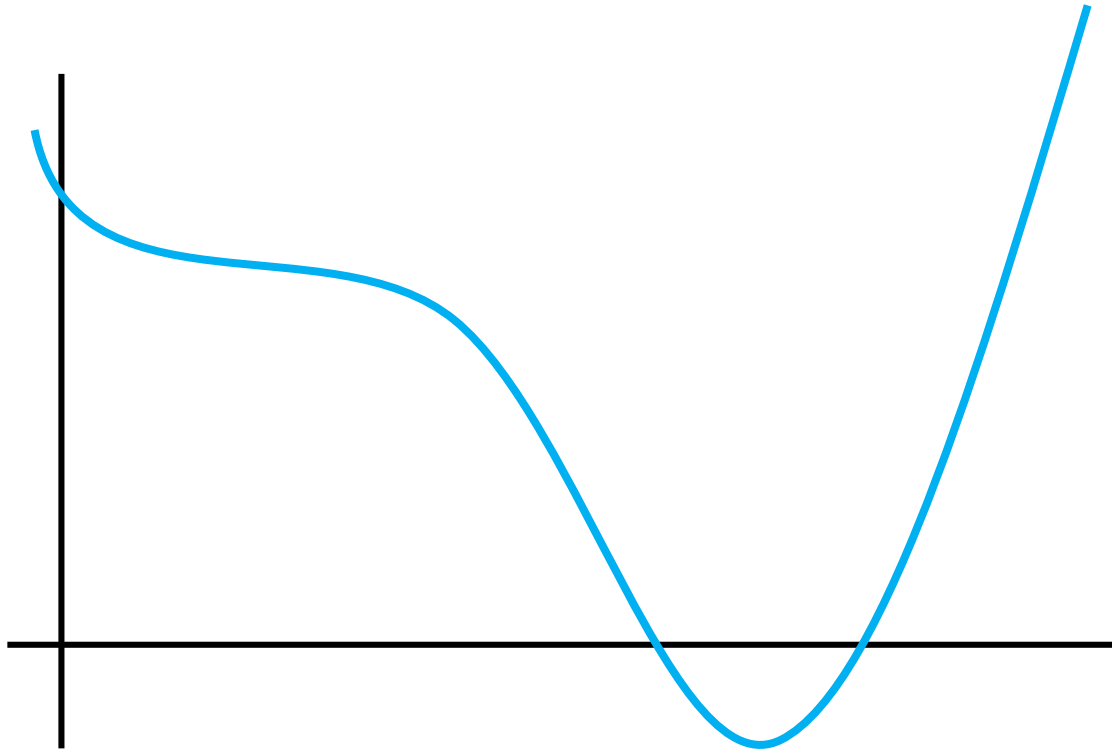
September 1, 2017



Newton Method

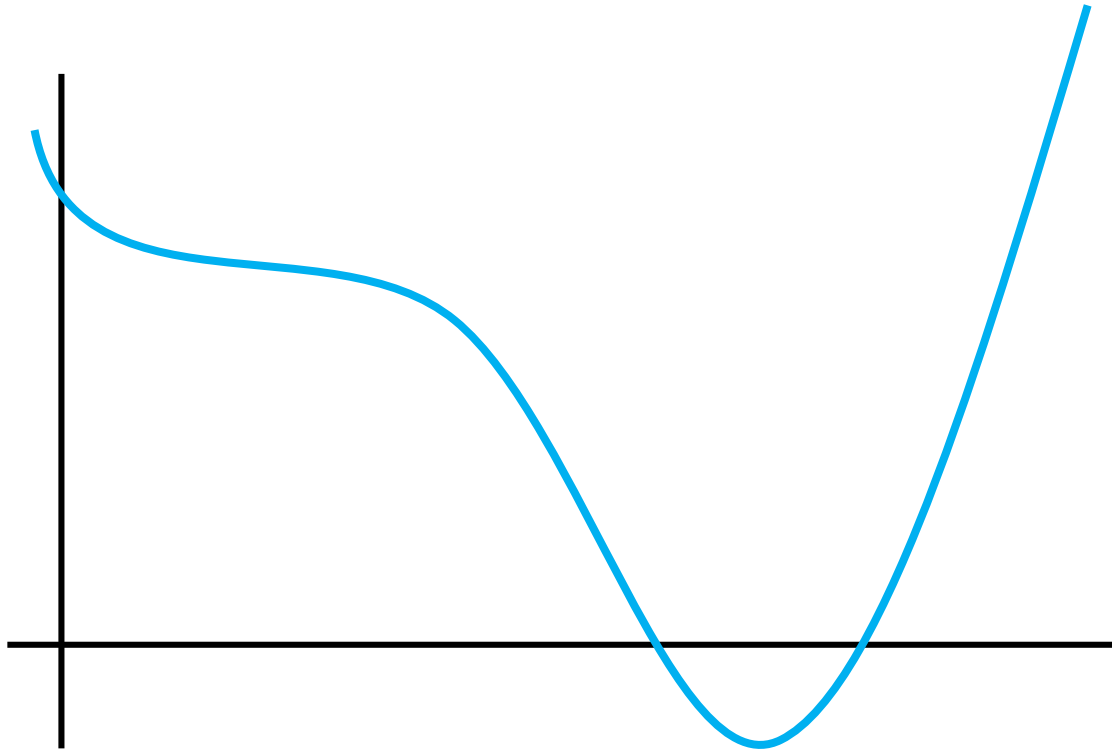


Newton Method



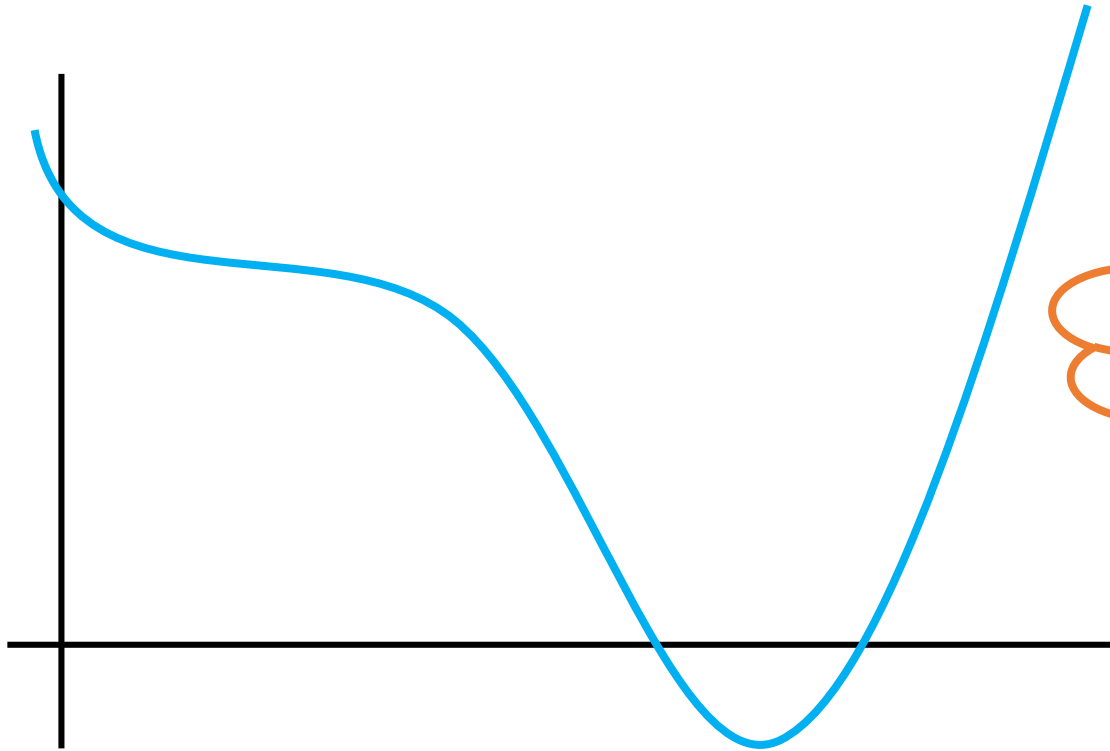
$$\mathbf{x} \in \arg \min f(\mathbf{x})$$

Newton Method



$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$

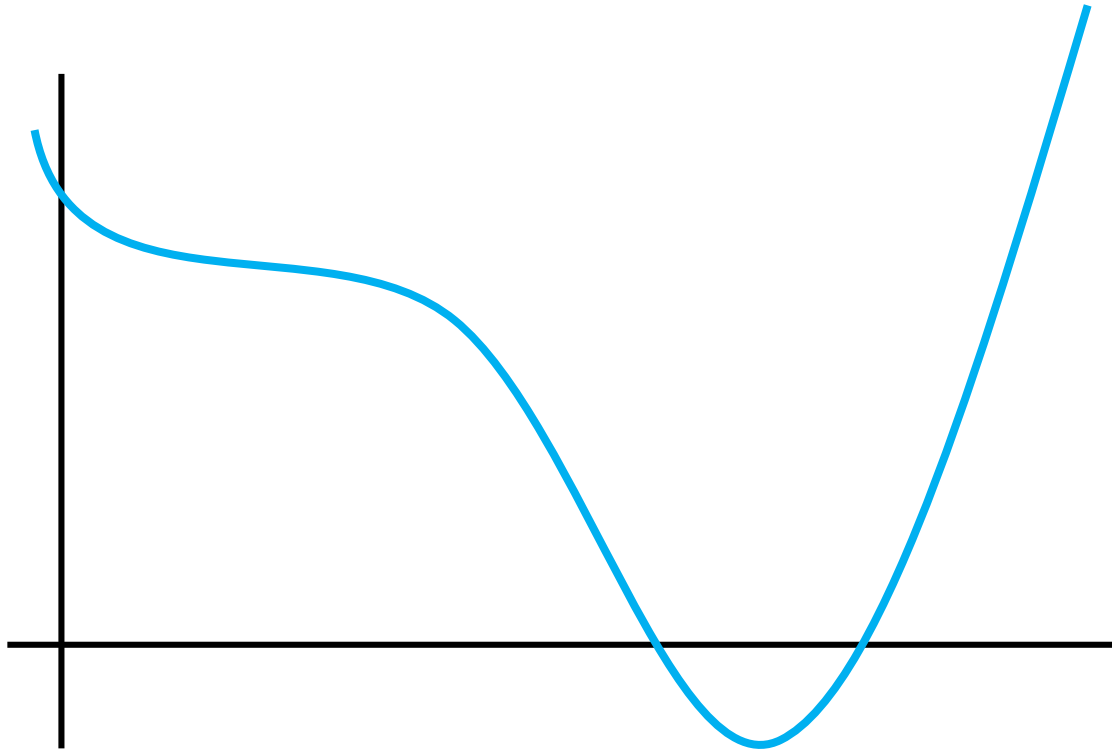
Newton Method



$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$

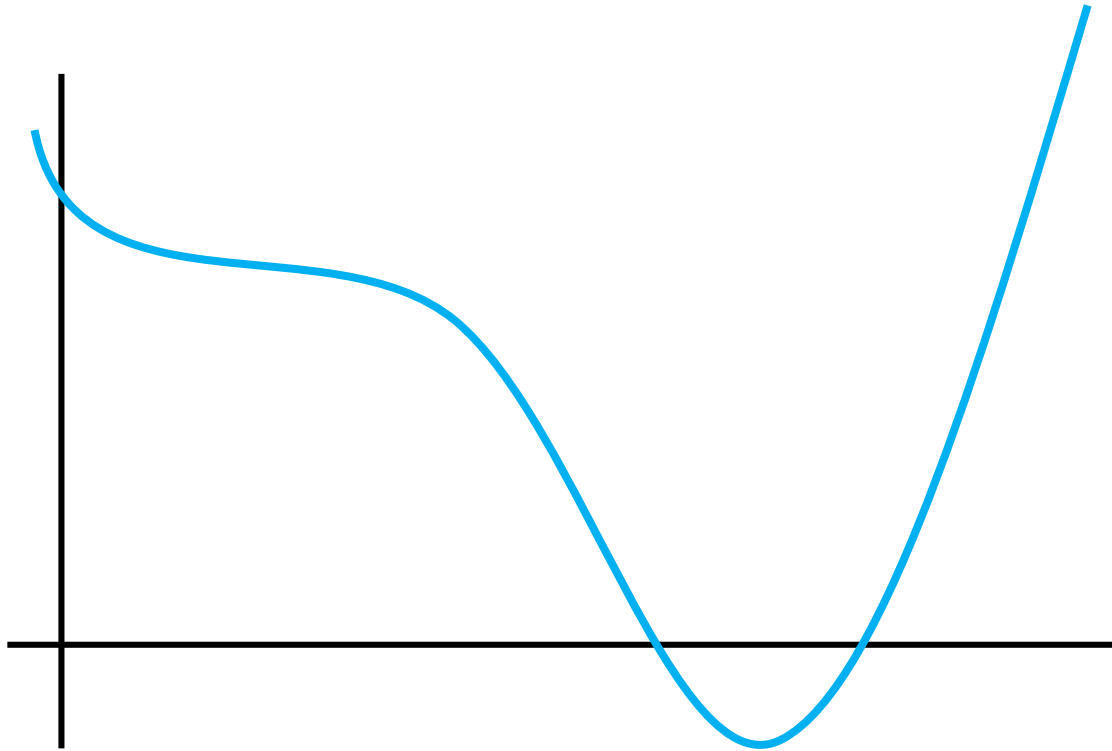
Instead of finding minima,
search for stable points
instead i.e. $\nabla f(\mathbf{x}) = \mathbf{0}$

Newton Method



$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$

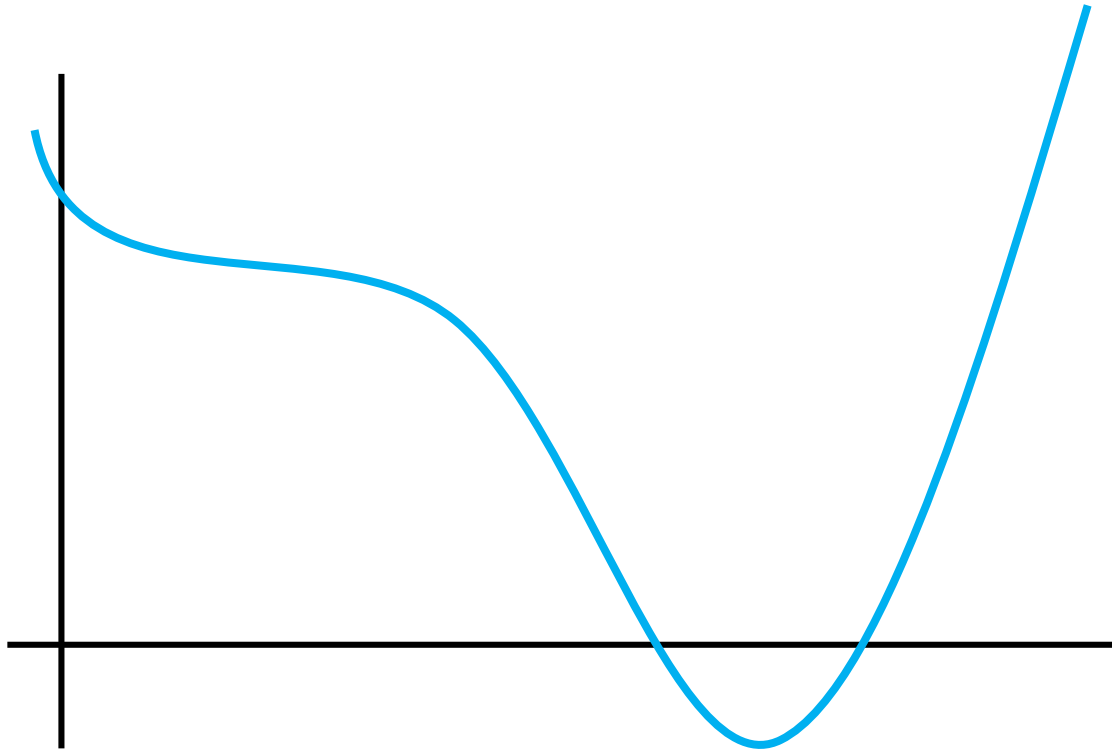
Newton Method



$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$

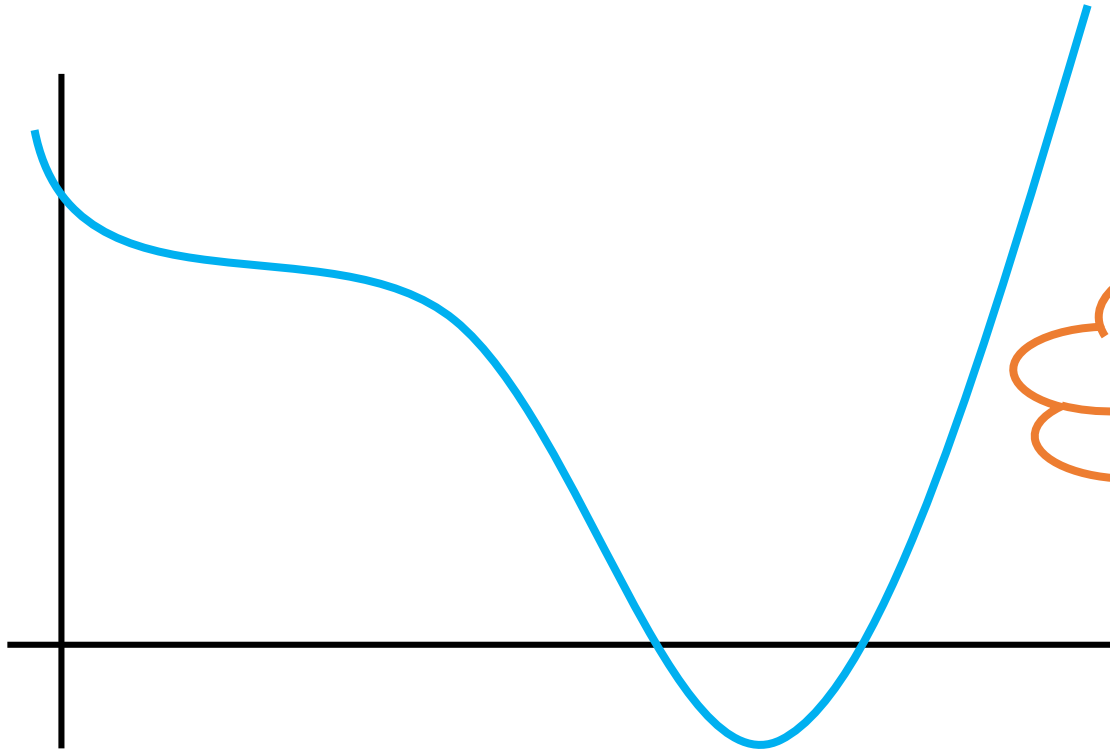
What is the
gradient of $\nabla f(\cdot)$?

Newton Method



$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$

Newton Method

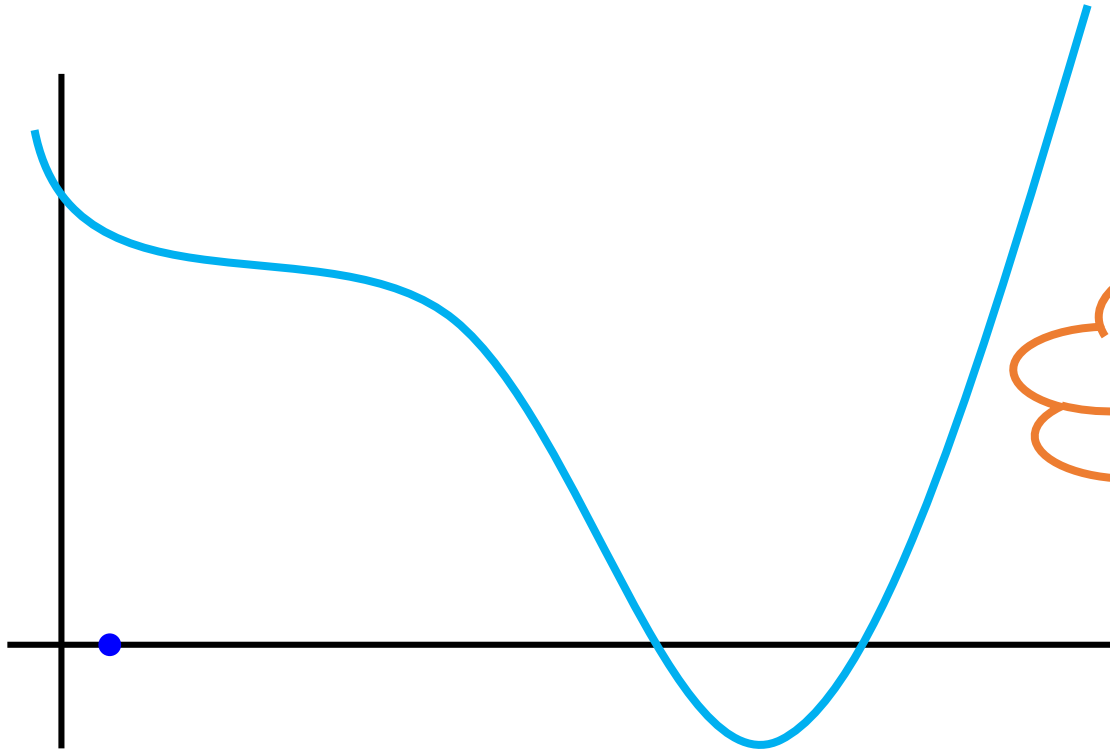


$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$

Better motivation: Finding
minima of quadratic fn. easy

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \langle \mathbf{b}, \mathbf{x} \rangle + c$$

Newton Method



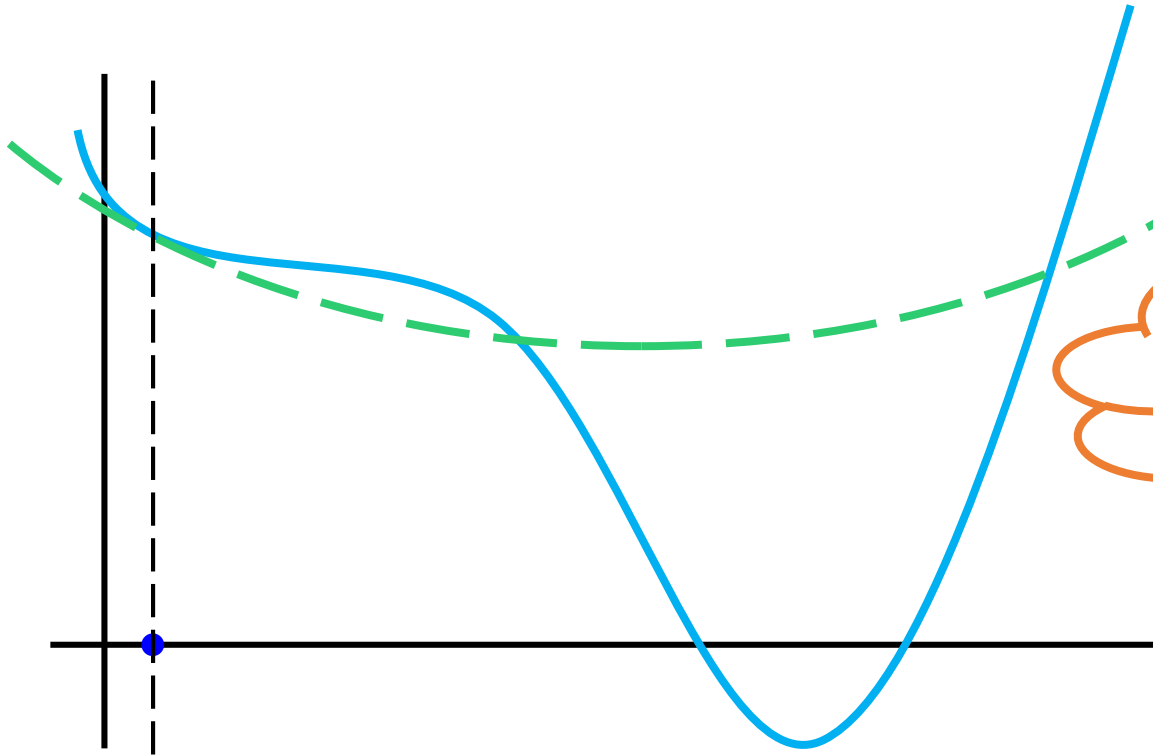
$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$

Better motivation: Finding
minima of quadratic fn. easy

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \langle \mathbf{b}, \mathbf{x} \rangle + c$$

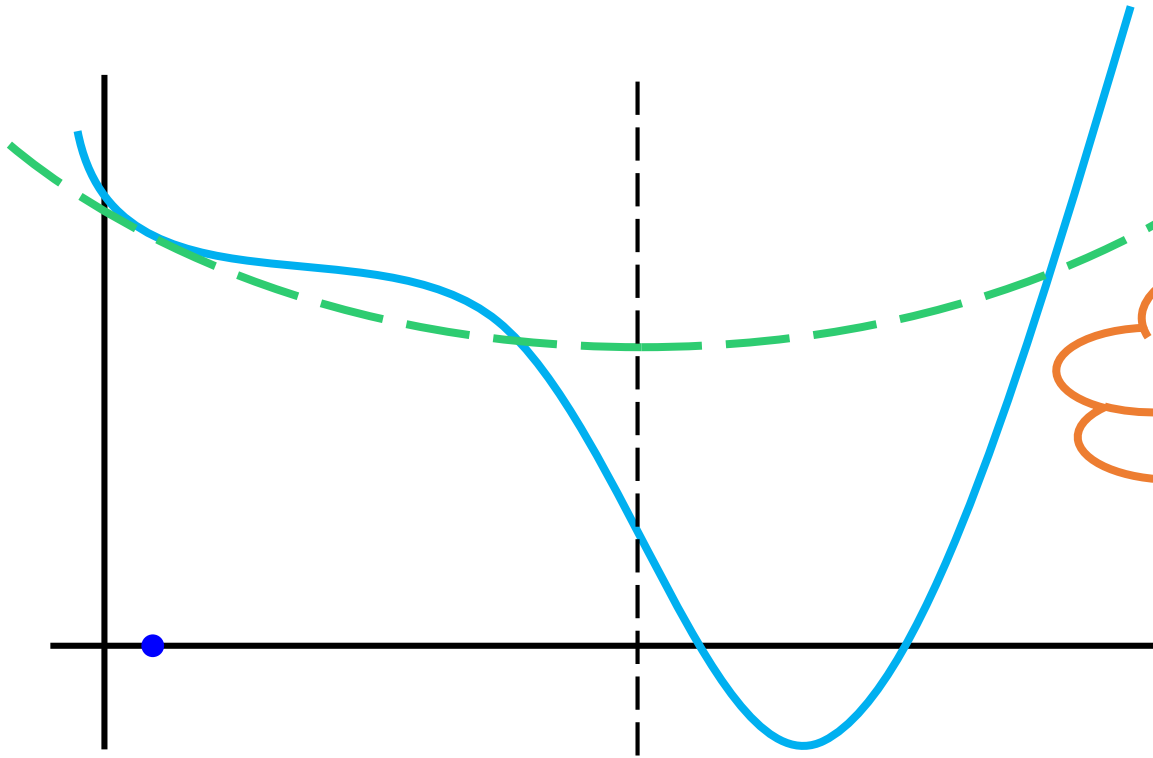
Newton Method

$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$



Better motivation: Finding
minima of quadratic fn. easy
 $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \langle \mathbf{b}, \mathbf{x} \rangle + c$

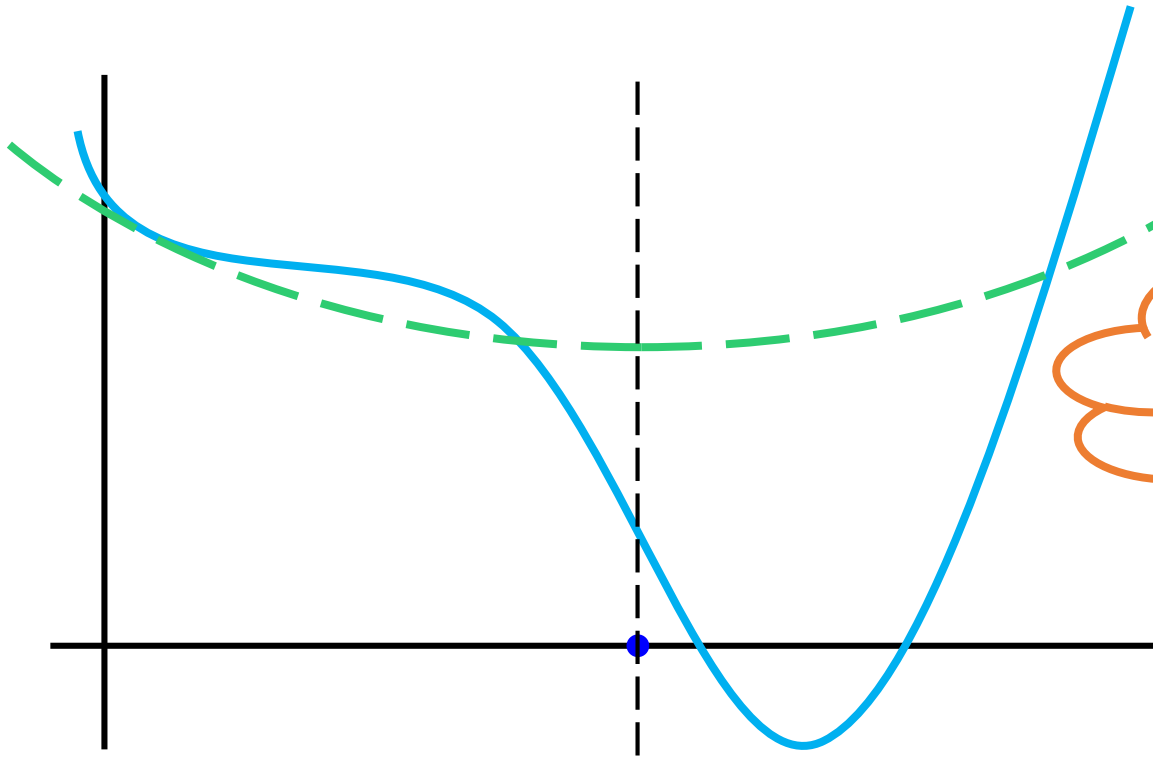
Newton Method



$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$

Better motivation: Finding
minima of quadratic fn. easy
 $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \langle \mathbf{b}, \mathbf{x} \rangle + c$

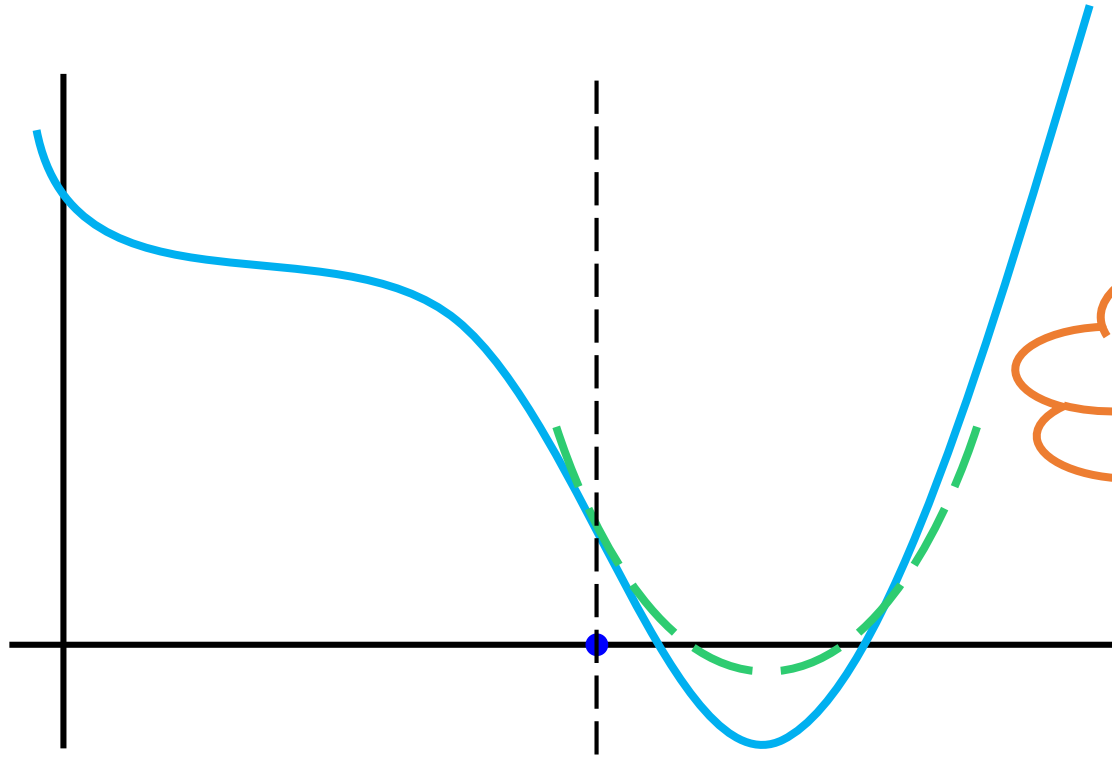
Newton Method



$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$

Better motivation: Finding
minima of quadratic fn. easy
 $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \langle \mathbf{b}, \mathbf{x} \rangle + c$

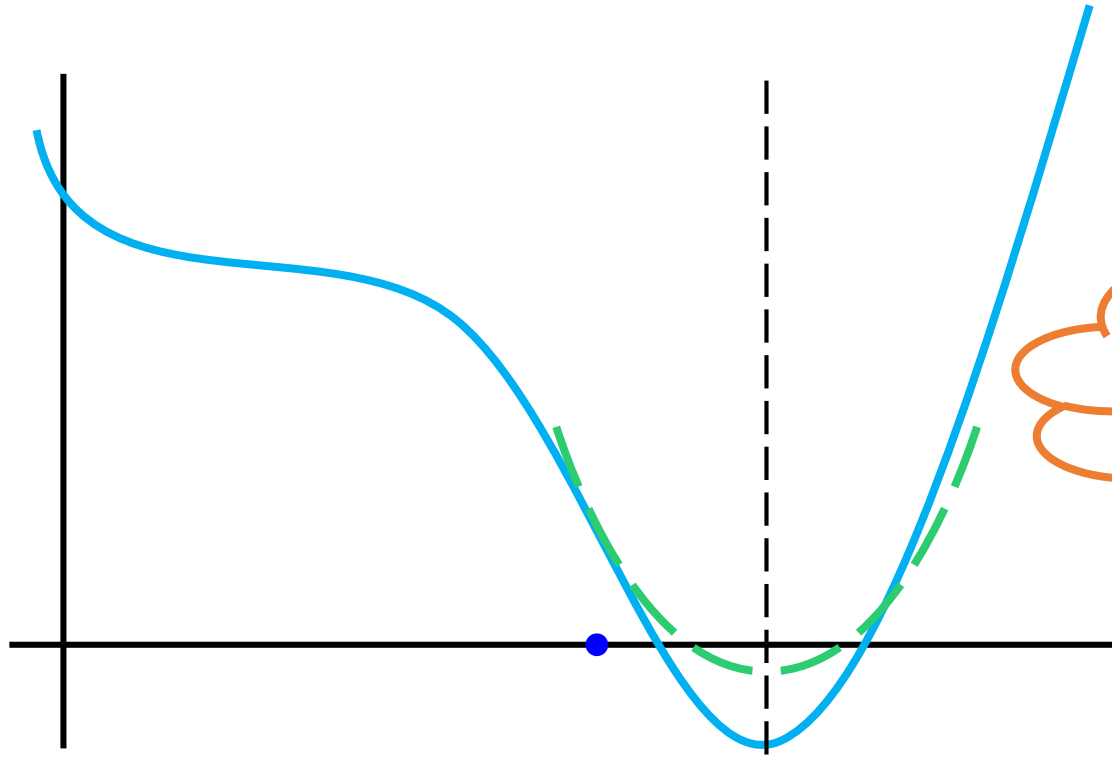
Newton Method



$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$

Better motivation: Finding
minima of quadratic fn. easy
 $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \langle \mathbf{b}, \mathbf{x} \rangle + c$

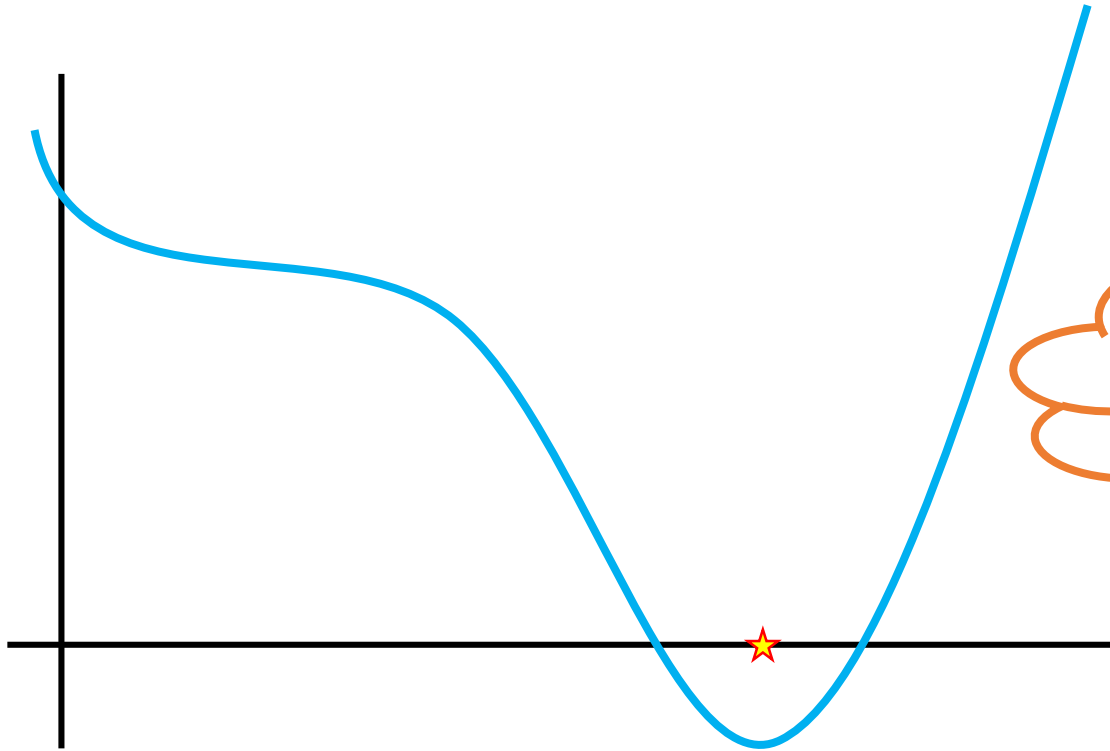
Newton Method



$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$

Better motivation: Finding
minima of quadratic fn. easy
 $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \langle \mathbf{b}, \mathbf{x} \rangle + c$

Newton Method

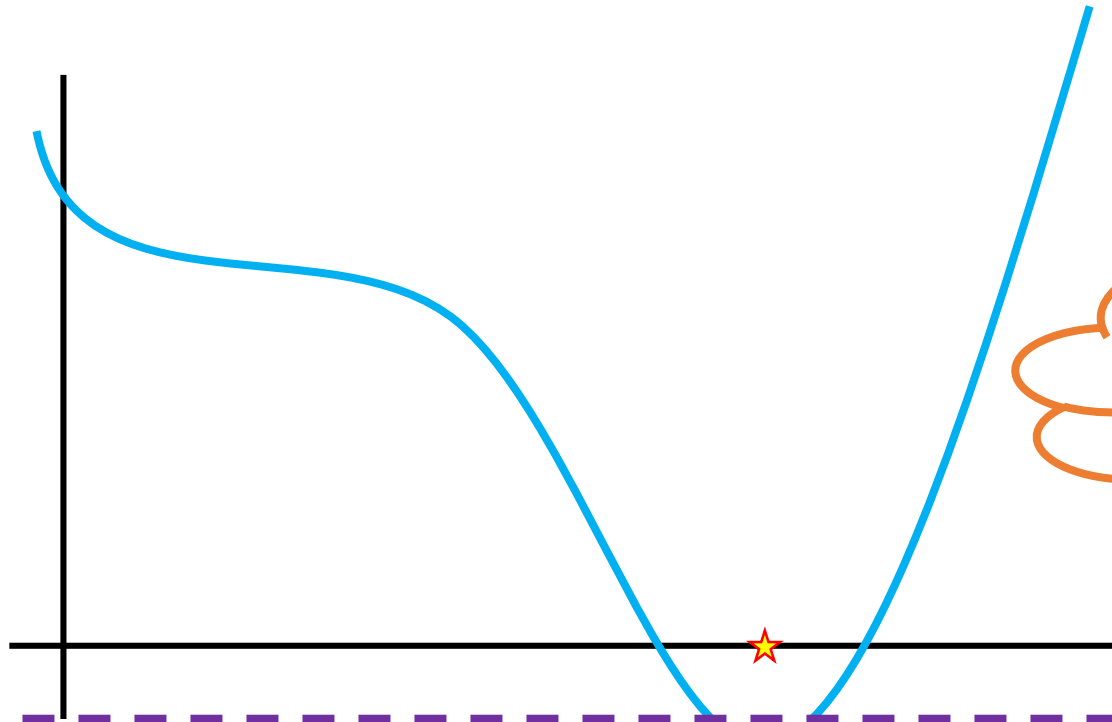


$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$

Better motivation: Finding
minima of quadratic fn. easy

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \langle \mathbf{b}, \mathbf{x} \rangle + c$$

Newton Method



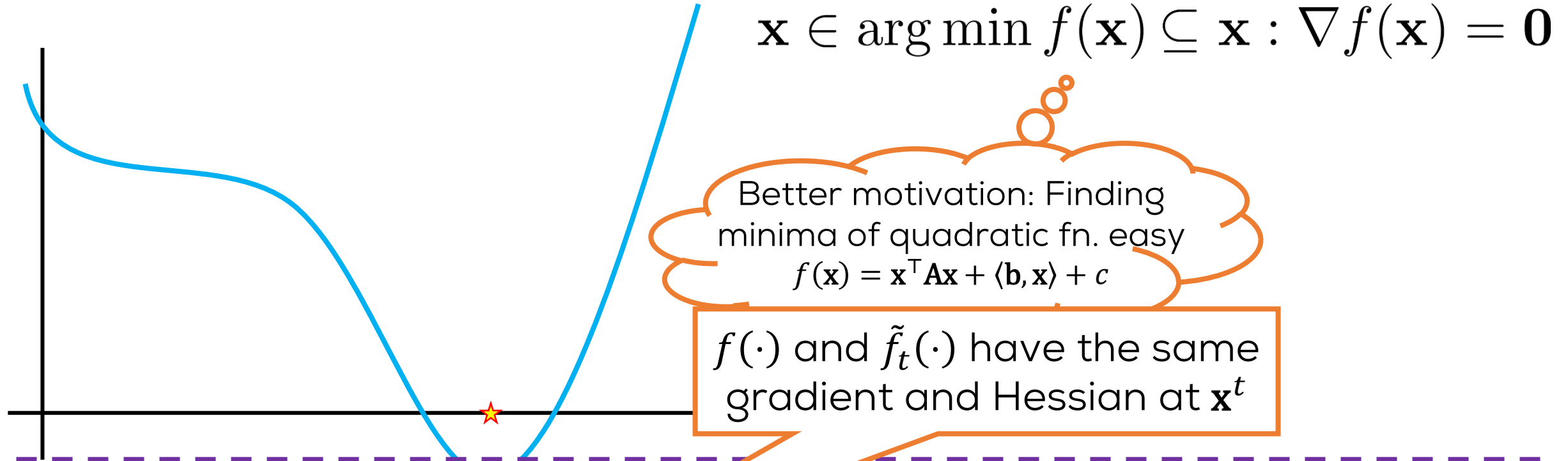
$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$

Better motivation: Finding
minima of quadratic fn. easy
 $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \langle \mathbf{b}, \mathbf{x} \rangle + c$

NEWTON METHOD

1. Initialize \mathbf{x}^0
2. Approx. $f(\cdot)$ by $\tilde{f}_t(\mathbf{x}) = f(\mathbf{x}^t) + \langle \nabla f(\mathbf{x}^t), \mathbf{x} - \mathbf{x}^t \rangle + \frac{1}{2} \langle \mathbf{x} - \mathbf{x}^t, \nabla^2 f(\mathbf{x}^t)(\mathbf{x} - \mathbf{x}^t) \rangle$
3. Update $\mathbf{x}^{t+1} \leftarrow \arg \min \tilde{f}_t(\mathbf{x}) = \mathbf{x}^t - (\nabla^2 f(\mathbf{x}^t))^{-1} \nabla f(\mathbf{x}^t)$
4. Repeat until convergence

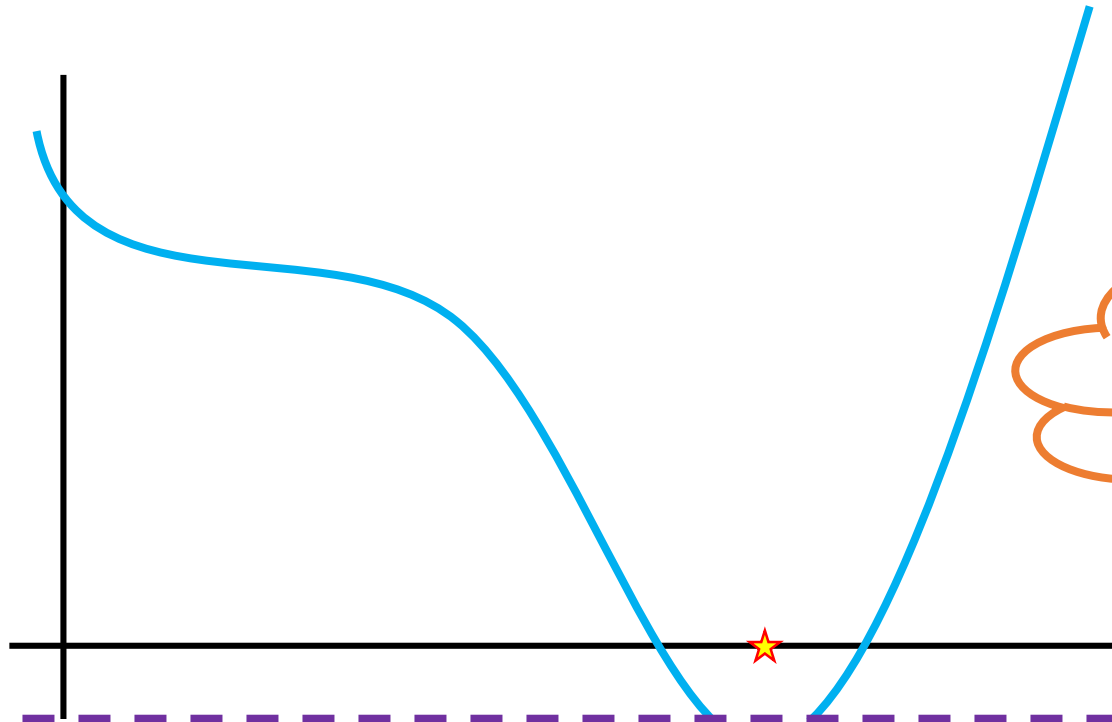
Newton Method



NEWTON METHOD

1. Initialize \mathbf{x}^0
2. Approx. $f(\cdot)$ by $\tilde{f}_t(\mathbf{x}) = f(\mathbf{x}^t) + \langle \nabla f(\mathbf{x}^t), \mathbf{x} - \mathbf{x}^t \rangle + \frac{1}{2} \langle \mathbf{x} - \mathbf{x}^t, \nabla^2 f(\mathbf{x}^t)(\mathbf{x} - \mathbf{x}^t) \rangle$
3. Update $\mathbf{x}^{t+1} \leftarrow \arg \min \tilde{f}_t(\mathbf{x}) = \mathbf{x}^t - (\nabla^2 f(\mathbf{x}^t))^{-1} \nabla f(\mathbf{x}^t)$
4. Repeat until convergence

Newton Method



$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$

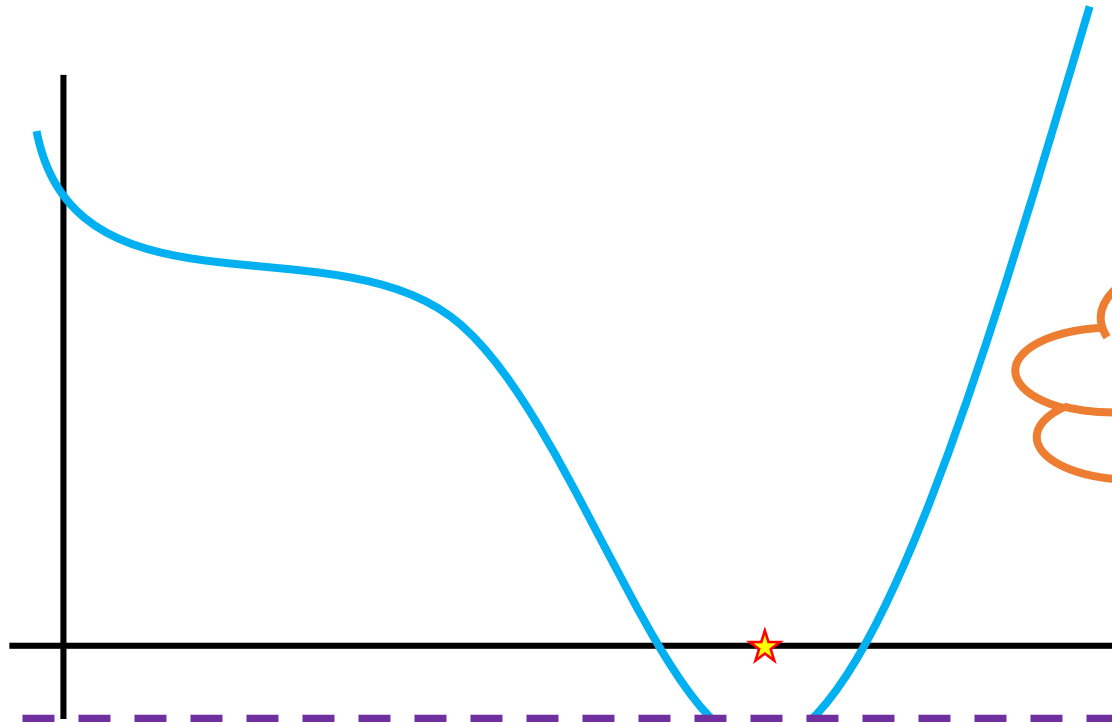
Better motivation: Finding minima of quadratic fn. easy
 $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \langle \mathbf{b}, \mathbf{x} \rangle + c$

Ridiculously fast for convex problems

NEWTON METHOD

1. Initialize \mathbf{x}^0
2. Approx. $f(\cdot)$ by $\tilde{f}_t(\mathbf{x}) = f(\mathbf{x}^t) + \langle \nabla f(\mathbf{x}^t), \mathbf{x} - \mathbf{x}^t \rangle + \frac{1}{2} \langle \mathbf{x} - \mathbf{x}^t, \nabla^2 f(\mathbf{x}^t)(\mathbf{x} - \mathbf{x}^t) \rangle$
3. Update $\mathbf{x}^{t+1} \leftarrow \arg \min \tilde{f}_t(\mathbf{x}) = \mathbf{x}^t - (\nabla^2 f(\mathbf{x}^t))^{-1} \nabla f(\mathbf{x}^t)$
4. Repeat until convergence

Newton Method



$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$

Better motivation: Finding minima of quadratic fn. easy
 $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \langle \mathbf{b}, \mathbf{x} \rangle + c$

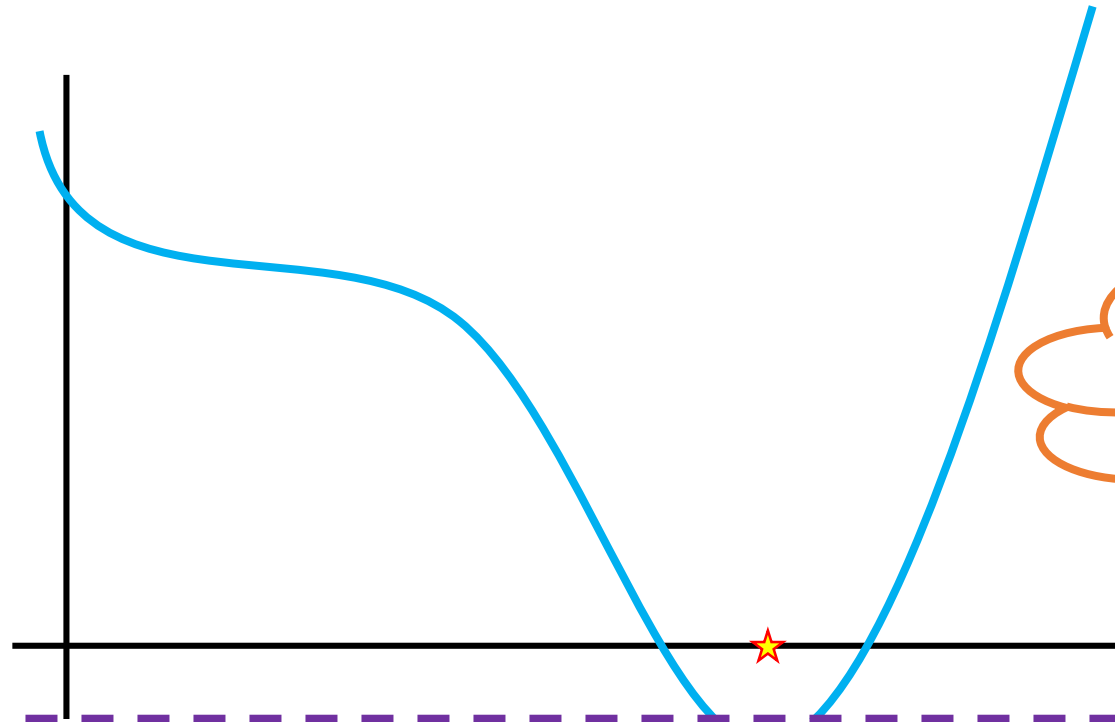
$O\left(\log \log \frac{1}{\epsilon}\right)$ steps

Ridiculously fast for convex problems

NEWTON METHOD

1. Initialize \mathbf{x}^0
2. Approx. $f(\cdot)$ by $\tilde{f}_t(\mathbf{x}) = f(\mathbf{x}^t) + \langle \nabla f(\mathbf{x}^t), \mathbf{x} - \mathbf{x}^t \rangle + \frac{1}{2} \langle \mathbf{x} - \mathbf{x}^t, \nabla^2 f(\mathbf{x}^t)(\mathbf{x} - \mathbf{x}^t) \rangle$
3. Update $\mathbf{x}^{t+1} \leftarrow \arg \min \tilde{f}_t(\mathbf{x}) = \mathbf{x}^t - (\nabla^2 f(\mathbf{x}^t))^{-1} \nabla f(\mathbf{x}^t)$
4. Repeat until convergence

Newton Method



$$\mathbf{x} \in \arg \min f(\mathbf{x}) \subseteq \mathbf{x} : \nabla f(\mathbf{x}) = \mathbf{0}$$

Better motivation: Finding minima of quadratic fn. easy
 $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \langle \mathbf{b}, \mathbf{x} \rangle + c$

Steps more costly $O(d^2)$

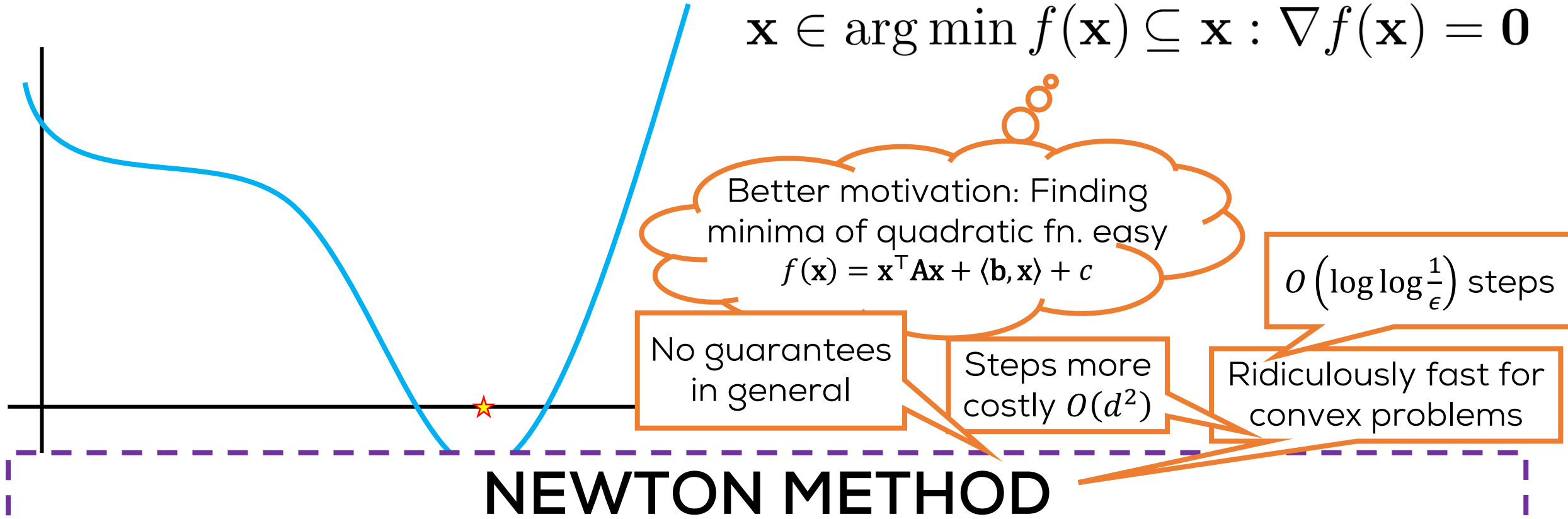
Ridiculously fast for convex problems

$O\left(\log \log \frac{1}{\epsilon}\right)$ steps

NEWTON METHOD

1. Initialize \mathbf{x}^0
2. Approx. $f(\cdot)$ by $\tilde{f}_t(\mathbf{x}) = f(\mathbf{x}^t) + \langle \nabla f(\mathbf{x}^t), \mathbf{x} - \mathbf{x}^t \rangle + \frac{1}{2} \langle \mathbf{x} - \mathbf{x}^t, \nabla^2 f(\mathbf{x}^t)(\mathbf{x} - \mathbf{x}^t) \rangle$
3. Update $\mathbf{x}^{t+1} \leftarrow \arg \min \tilde{f}_t(\mathbf{x}) = \mathbf{x}^t - (\nabla^2 f(\mathbf{x}^t))^{-1} \nabla f(\mathbf{x}^t)$
4. Repeat until convergence

Newton Method



1. Initialize \mathbf{x}^0
2. Approx. $f(\cdot)$ by $\tilde{f}_t(\mathbf{x}) = f(\mathbf{x}^t) + \langle \nabla f(\mathbf{x}^t), \mathbf{x} - \mathbf{x}^t \rangle + \frac{1}{2} \langle \mathbf{x} - \mathbf{x}^t, \nabla^2 f(\mathbf{x}^t)(\mathbf{x} - \mathbf{x}^t) \rangle$
3. Update $\mathbf{x}^{t+1} \leftarrow \arg \min \tilde{f}_t(\mathbf{x}) = \mathbf{x}^t - (\nabla^2 f(\mathbf{x}^t))^{-1} \nabla f(\mathbf{x}^t)$
4. Repeat until convergence

Duality

How to take the problem you want to solve
... and convert it to a problem you can solve

Fenchel, Lagrange, Wolfe, Pontryagin

Duality

How to take the problem you want to solve
... and convert it to a problem you can solve

Fenchel, Lagrange, Wolfe, Pontryagin

Duality

How to take the problem you want to solve
... and convert it to a problem you can solve

The Lagrangian

September 1, 2017



The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2$$
$$\text{s.t. } \|\mathbf{w}\|_2 \leq r$$

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

Want $g(\mathbf{w}) \geq 0$?
 $-g(\mathbf{w}) \leq 0$

The Lagrangian

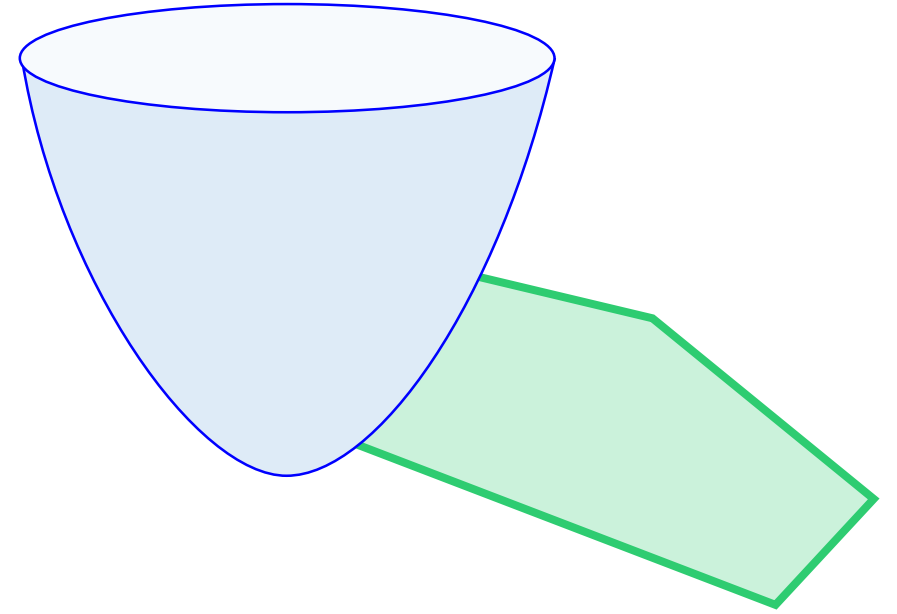
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

The Lagrangian

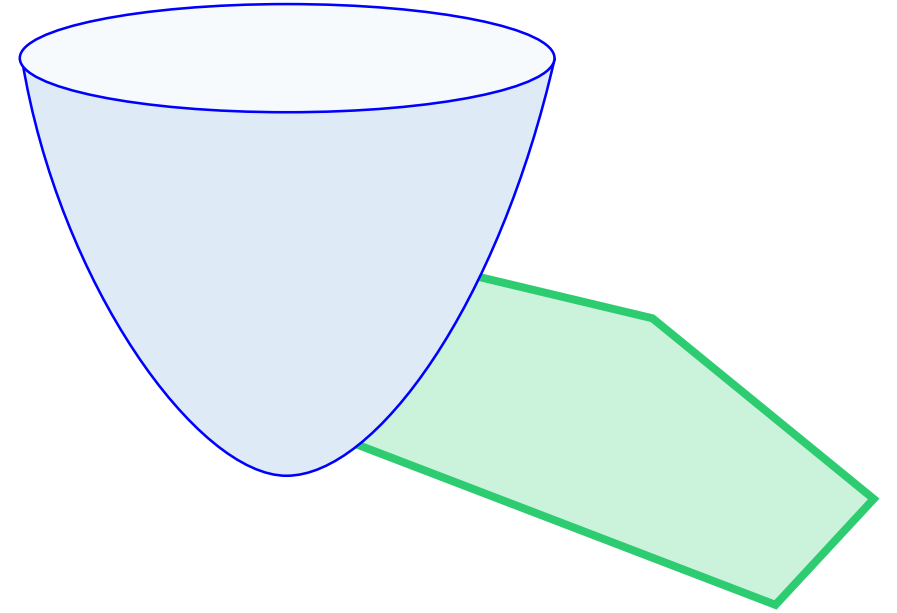
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



The Lagrangian

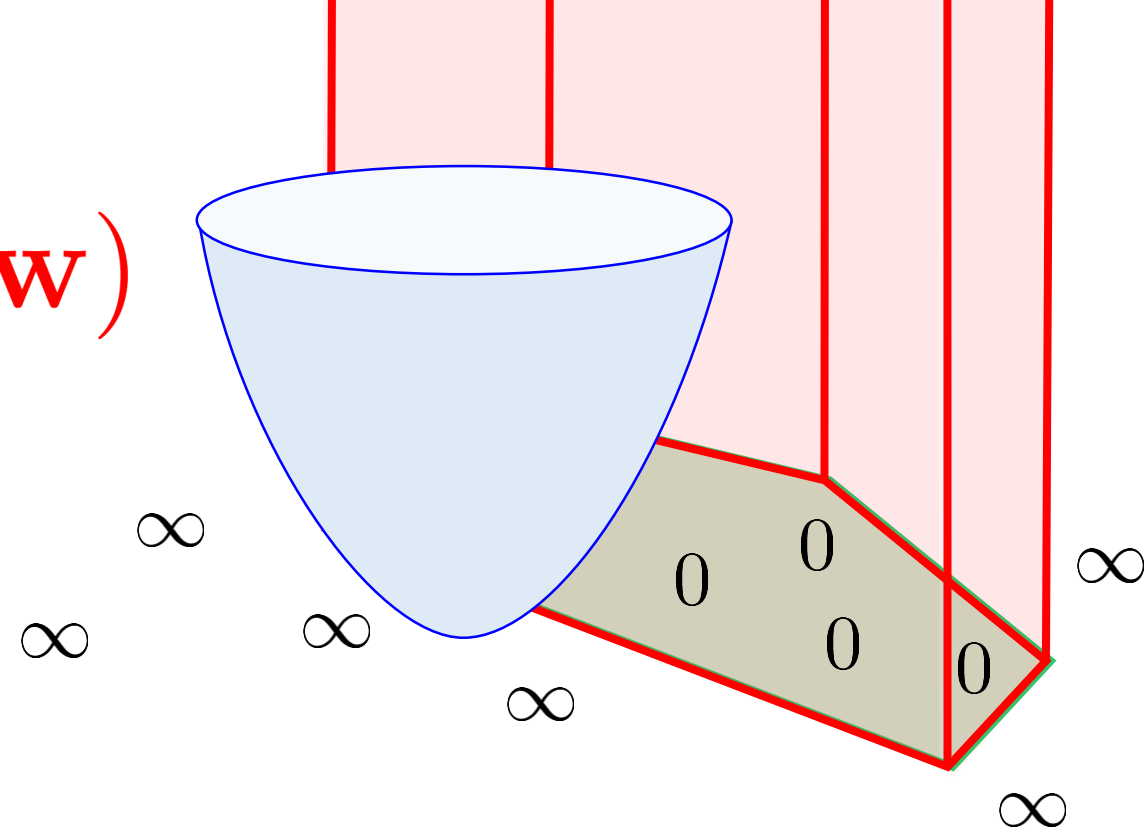
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$



The Lagrangian

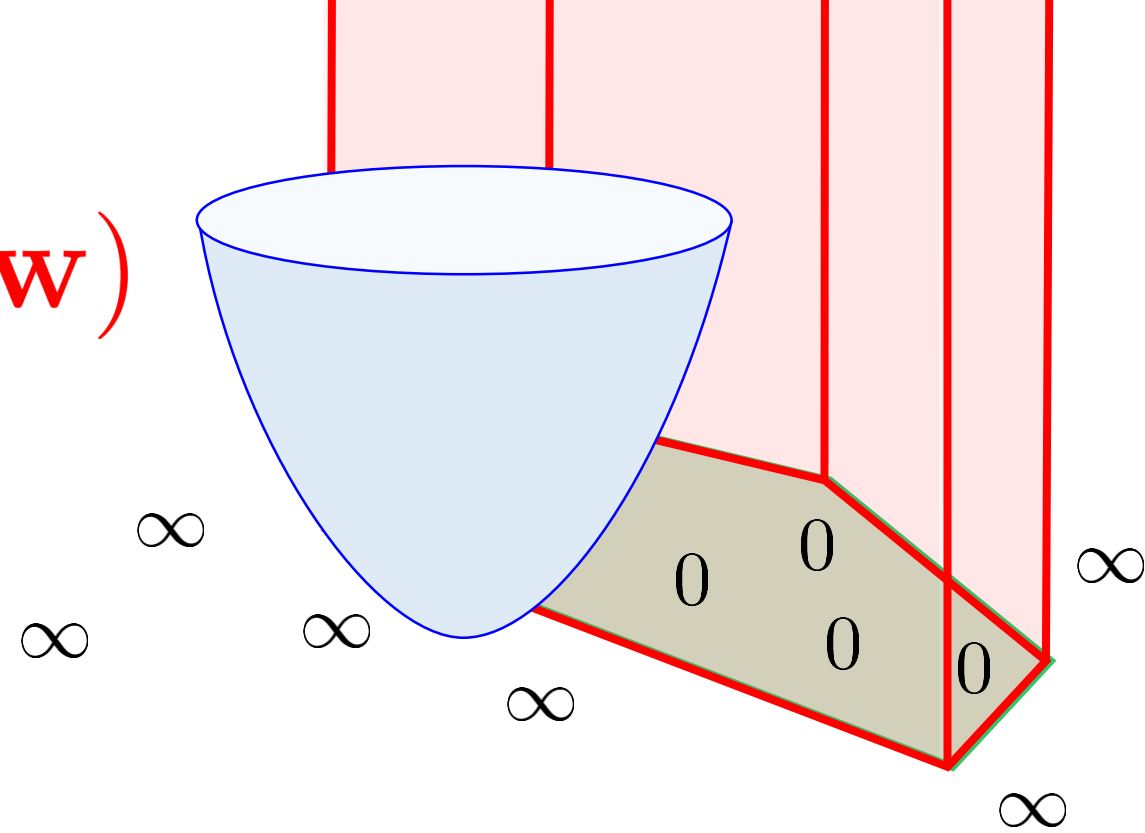
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$

Barrier Function



The Lagrangian

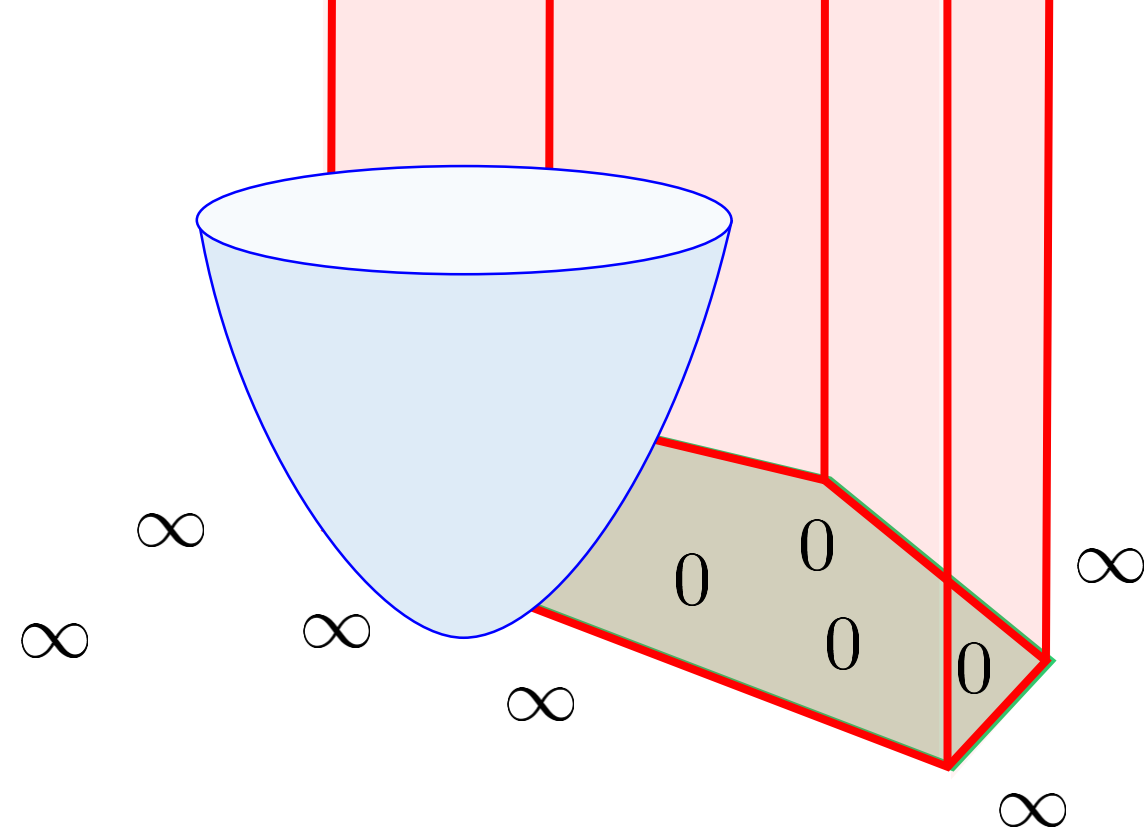
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$



The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

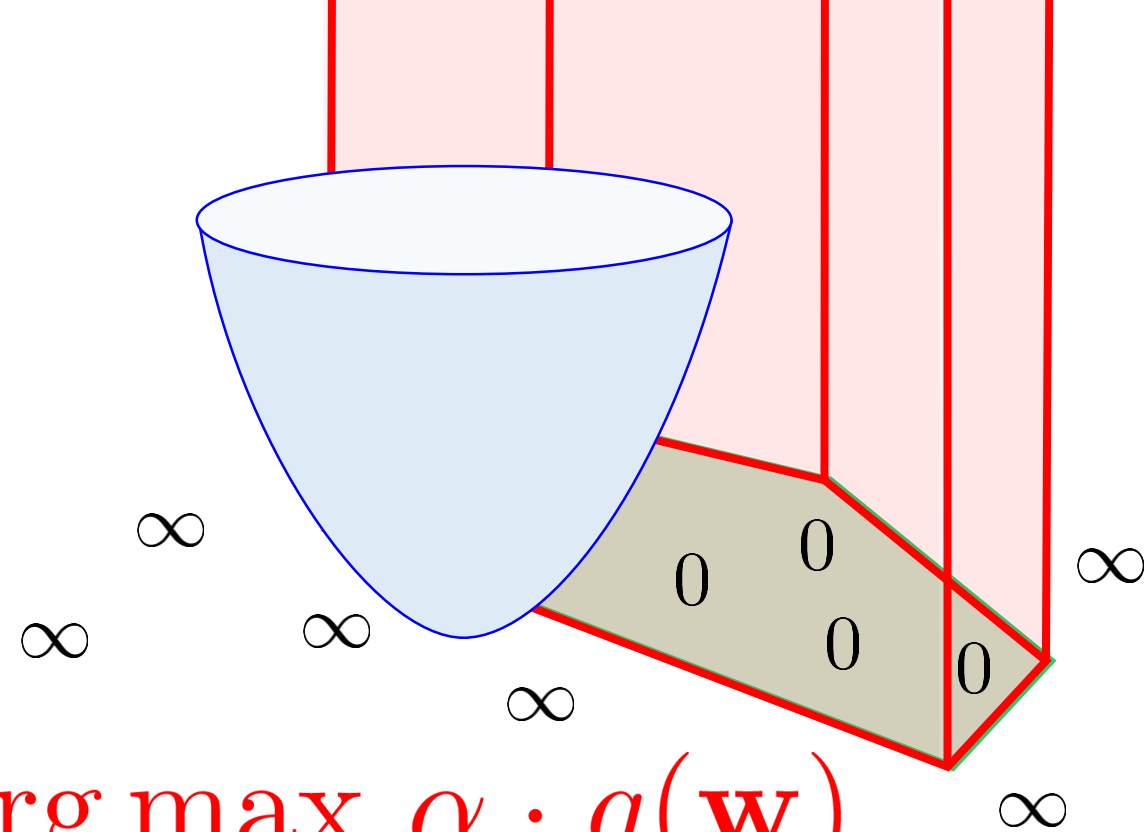


The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\arg \max_{\alpha \geq 0} \alpha \cdot g(\mathbf{w})$$



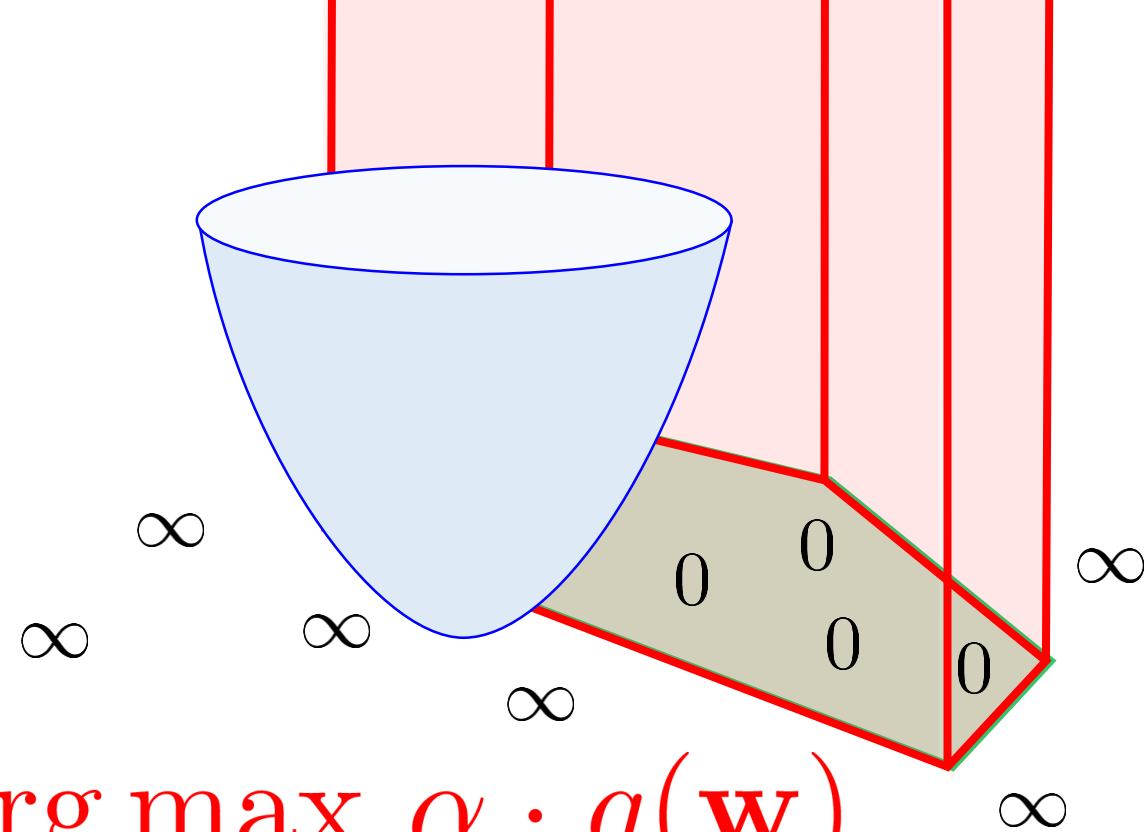
The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\arg \max_{\alpha \geq 0} \alpha \cdot g(\mathbf{w})$$

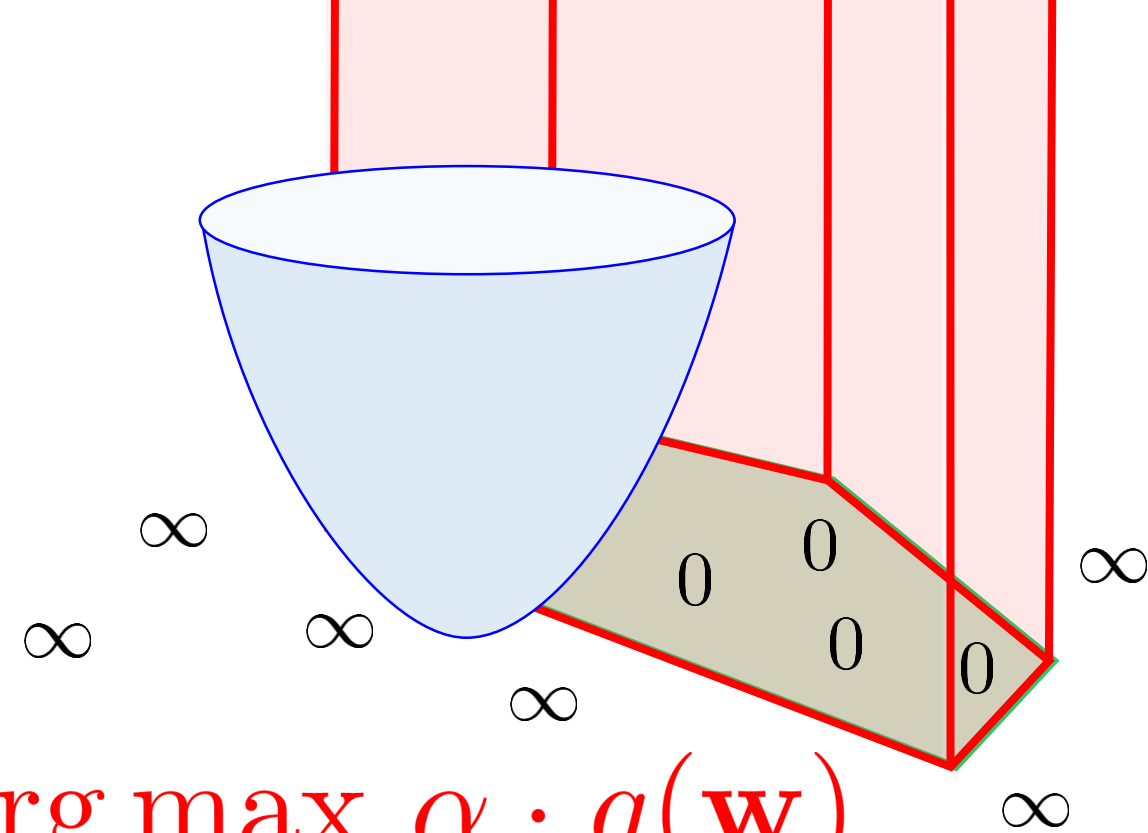
∞ if $g(\mathbf{w}) > 0$



The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\arg \max_{\alpha \geq 0} \alpha \cdot g(\mathbf{w})$$

∞ if $g(\mathbf{w}) > 0$

0 if $g(\mathbf{w}) \leq 0$

The Lagrangian

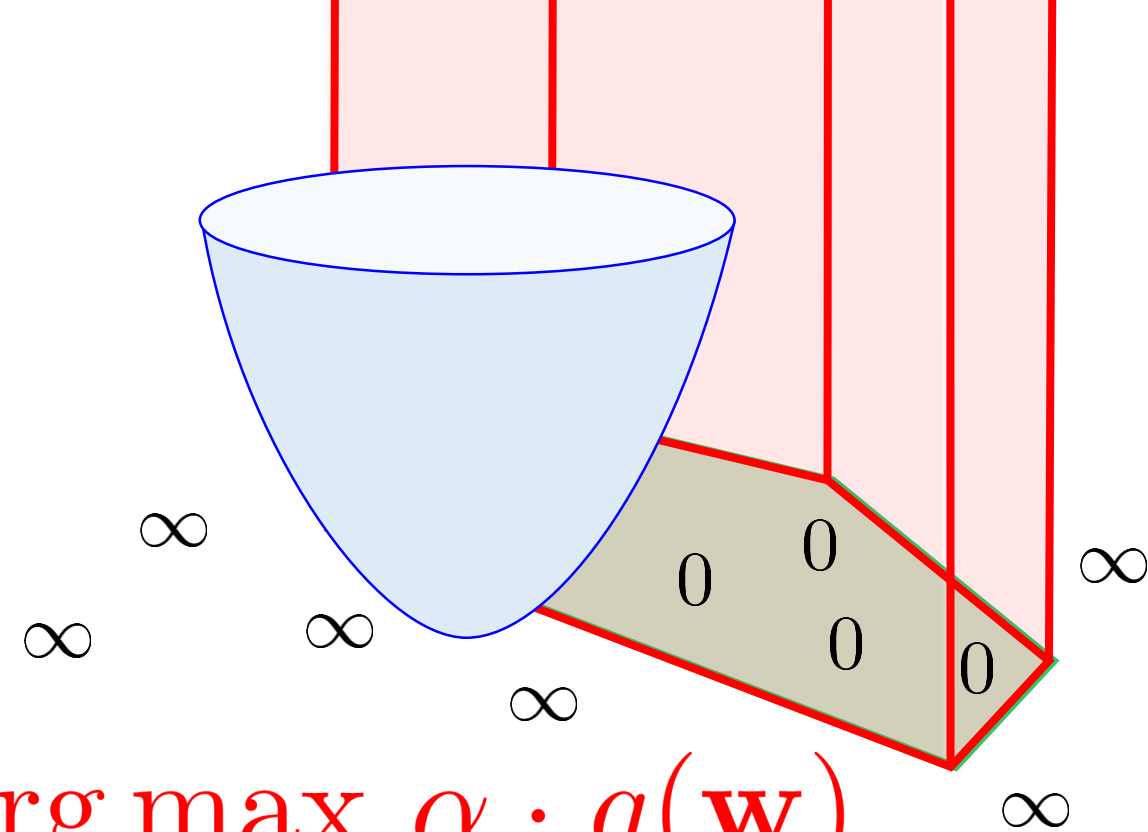
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$f(\mathbf{w}) + \arg \max_{\alpha \geq 0} \alpha \cdot g(\mathbf{w})$$

∞ if $g(\mathbf{w}) > 0$

0 if $g(\mathbf{w}) \leq 0$



The Lagrangian

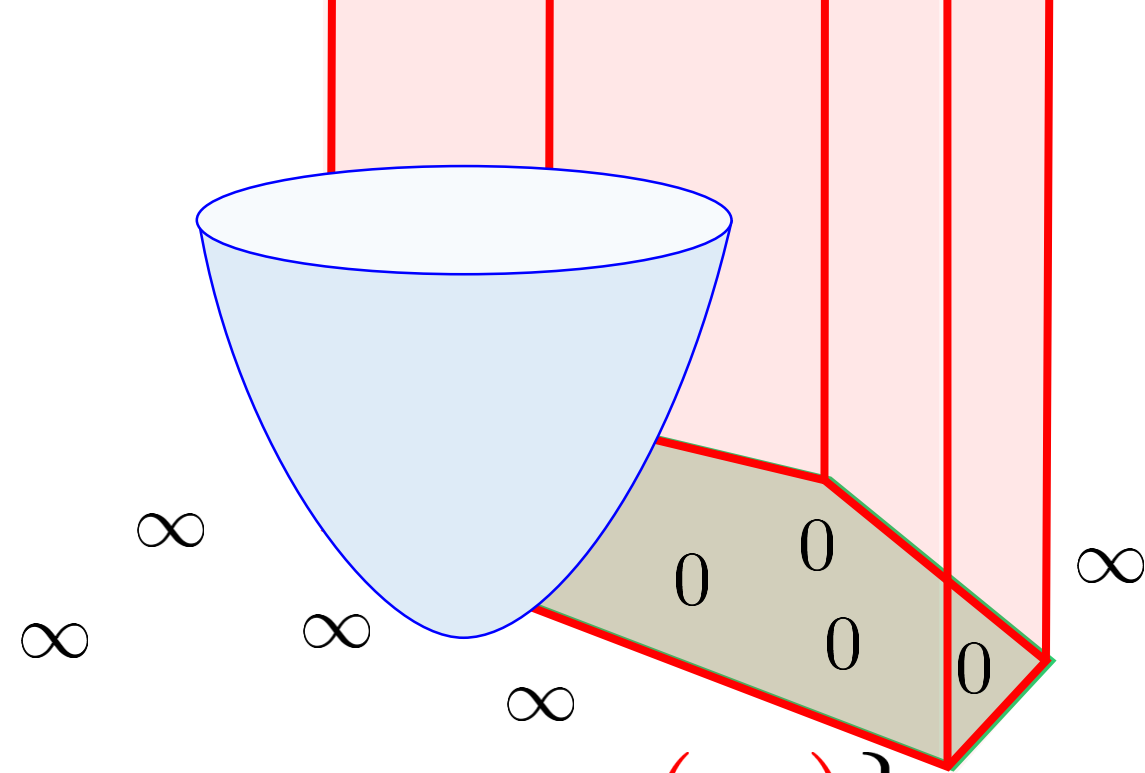
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}) + \arg \max_{\alpha \geq 0} \alpha \cdot g(\mathbf{w}) \right\}$$

∞ if $g(\mathbf{w}) > 0$

0 if $g(\mathbf{w}) \leq 0$



The Lagrangian

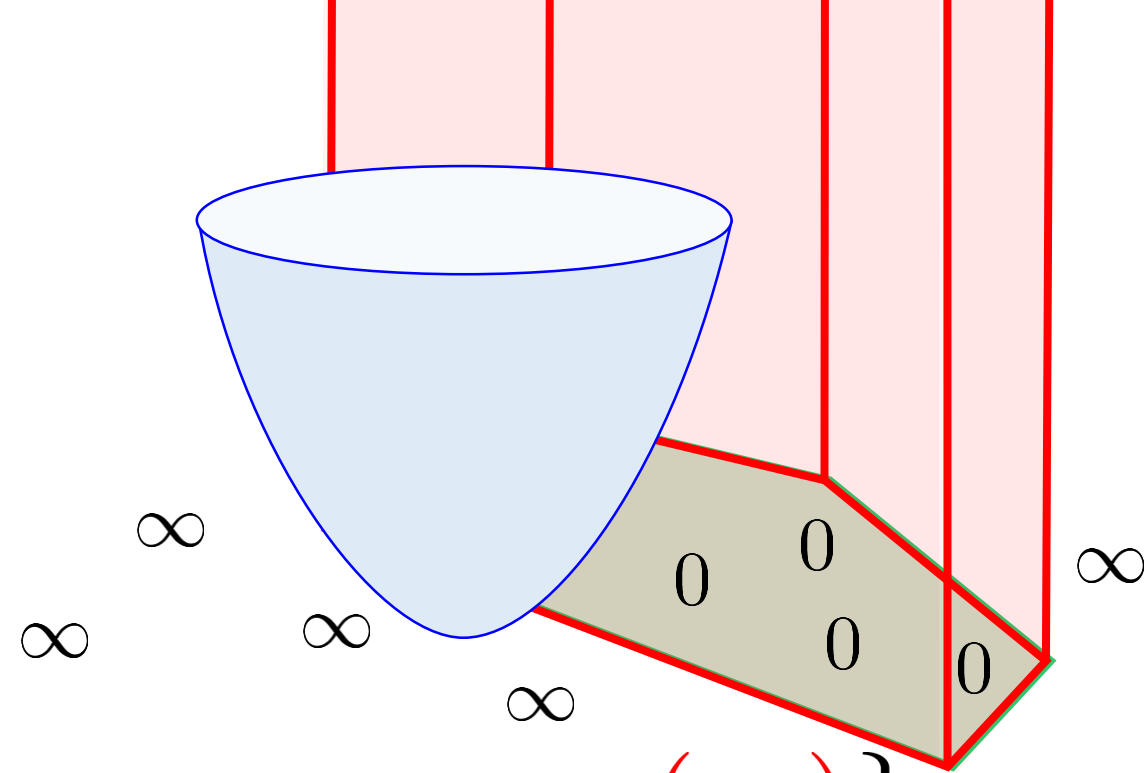
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}) + \arg \max_{\alpha \geq 0} \alpha \cdot g(\mathbf{w}) \right\}$$

∞ if $g(\mathbf{w}) > 0$

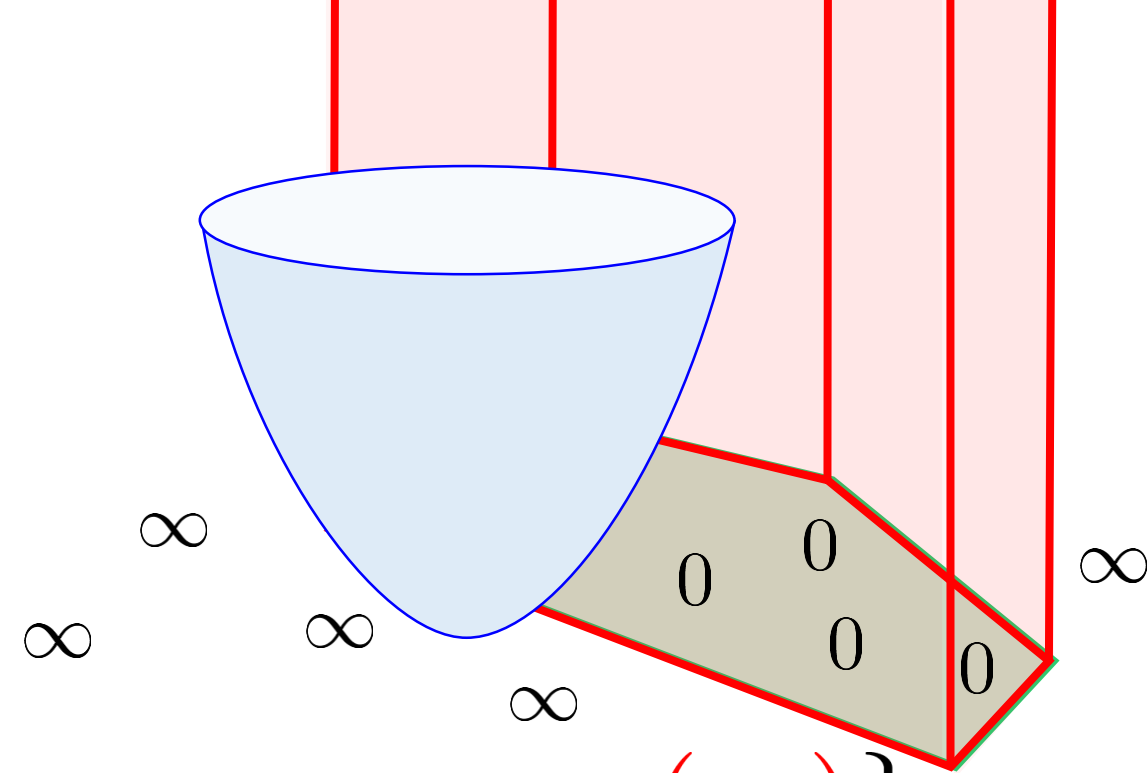
0 if $g(\mathbf{w}) \leq 0$



The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}) + \arg \max_{\alpha \geq 0} \alpha \cdot g(\mathbf{w}) \right\}$$

∞ if $g(\mathbf{w}) > 0$

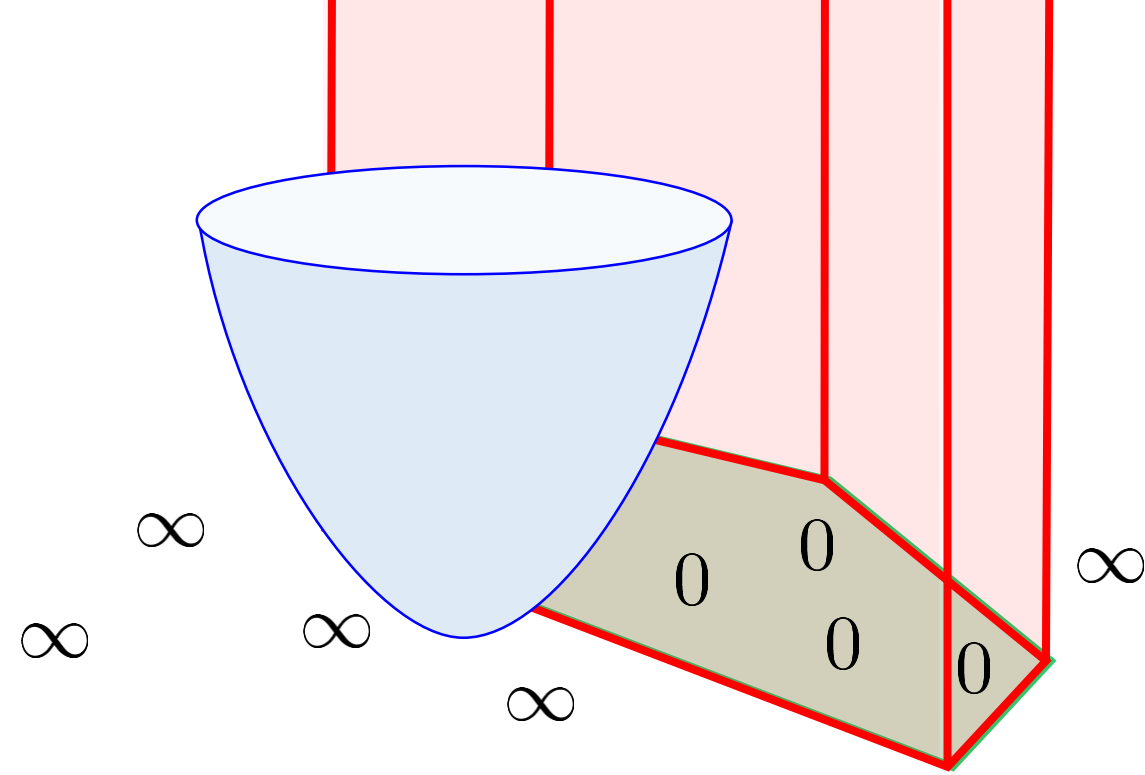
0 if $g(\mathbf{w}) \leq 0$

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \left\{ f(\mathbf{w}) + \alpha \cdot g(\mathbf{w}) \right\} \right\}$$

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

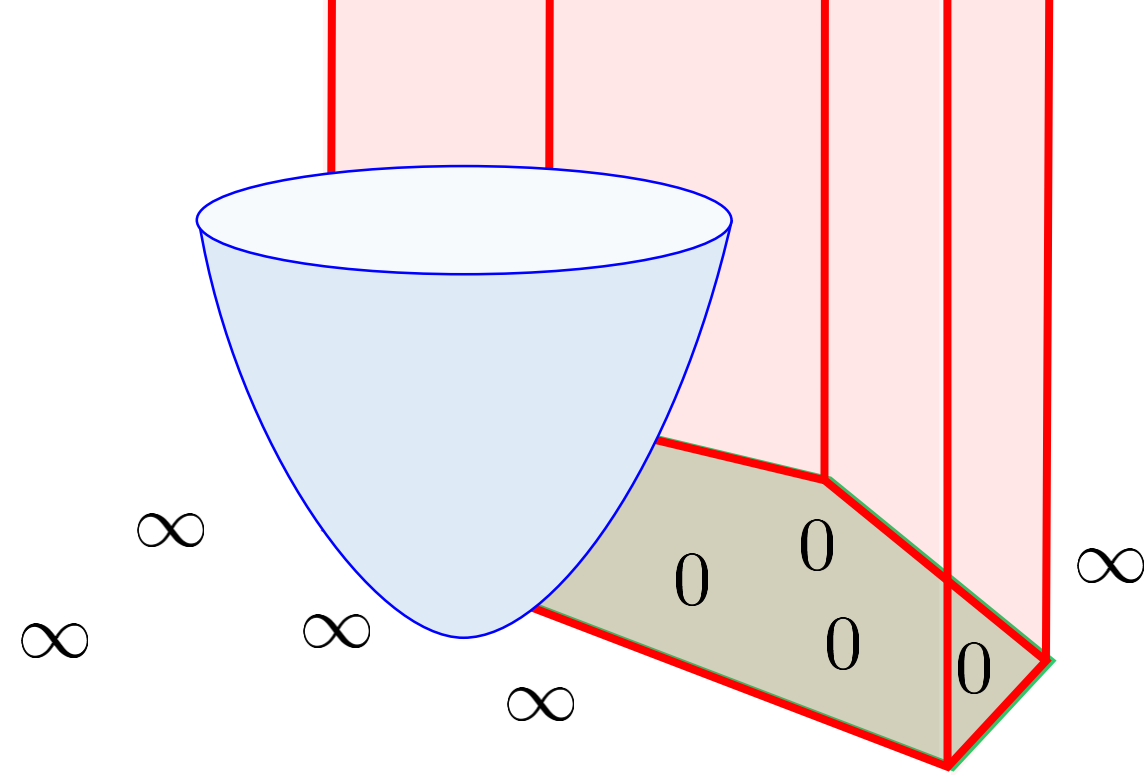


$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \{f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})\} \right\}$$

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



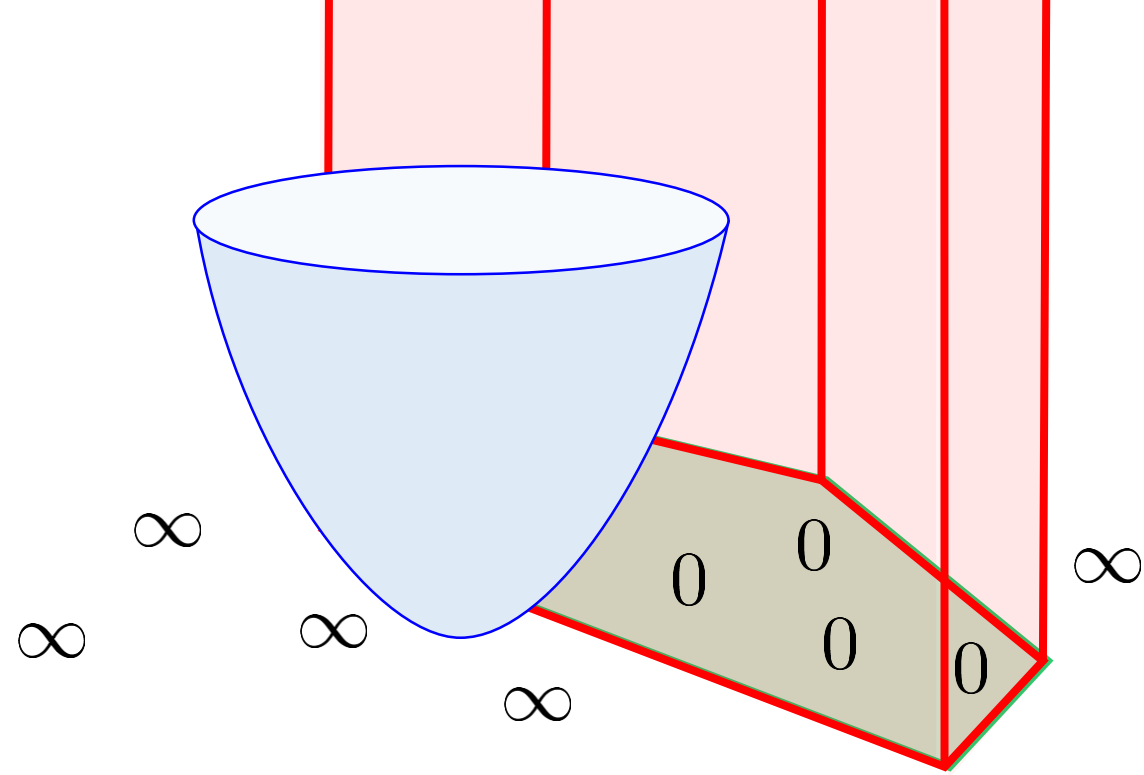
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \{f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})\} \right\}$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \{f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})\} \right\}^\infty$$

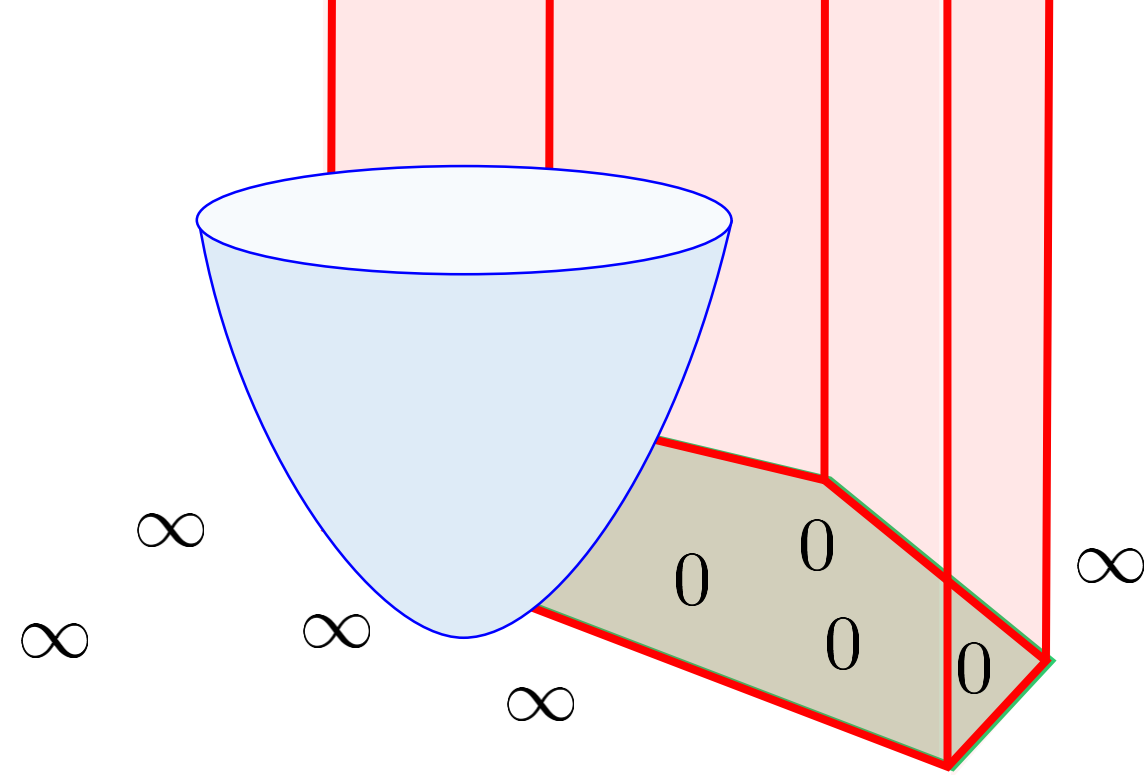
$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

Lagrangian

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \{f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})\} \right\}$$

Lagrange multiplier

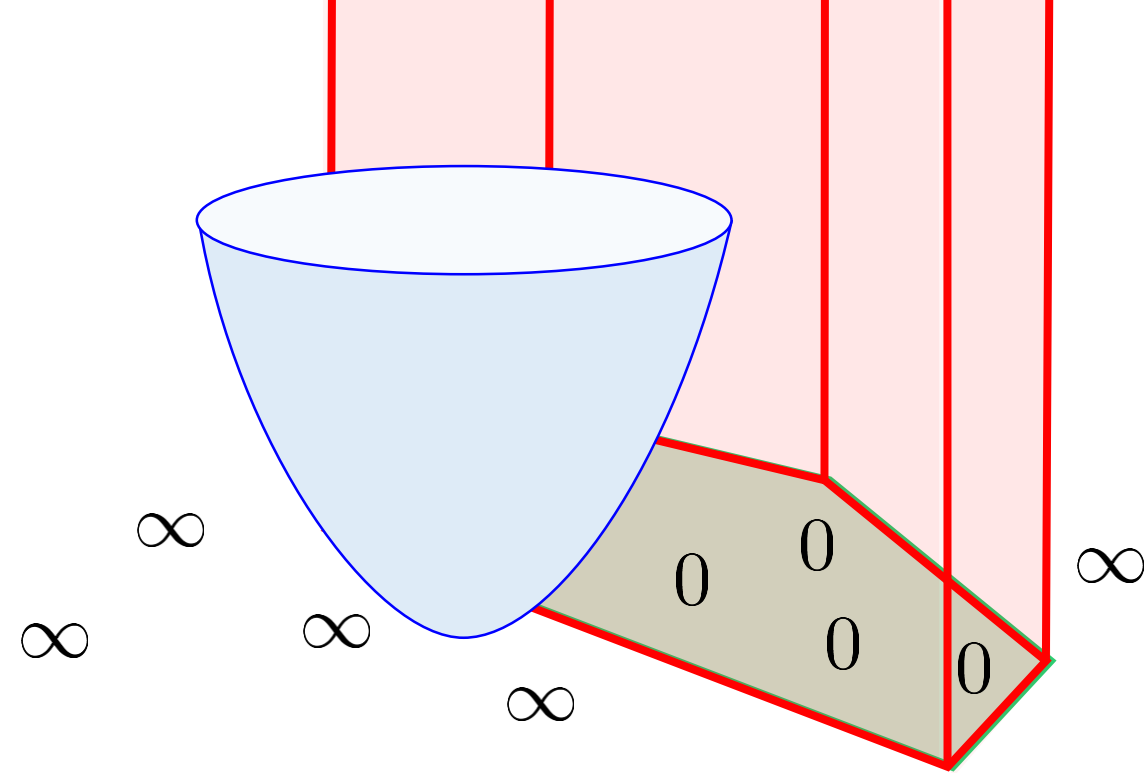
$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

Lagrangian

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \left\{ \mathcal{L}(\mathbf{w}, \alpha) \right\} \right\}$$

Lagrange multiplier

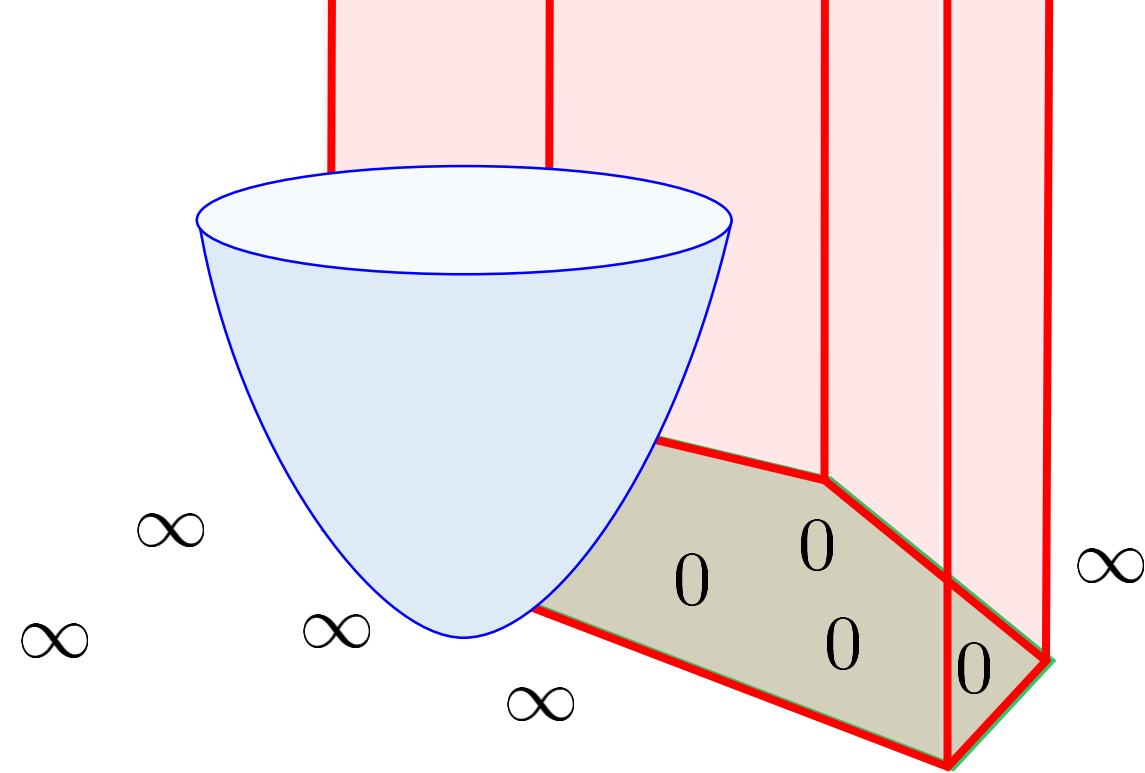
$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

Lagrangian

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \left\{ \mathcal{L}(\mathbf{w}, \alpha) \right\} \right\}$$

Lagrange multiplier

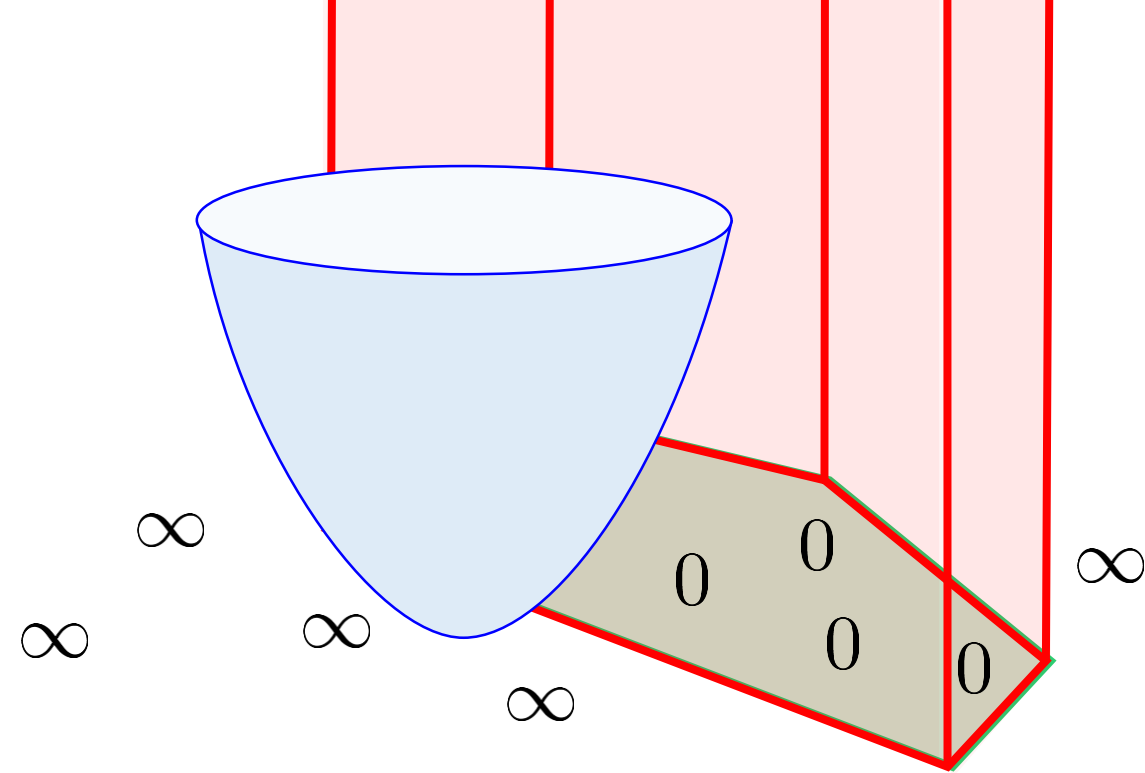
$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

Lagrangian

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \left\{ \mathcal{L}(\mathbf{w}, \alpha) \right\} \right\}$$

Lagrange multiplier

Primal problem

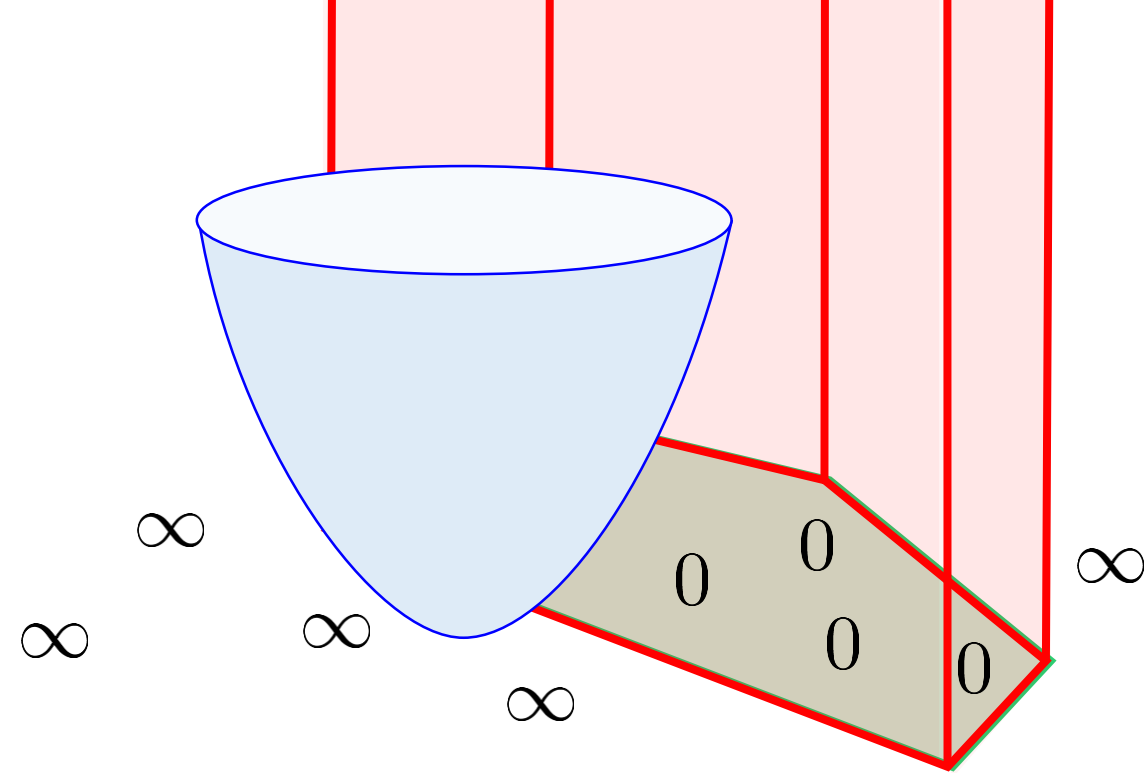
$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

Lagrangian

The Lagrangian

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \left\{ \mathcal{L}(\mathbf{w}, \alpha) \right\} \right\}$$

Lagrange multiplier

Primal problem

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

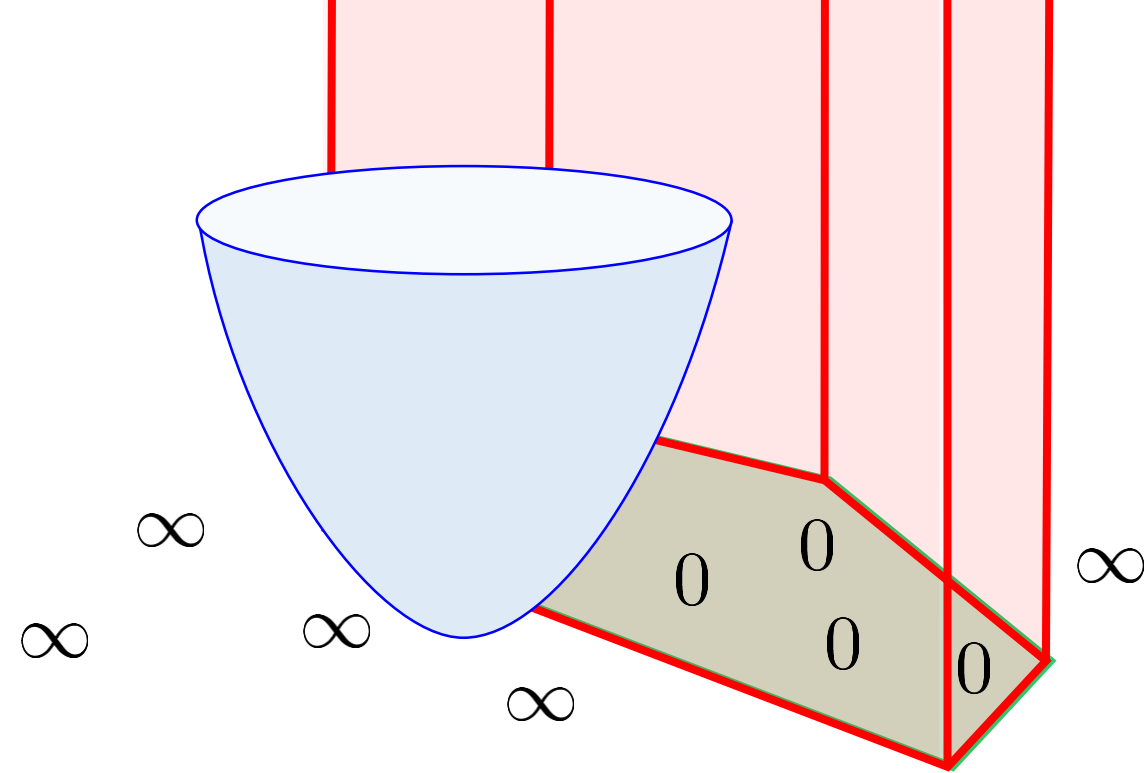
Lagrangian

The Lagrangian

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

Primal problem

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \left\{ \mathcal{L}(\mathbf{w}, \alpha) \right\} \right\}$$

Lagrange multiplier

Primal problem

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

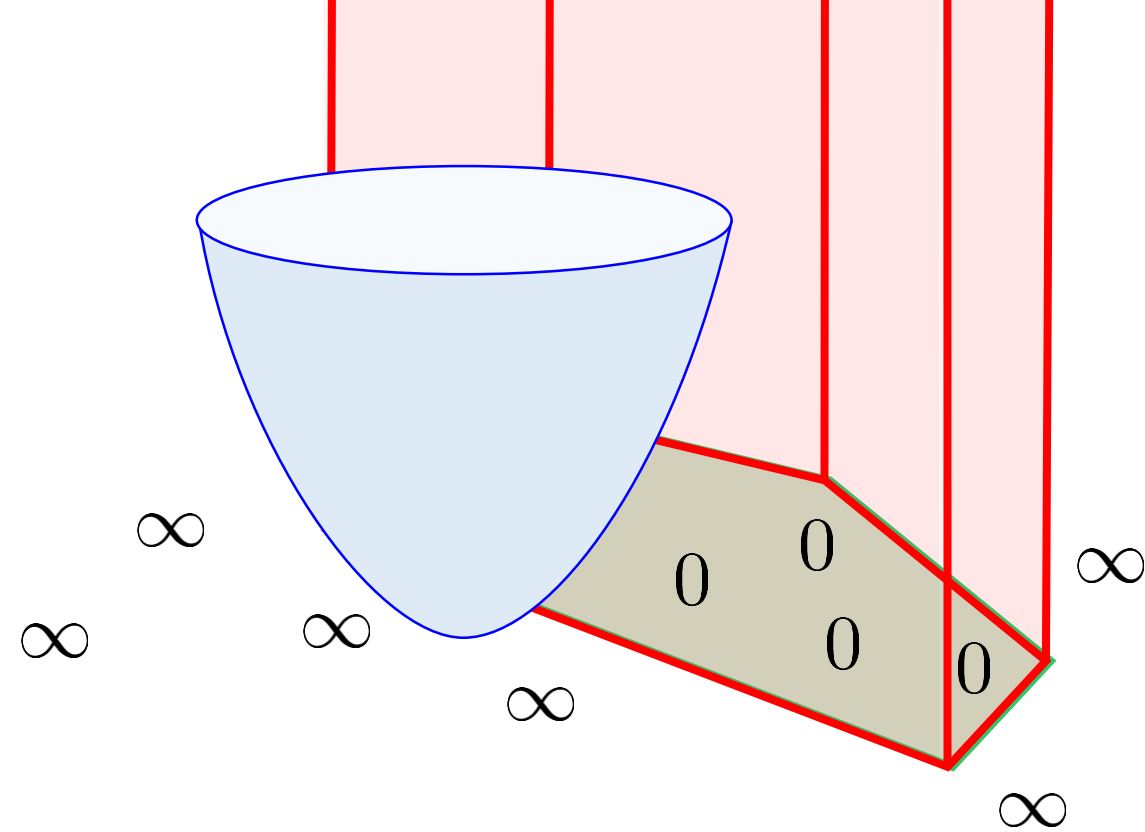
Lagrangian

The Lagrangian

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

Primal
problem

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

Lagrange multiplier

Primal problem

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

Lagrangian

The Lagrangian

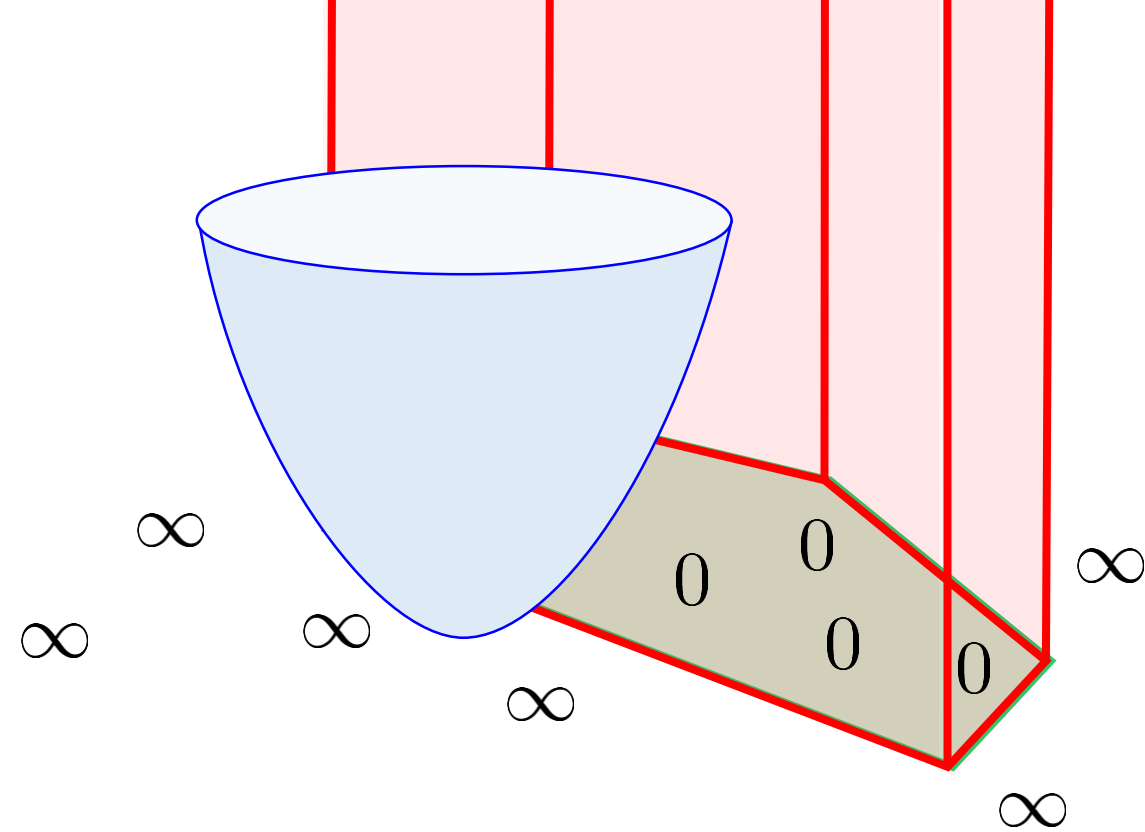
$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

Primal
problem

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

Primal problem

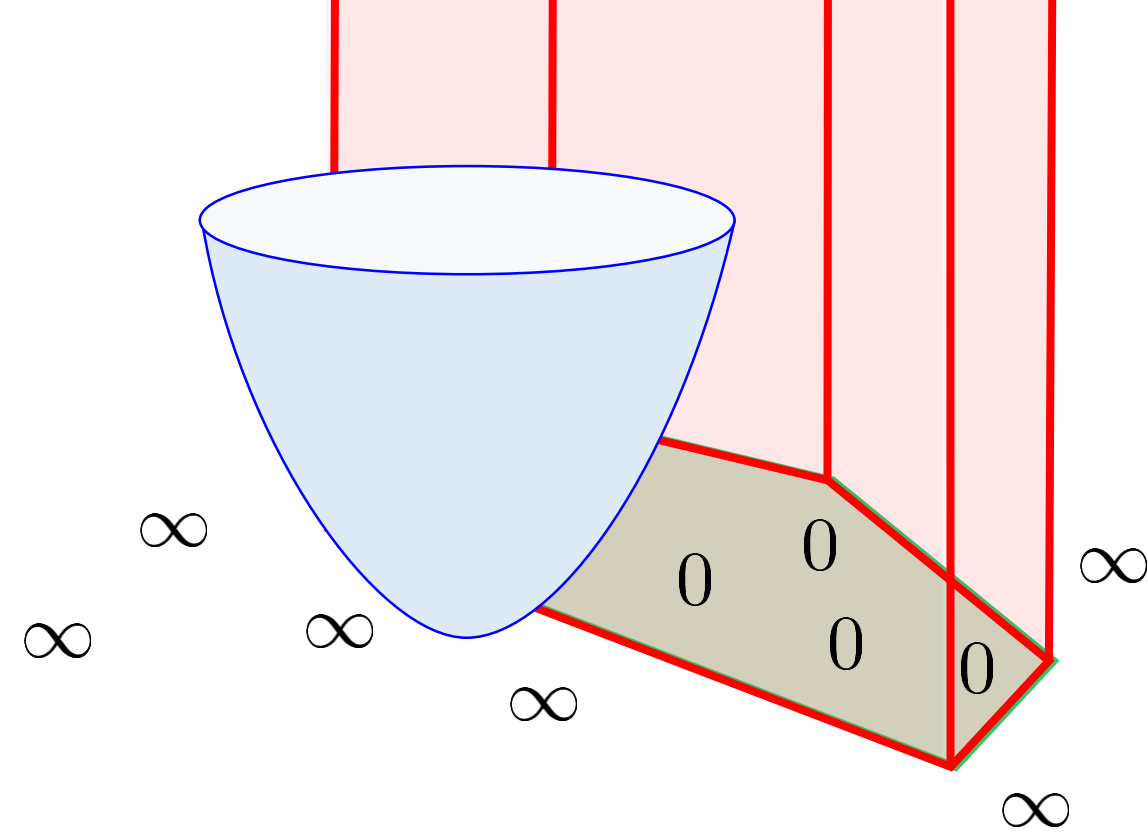


The Lagrangian

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

Primal
problem



$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

Primal problem

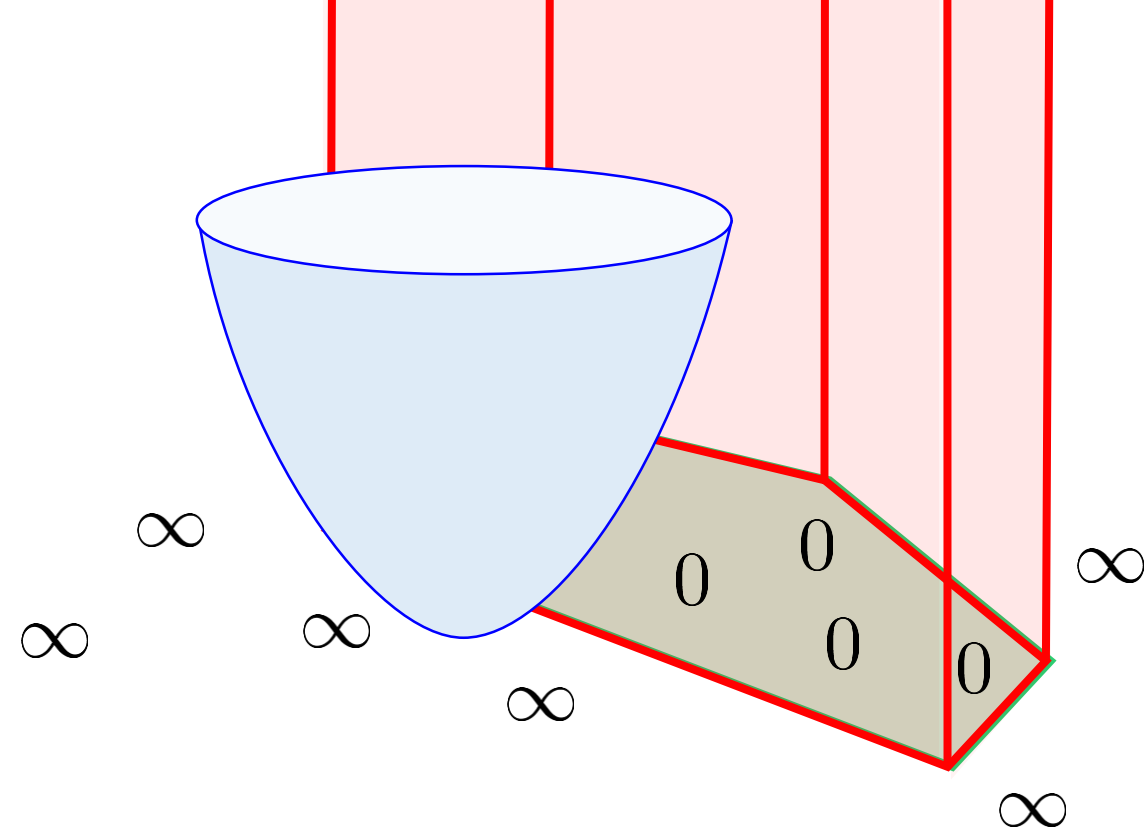
$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

The Lagrangian

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

Primal
problem

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

Primal problem

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

Dual problem

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

$$\hat{\mathbf{w}}_P \stackrel{?}{=} \hat{\mathbf{w}}_D$$

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\hat{\mathbf{w}}_P \stackrel{?}{=} \hat{\mathbf{w}}_D$$

$\hat{\mathbf{w}}_p = \hat{\mathbf{w}}_d$ for nice problems*

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

$$\hat{\mathbf{w}}_P \stackrel{?}{=} \hat{\mathbf{w}}_D$$

$\hat{\mathbf{w}}_p = \hat{\mathbf{w}}_d$ for nice problems*

E.g. if $f(\cdot)$ is convex and $\{\mathbf{w} : g(\mathbf{w}) \leq 0\}$ is convex

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

$$\alpha_D \cdot g(\hat{\mathbf{w}}_D) = 0$$

$$\hat{\mathbf{w}}_P \stackrel{?}{=} \hat{\mathbf{w}}_D$$

$\hat{\mathbf{w}}_p = \hat{\mathbf{w}}_d$ for nice problems*

E.g. if $f(\cdot)$ is convex and $\{\mathbf{w} : g(\mathbf{w}) \leq 0\}$ is convex

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

$$\alpha_D \cdot g(\hat{\mathbf{w}}_D) = 0$$

$$\hat{\mathbf{w}}_P \stackrel{?}{=} \hat{\mathbf{w}}_D$$

$\hat{\mathbf{w}}_p = \hat{\mathbf{w}}_d$ for nice problems*

E.g. if $f(\cdot)$ is convex and $\{\mathbf{w} : g(\mathbf{w}) \leq 0\}$ is convex

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

The
maximizer

$$\alpha_D \cdot g(\hat{\mathbf{w}}_D) = 0$$

$$\hat{\mathbf{w}}_P \stackrel{?}{=} \hat{\mathbf{w}}_D$$

$\hat{\mathbf{w}}_p = \hat{\mathbf{w}}_d$ for nice
problems*

E.g. if $f(\cdot)$ is convex and
 $\{\mathbf{w} : g(\mathbf{w}) \leq 0\}$ is convex

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

The
maximizer

$$\alpha_D \cdot g(\hat{\mathbf{w}}_D) = 0$$

$$\hat{\mathbf{w}}_P \stackrel{?}{=} \hat{\mathbf{w}}_D$$

$\hat{\mathbf{w}}_p = \hat{\mathbf{w}}_d$ for nice
problems*

E.g. if $f(\cdot)$ is convex and
 $\{\mathbf{w} : g(\mathbf{w}) \leq 0\}$ is convex

Complimentary slackness
KKT Condition

Please give your Feedback

<http://tinyurl.com/ml17-18afb>