

Function Approximation Methods-III

CS771: Introduction to Machine Learning
Purushottam Kar



Discussion Session

Sunday, 03 September, 2017

1800–1930 hrs, RM101, CSE department

Post query on <http://tinyurl.com/ml17-18ads1> before coming!

Recap

Multi-classification Loss Functions

One-vs-All (OVA)

$$\hat{y}^i = \arg \max_{j \in [K]} \langle \mathbf{w}^j, \mathbf{x}^i \rangle$$

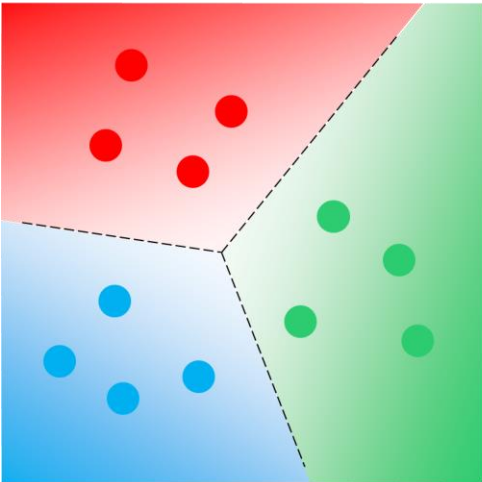
$$\mathbf{W} = [\mathbf{w}^1, \mathbf{w}^2, \dots, \mathbf{w}^K]$$

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W}} \sum_{k=1}^K \|\mathbf{w}^k\|_2^2 + \sum_{i=1}^n \ell_{\text{cs}}(y^i, \boldsymbol{\eta}^i)$$

$$\boldsymbol{\eta}^i = \langle \mathbf{W}, \mathbf{x}^i \rangle = [\eta_1^i, \eta_2^i, \dots, \eta_K^i]$$

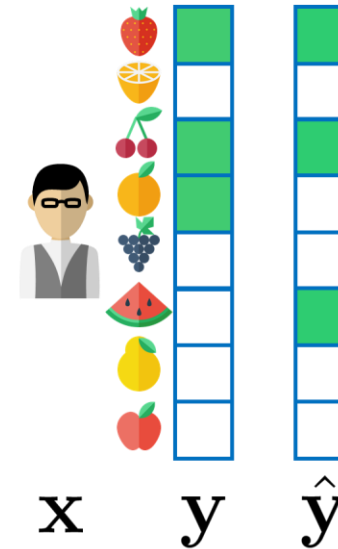
$$\ell_{\text{cs}}(y, \boldsymbol{\eta}) = [1 + \max_{k \neq y} \eta_k - \eta_y]_+$$

Crammer-Singer
loss function



August 23, 2017

Multi-label Classification Loss Functions



August 23, 2017

$$\frac{\left| \left\{ i : \begin{matrix} \hat{y}_i = 1 \\ y_i = 1 \end{matrix} \right\} \right|}{\left| \left\{ i : \begin{matrix} \hat{y}_i = -1 \\ y_i = 1 \end{matrix} \right\} \right|} \quad \left| \left\{ i : \begin{matrix} \hat{y}_i = 1 \\ y_i = -1 \end{matrix} \right\} \right| \quad \left| \left\{ i : \begin{matrix} \hat{y}_i = -1 \\ y_i = -1 \end{matrix} \right\} \right|$$

$$\ell_{\text{Hamming}}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{b + c}{a + b + c + d}$$

$$r_{\text{Precision}}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{a}{a + b} \quad r_{\text{Recall}}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{a}{a + c}$$

$$F(\mathbf{y}, \hat{\mathbf{y}}) = \frac{2r_{\text{Prec}} \cdot r_{\text{Rec}}}{r_{\text{Prec}} + r_{\text{Rec}}} = \frac{2a}{2a + b + c}$$

Can be used for evaluation and training

Historically

More recent

August 30, 2017



First-order Optimality

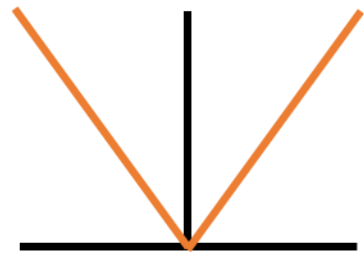
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$

$$\nabla f(\mathbf{w}) + \nabla r(\mathbf{w}) = \mathbf{0}$$

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

$$r(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$$

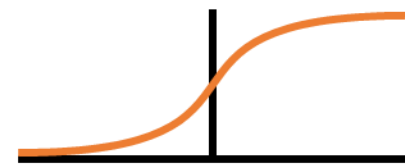
$$\mathbf{w} = (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$$



$$f(\mathbf{w}) = \sum_{i=1}^n \log (1 + \exp(-y^i \langle \mathbf{w}, \mathbf{x}^i \rangle))$$

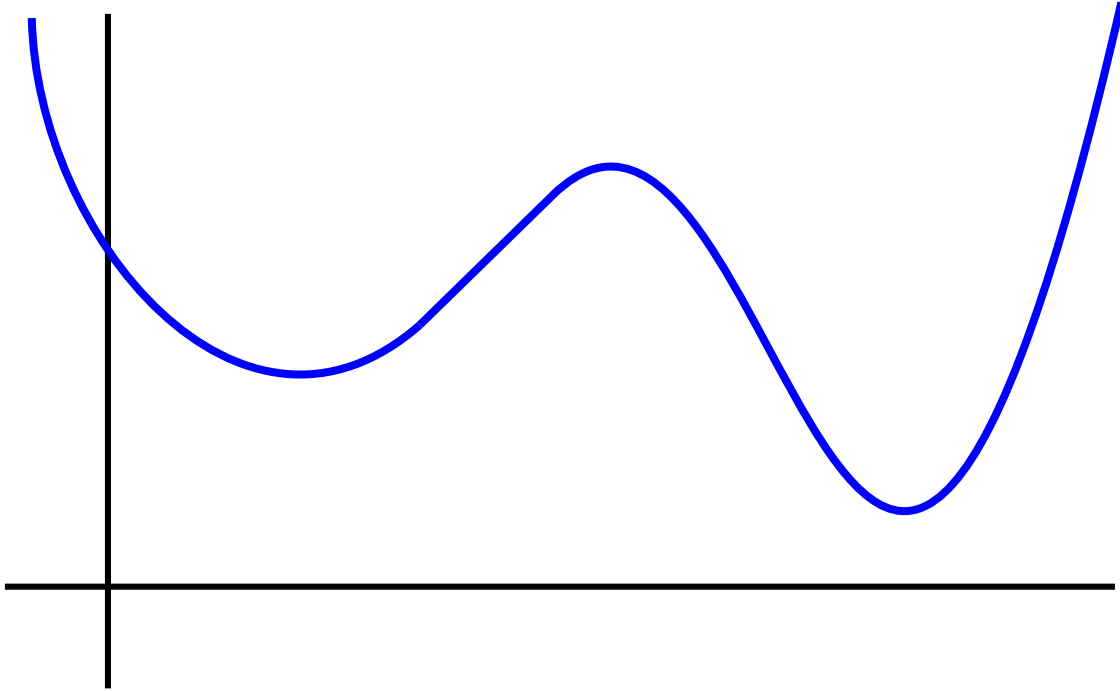
$$r(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$$

?

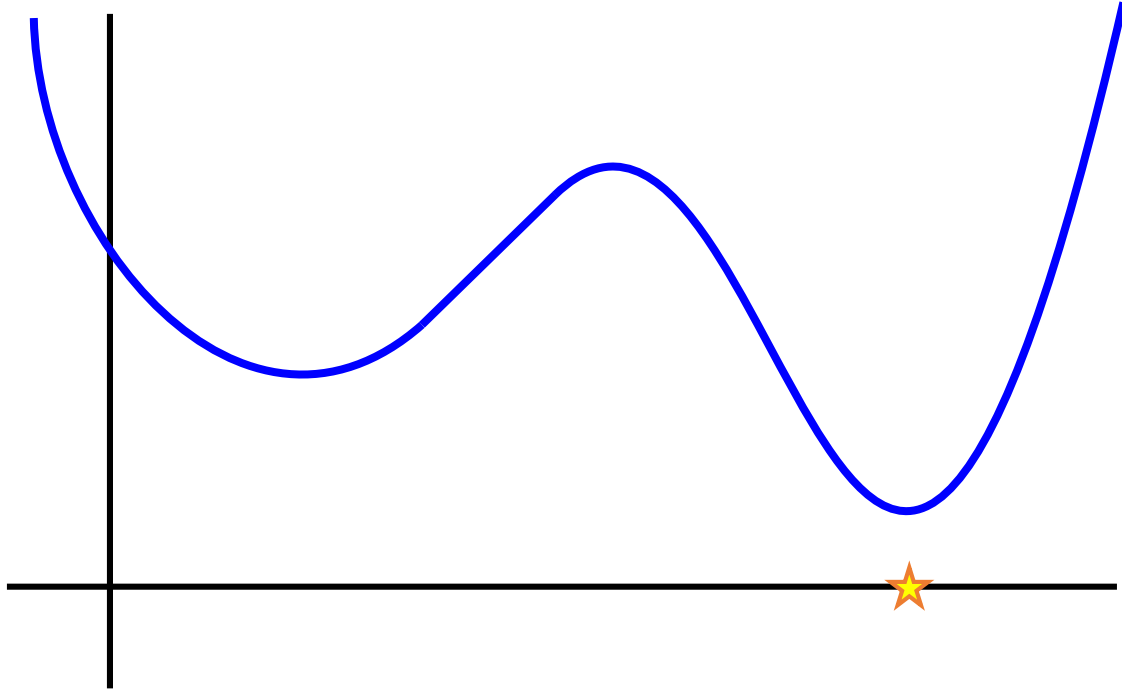


$$\sum_{i=1}^n (1 - \sigma(y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)) y^i \cdot \mathbf{x}^i + 2\lambda \cdot \mathbf{w} = \mathbf{0}$$

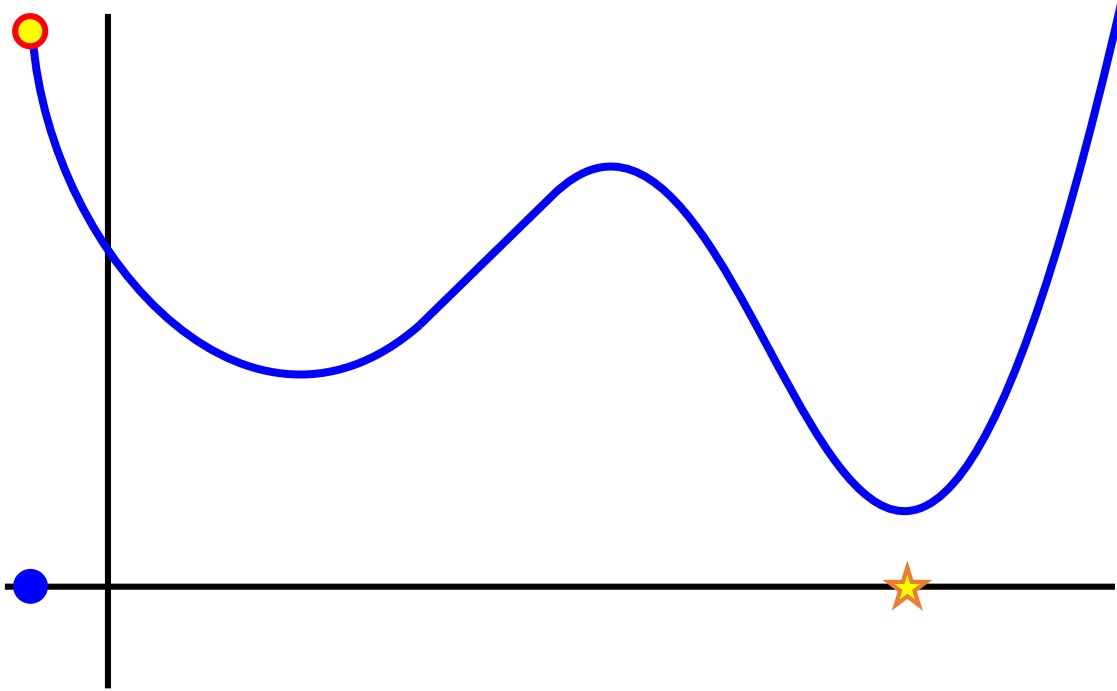
Gradient Descent



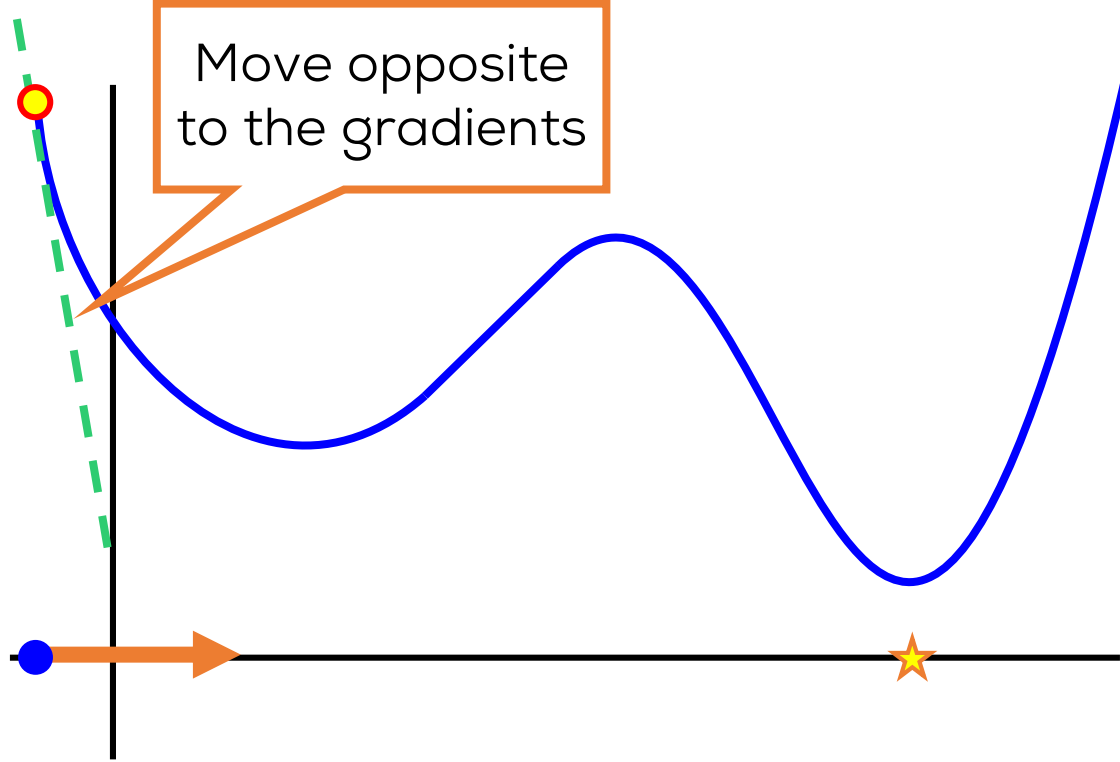
Gradient Descent



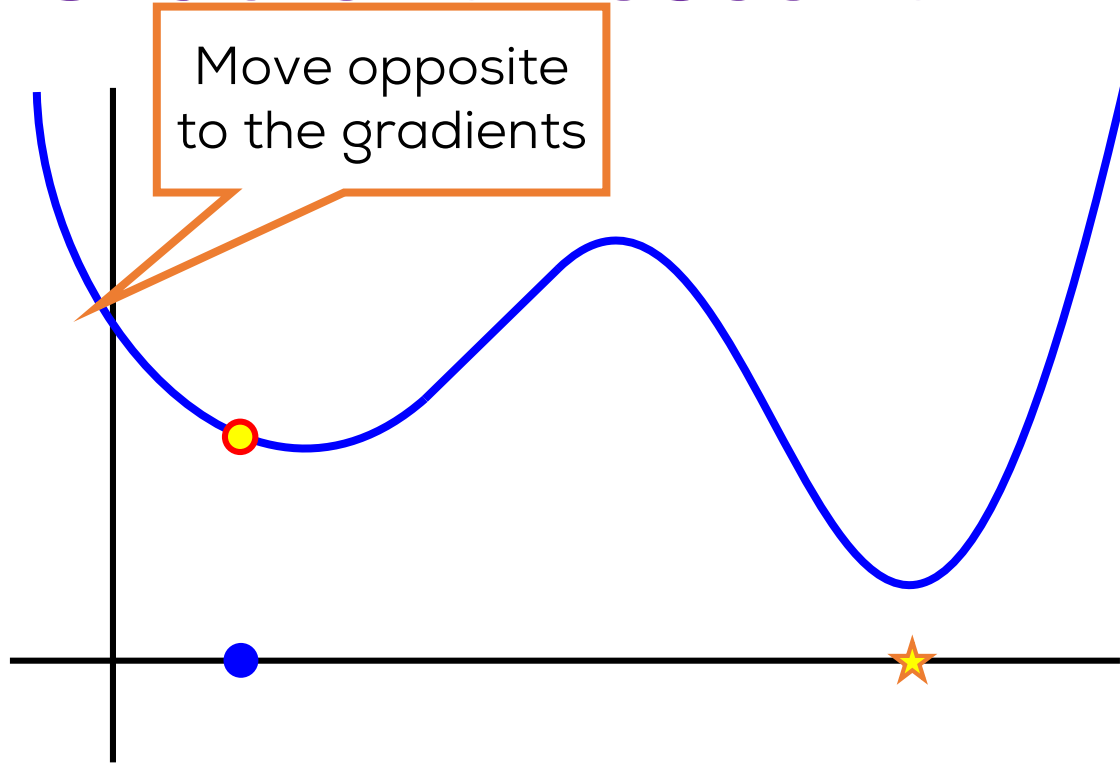
Gradient Descent



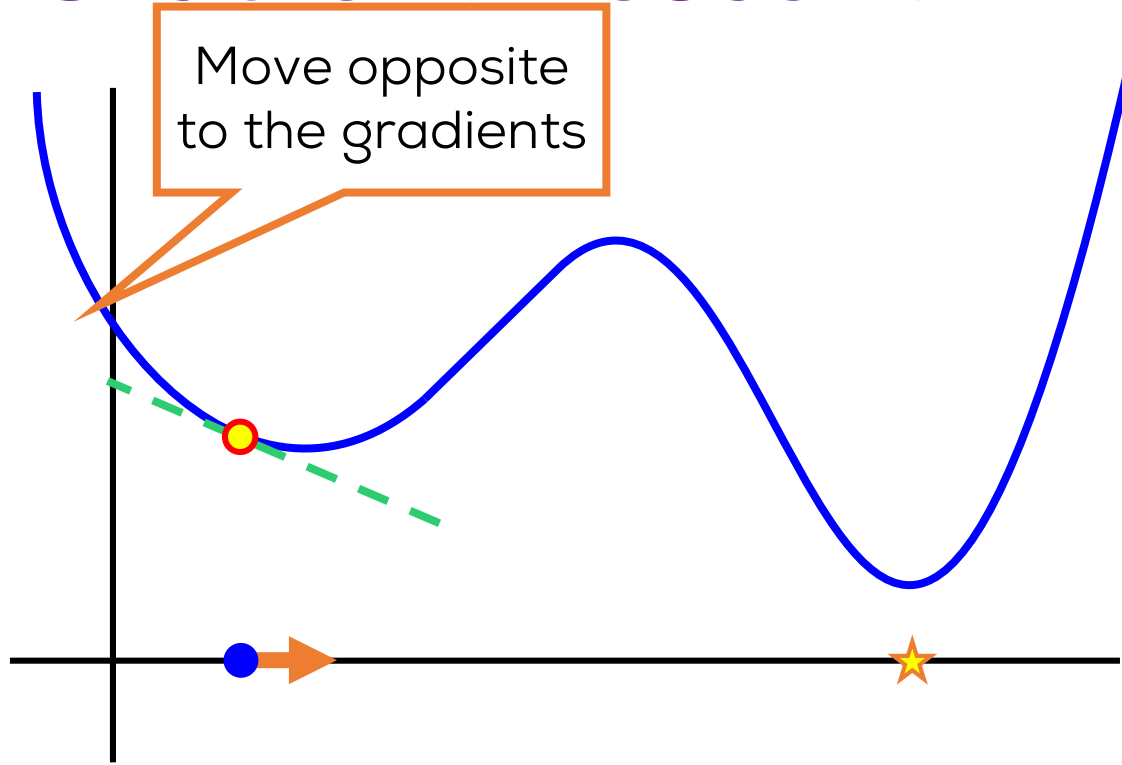
Gradient Descent



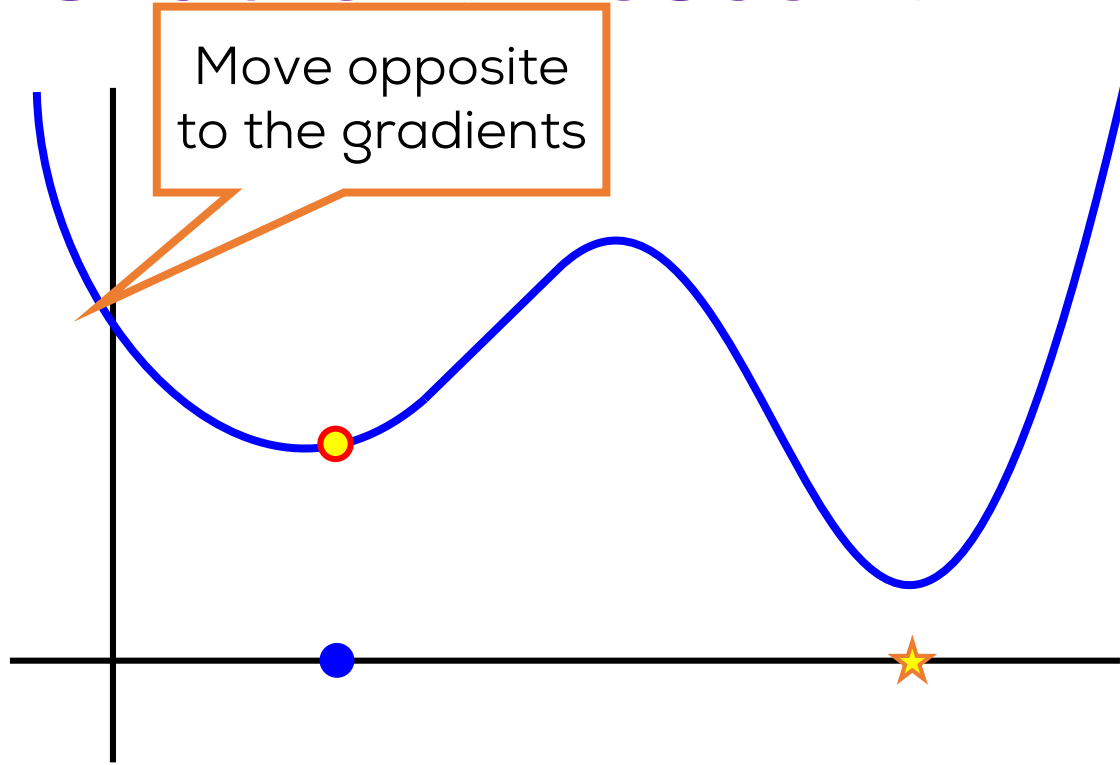
Gradient Descent



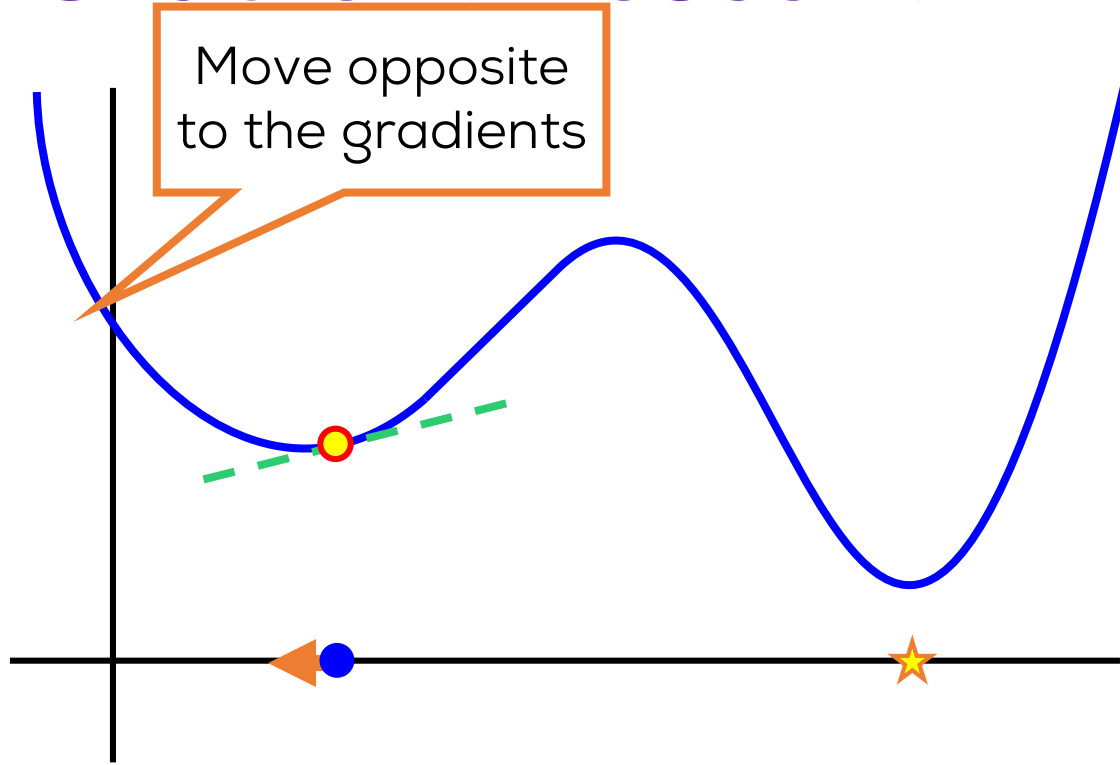
Gradient Descent



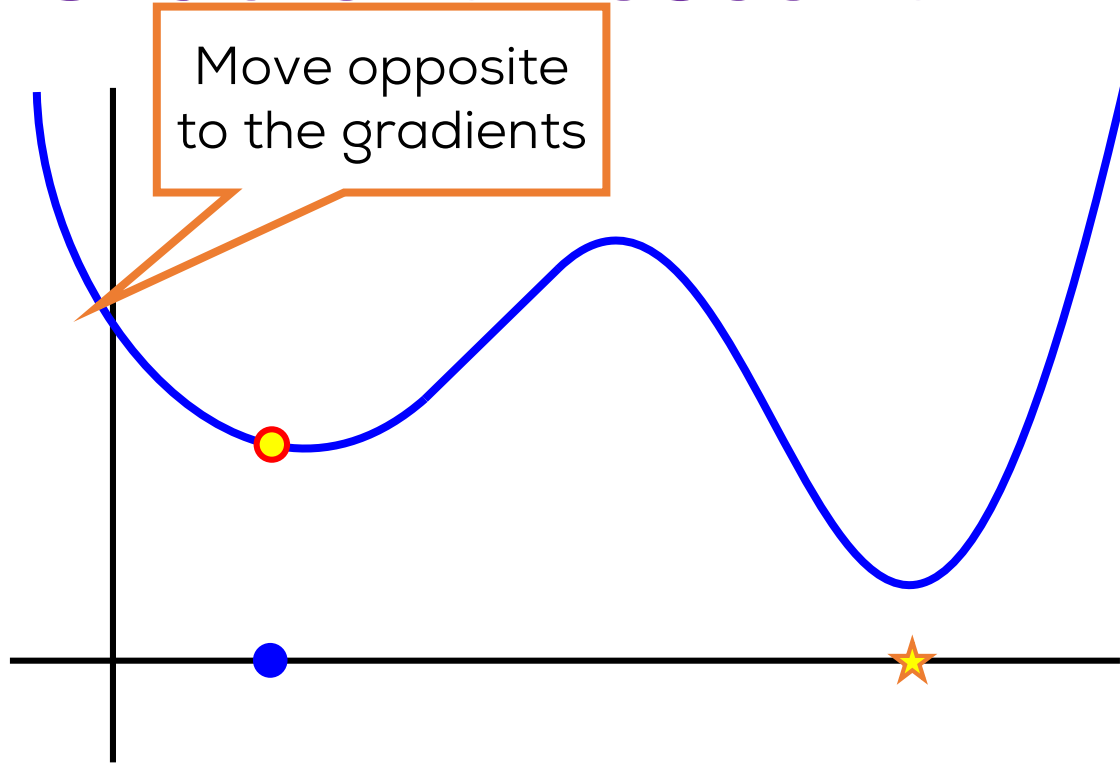
Gradient Descent



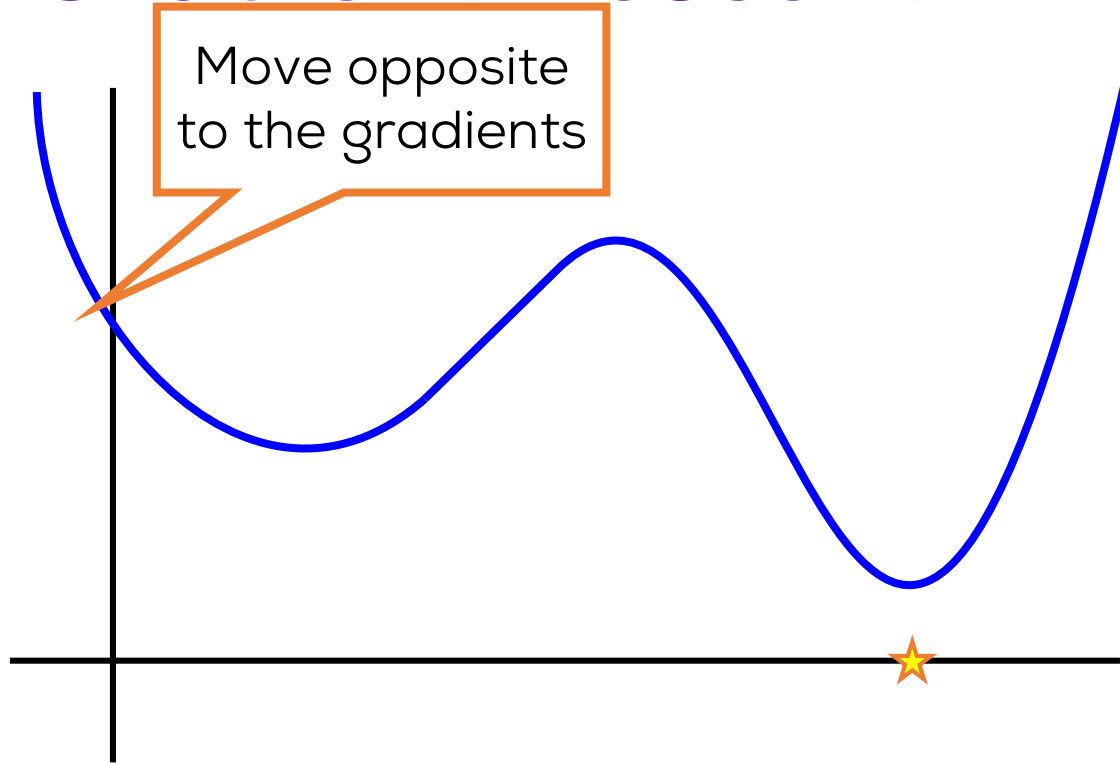
Gradient Descent



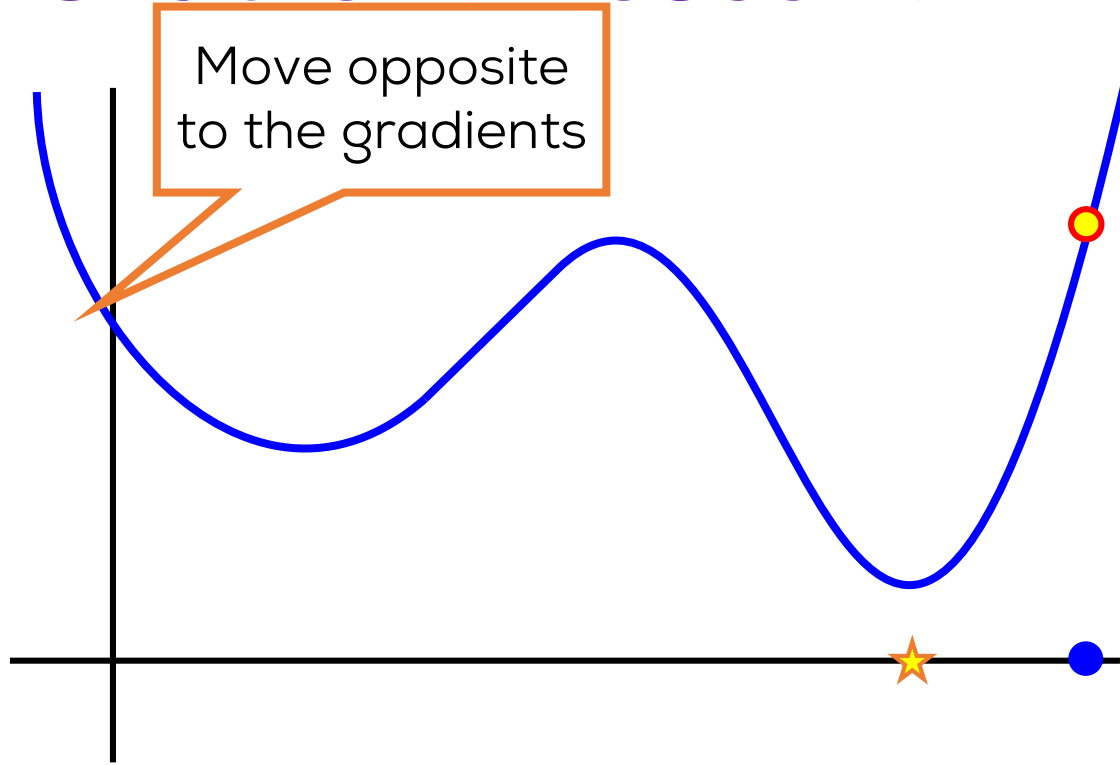
Gradient Descent



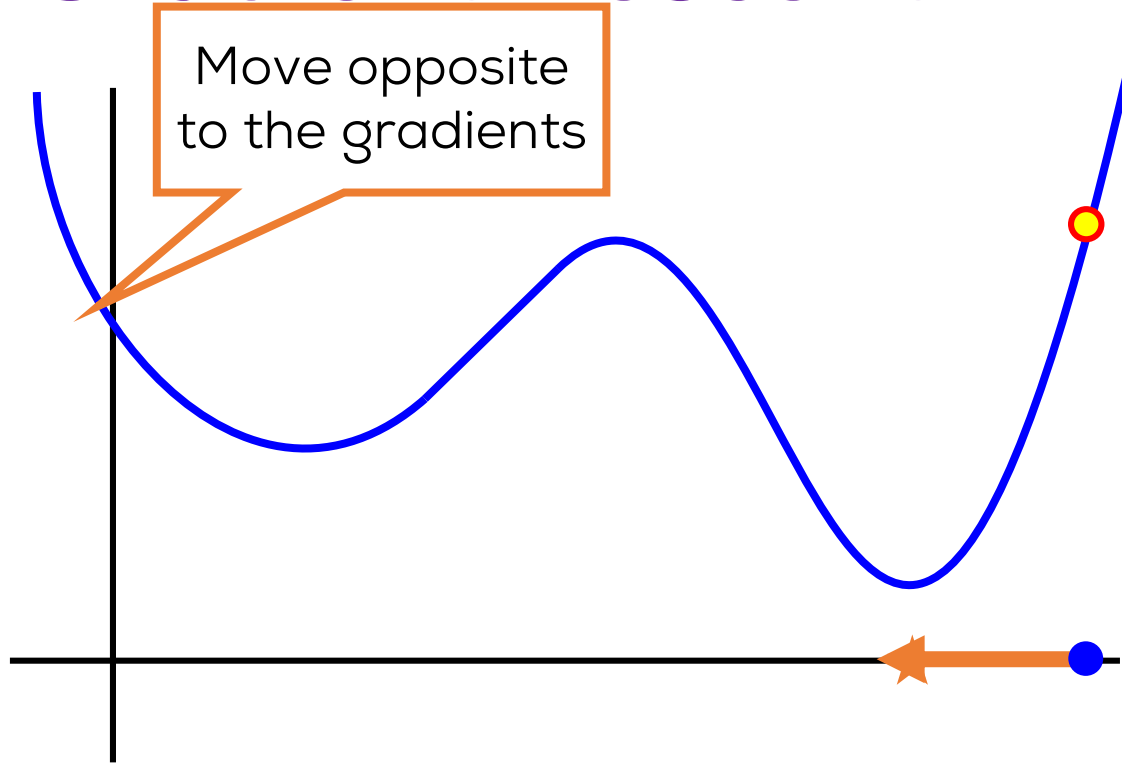
Gradient Descent



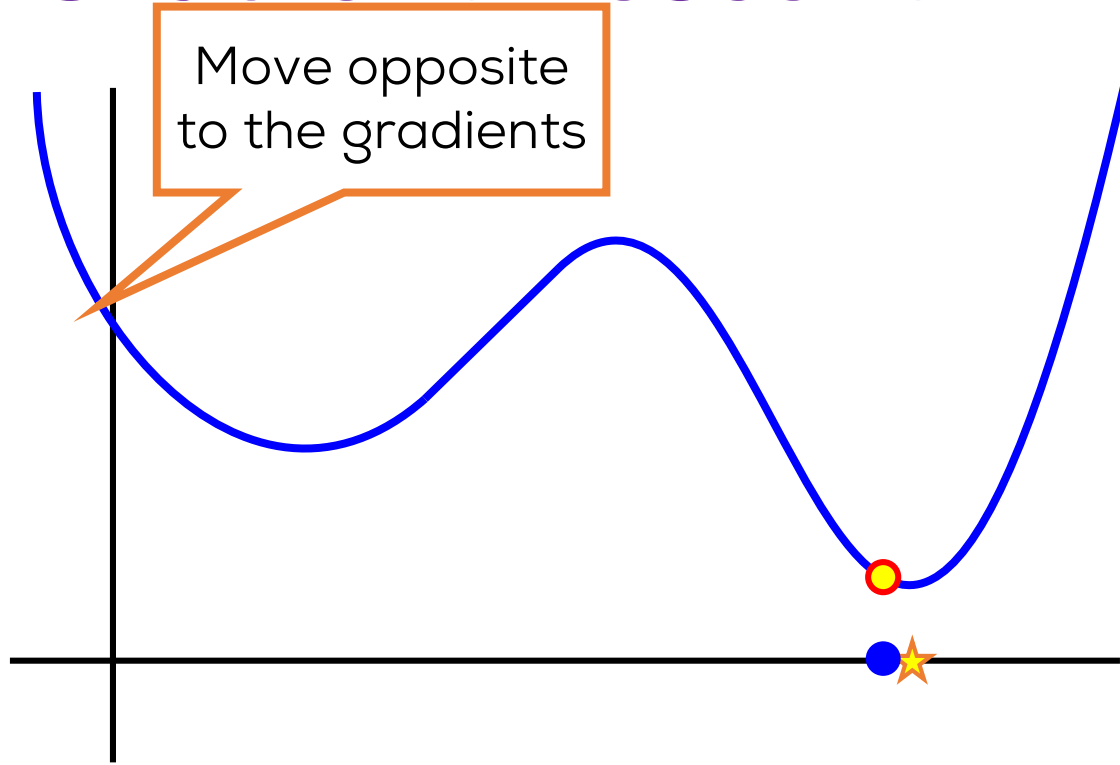
Gradient Descent



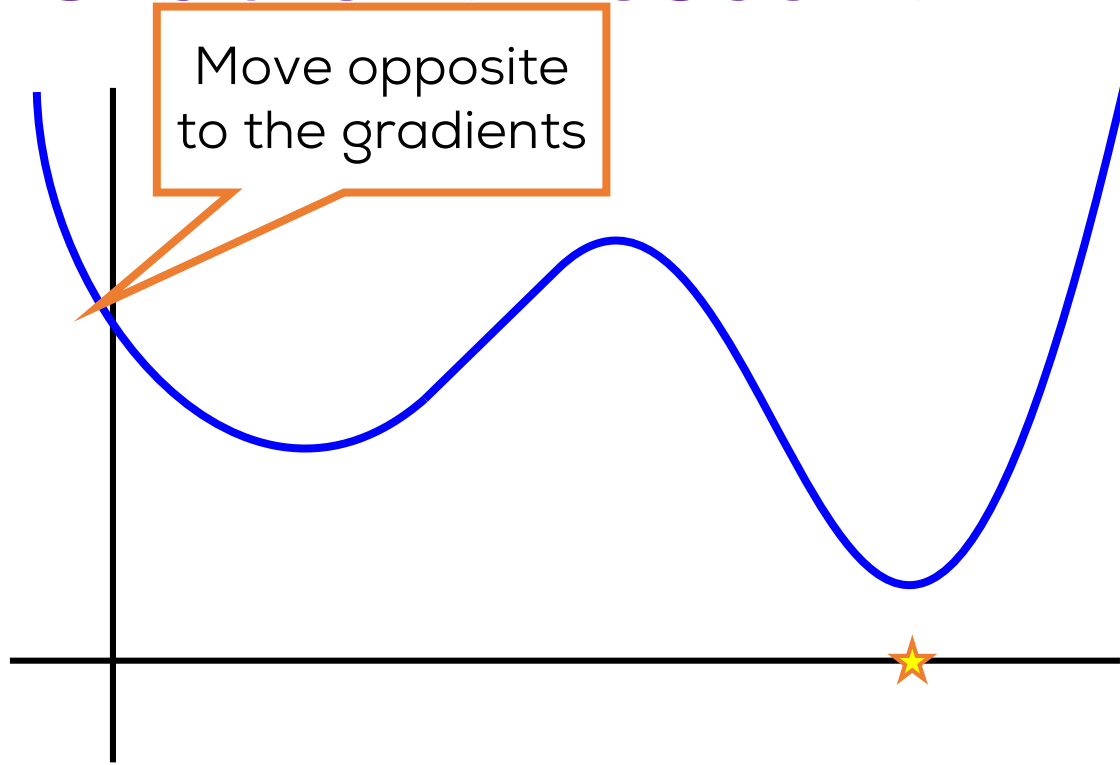
Gradient Descent



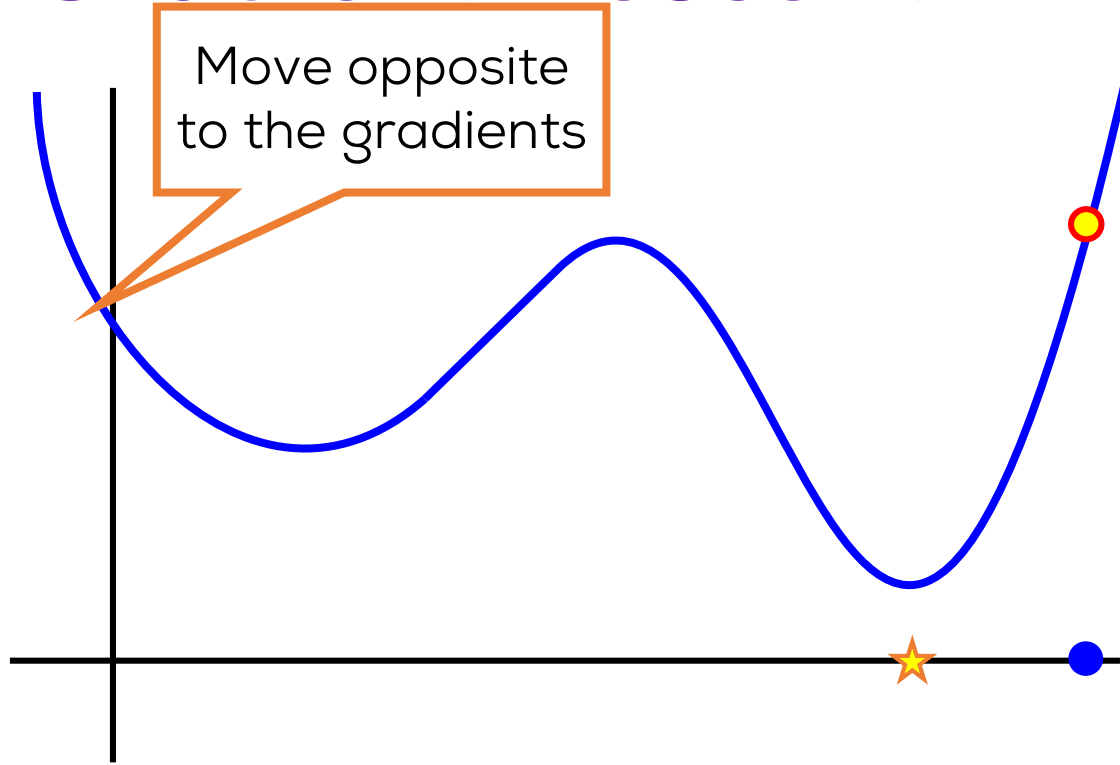
Gradient Descent



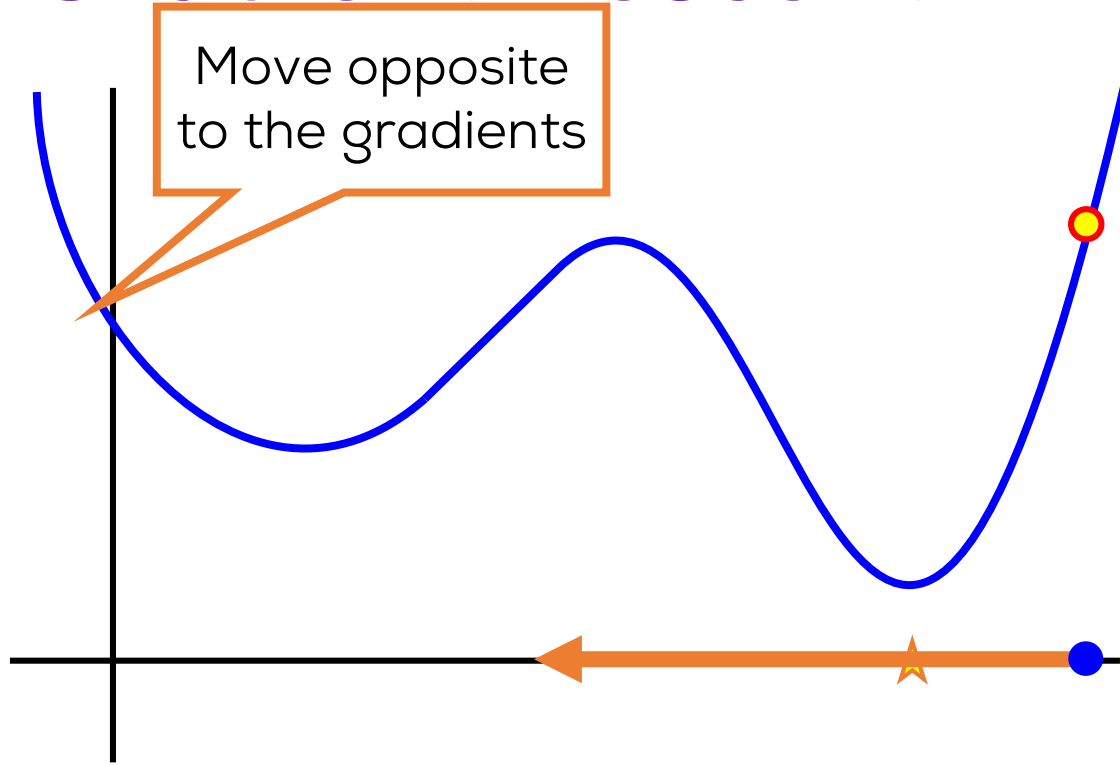
Gradient Descent



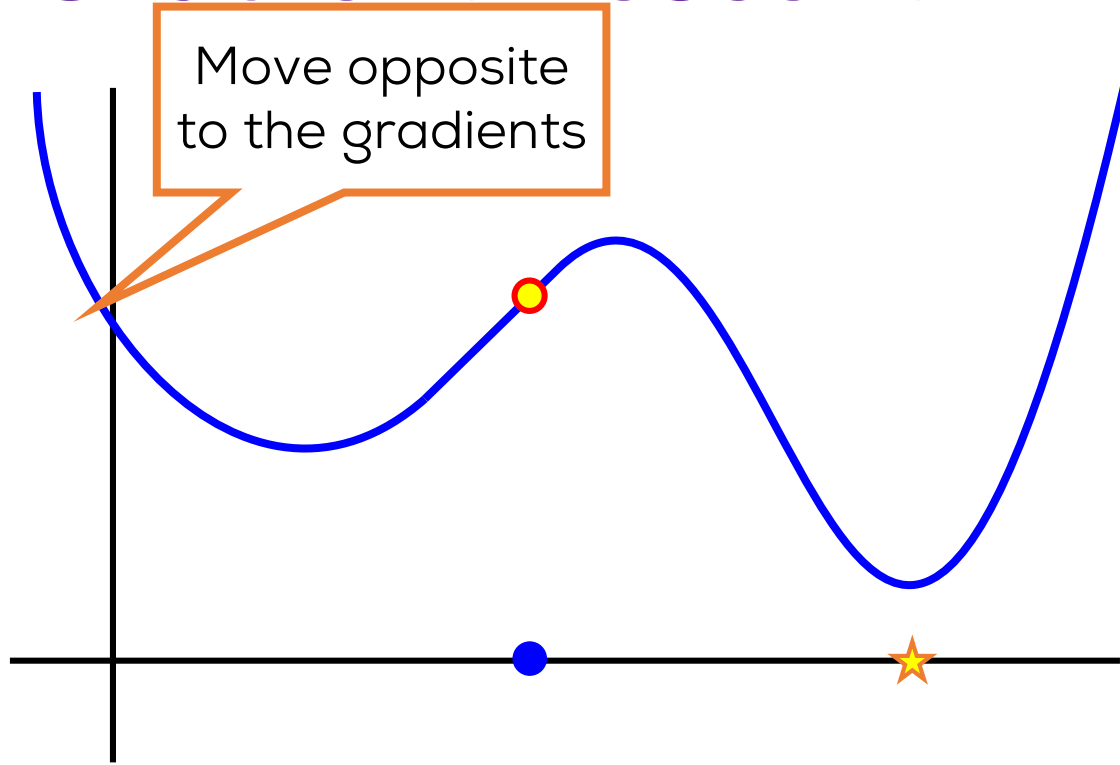
Gradient Descent



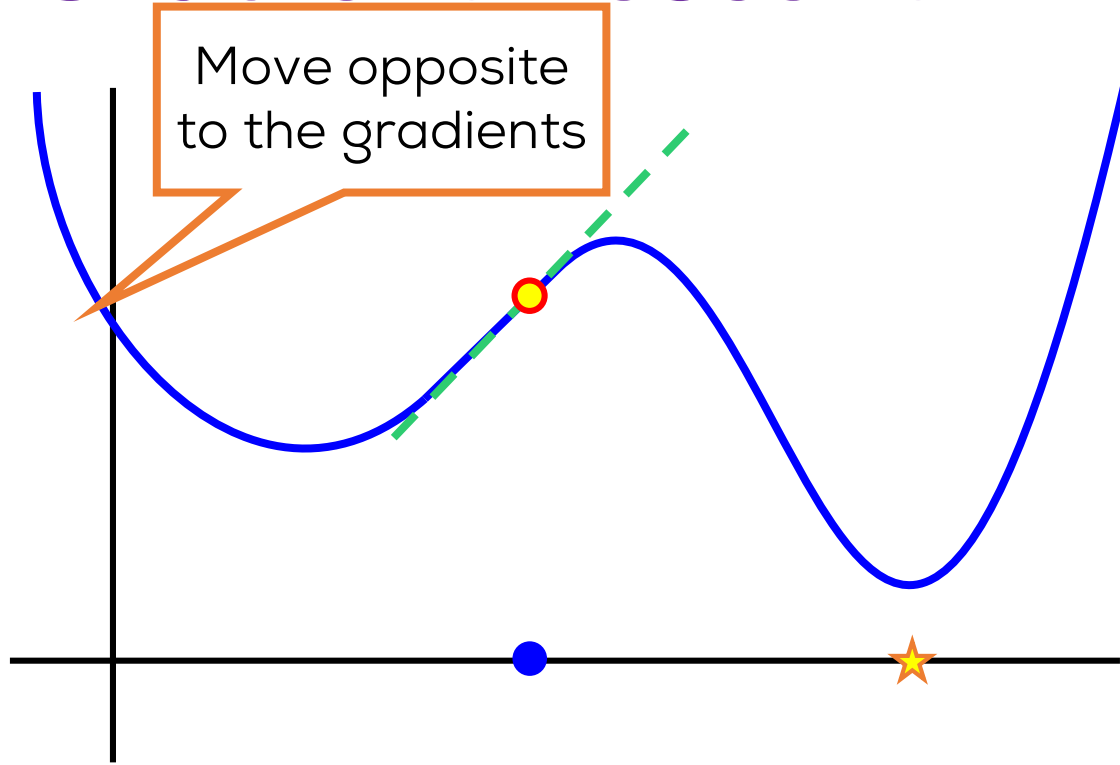
Gradient Descent



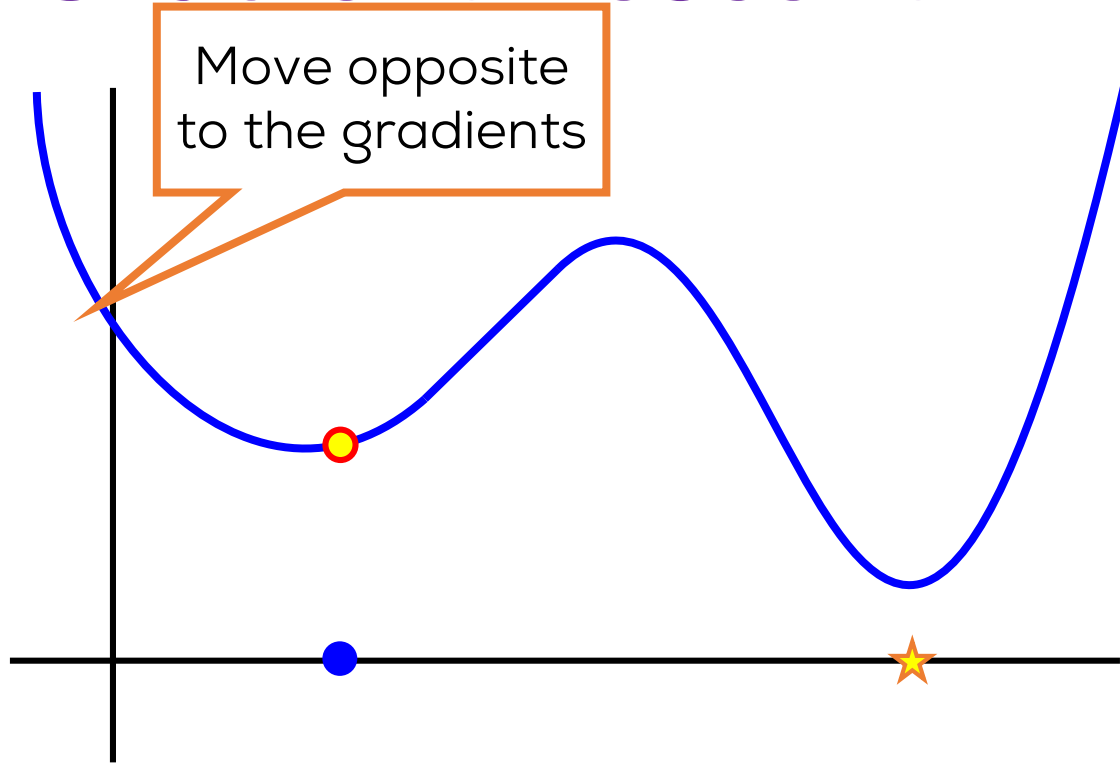
Gradient Descent



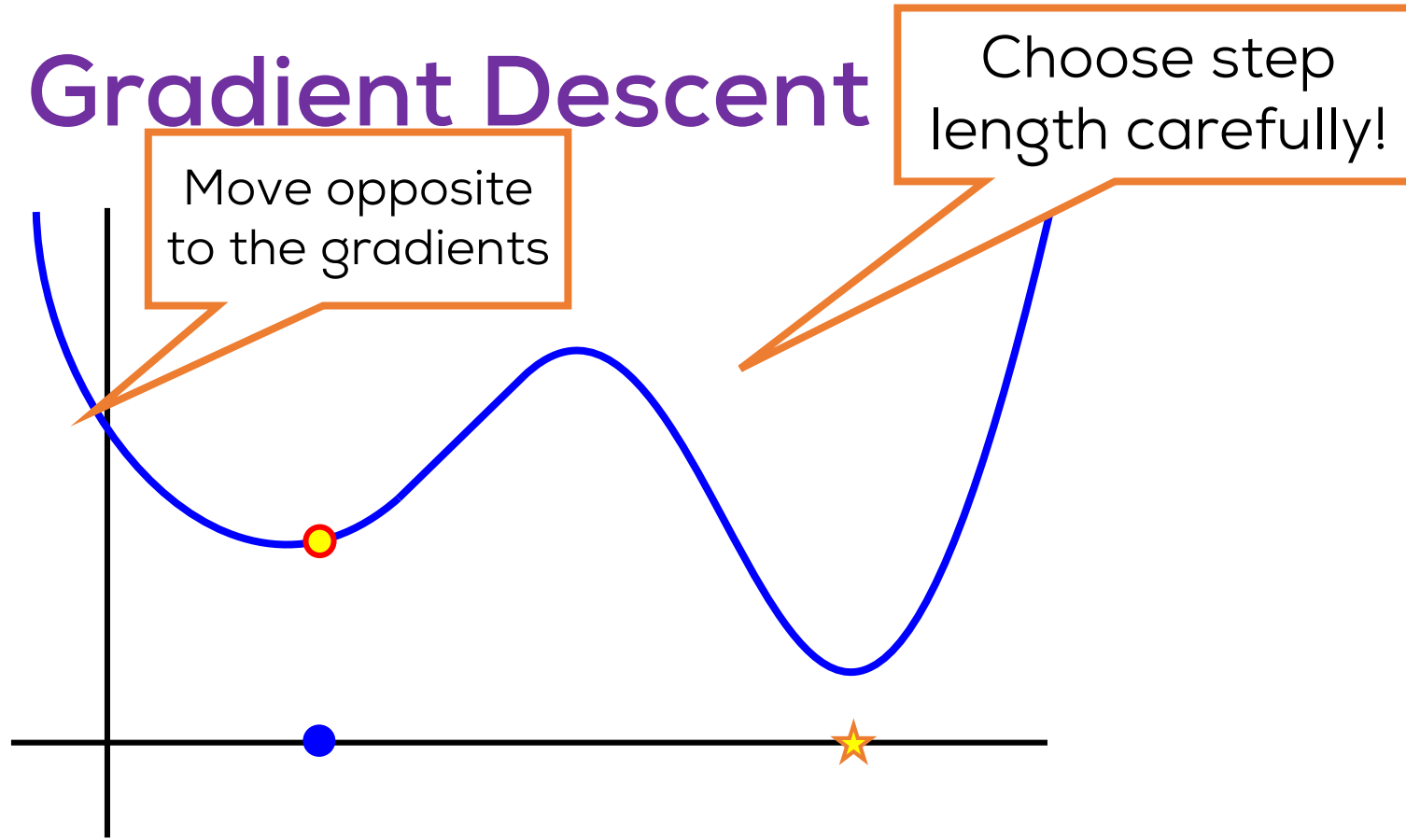
Gradient Descent



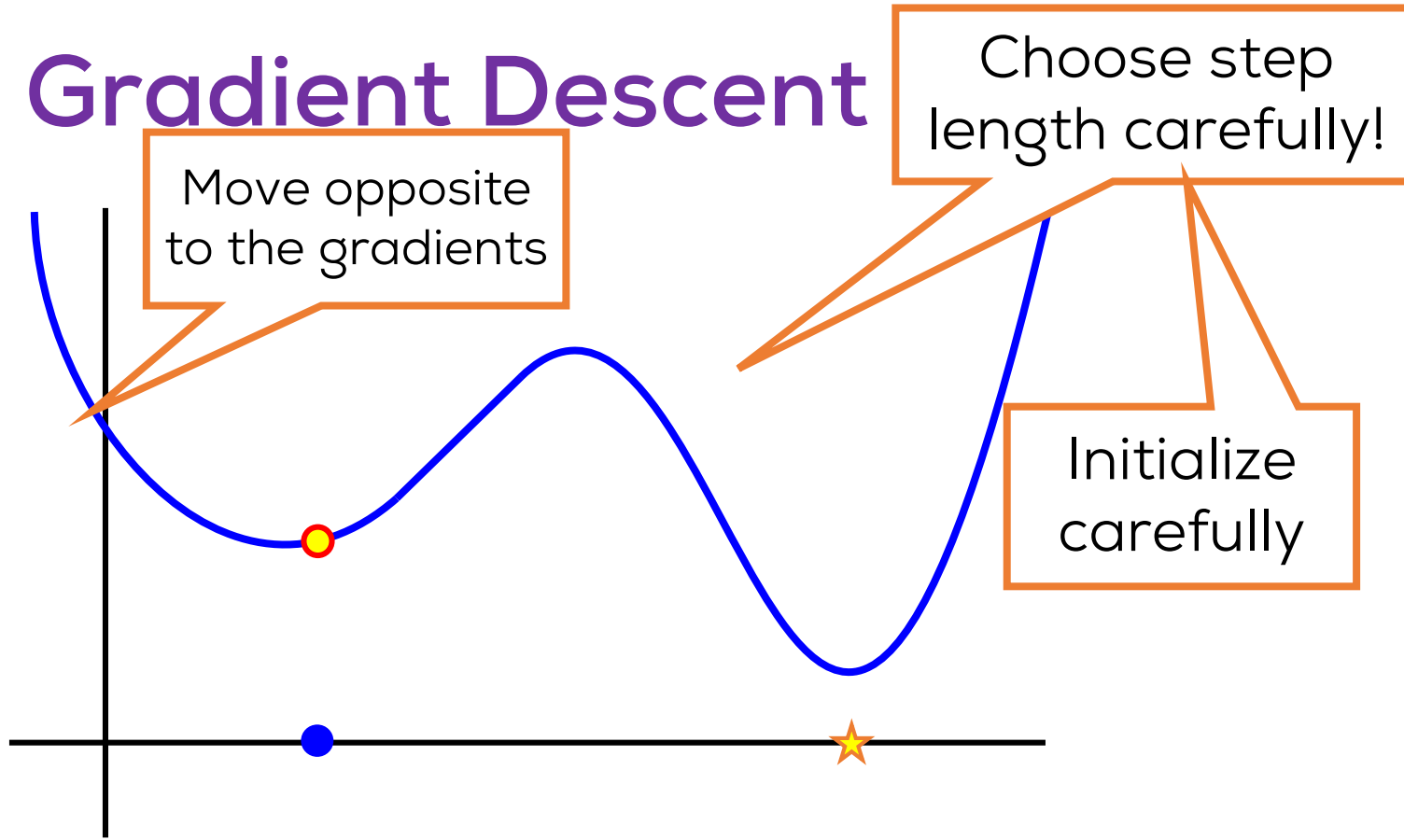
Gradient Descent



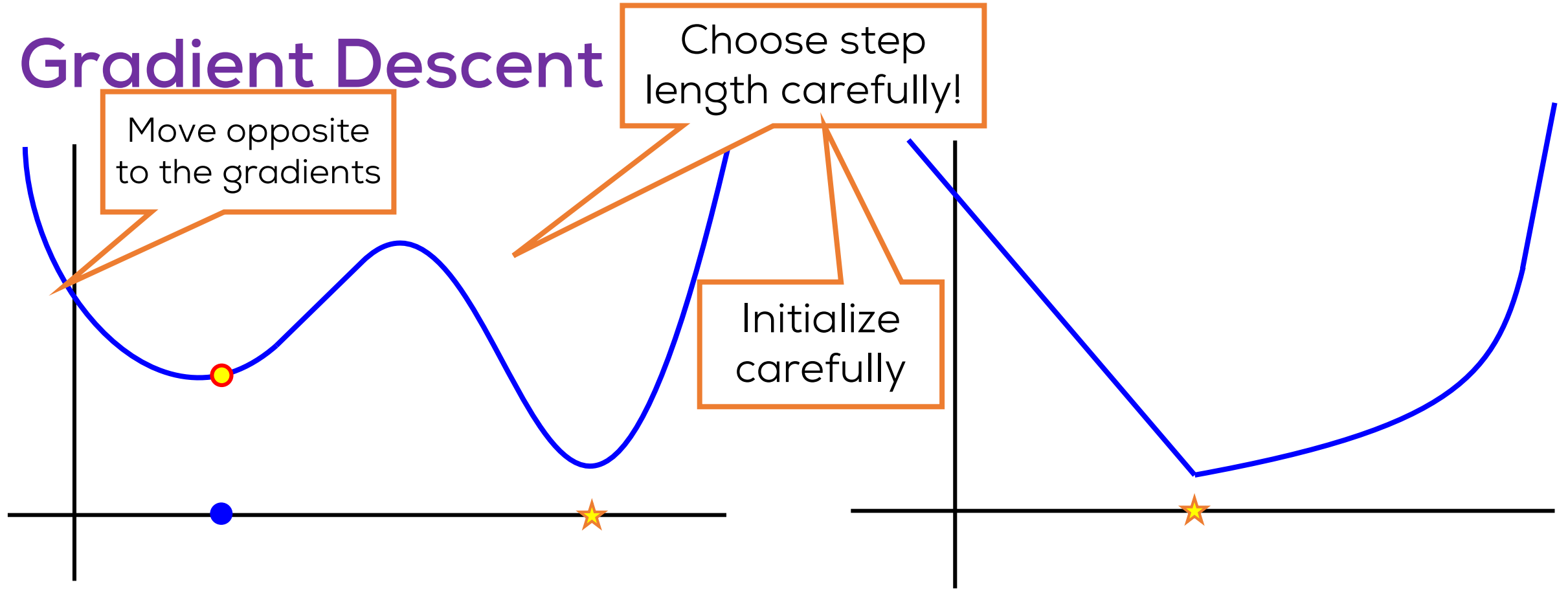
Gradient Descent



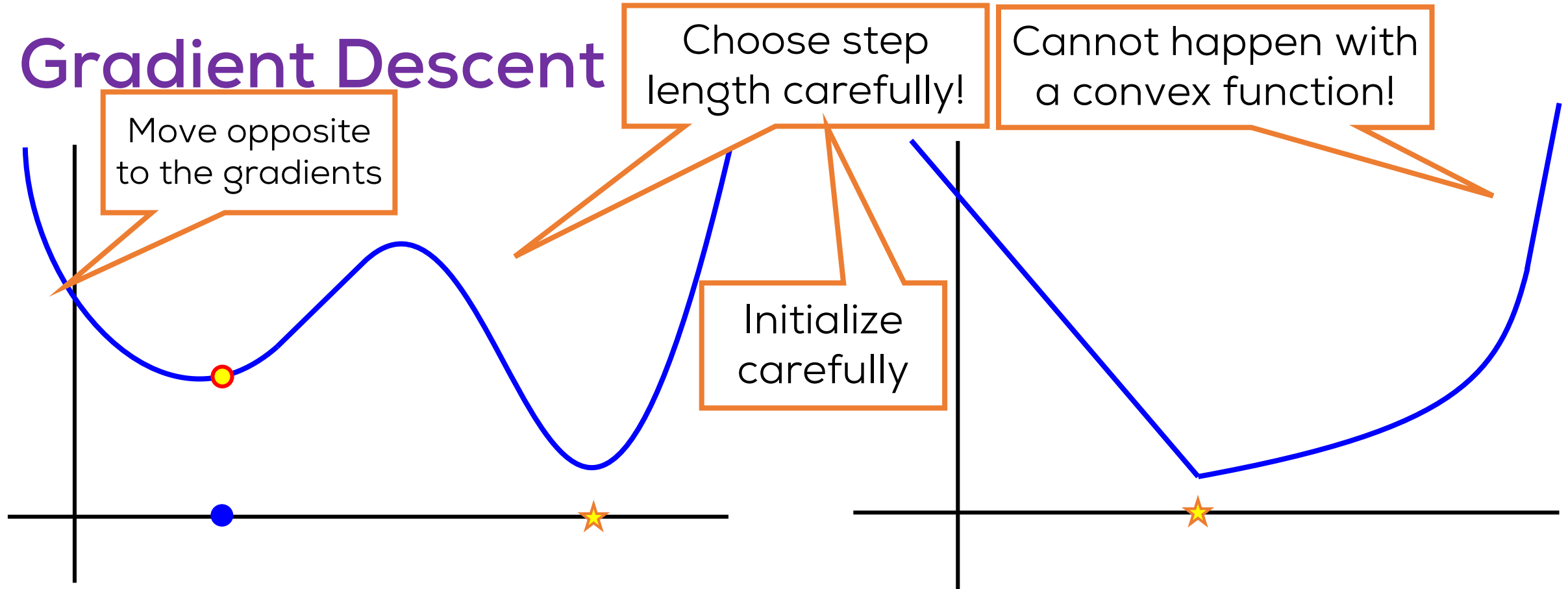
Gradient Descent



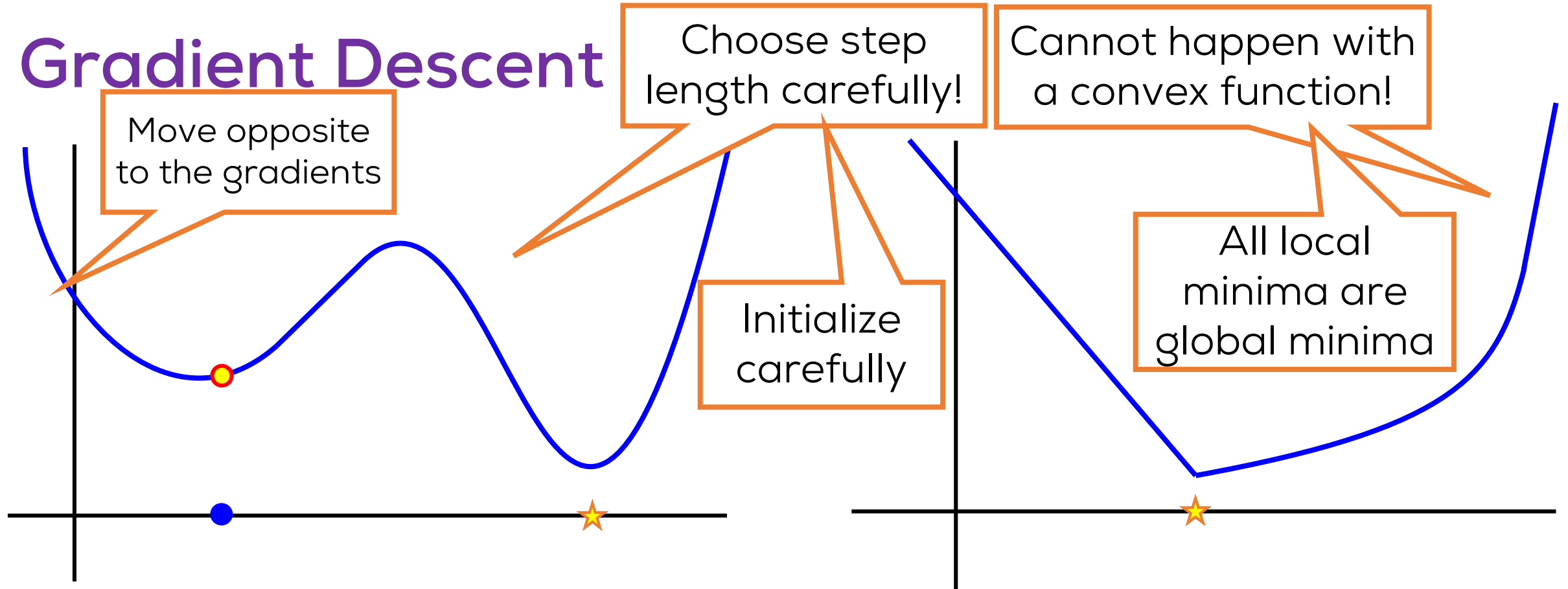
Gradient Descent



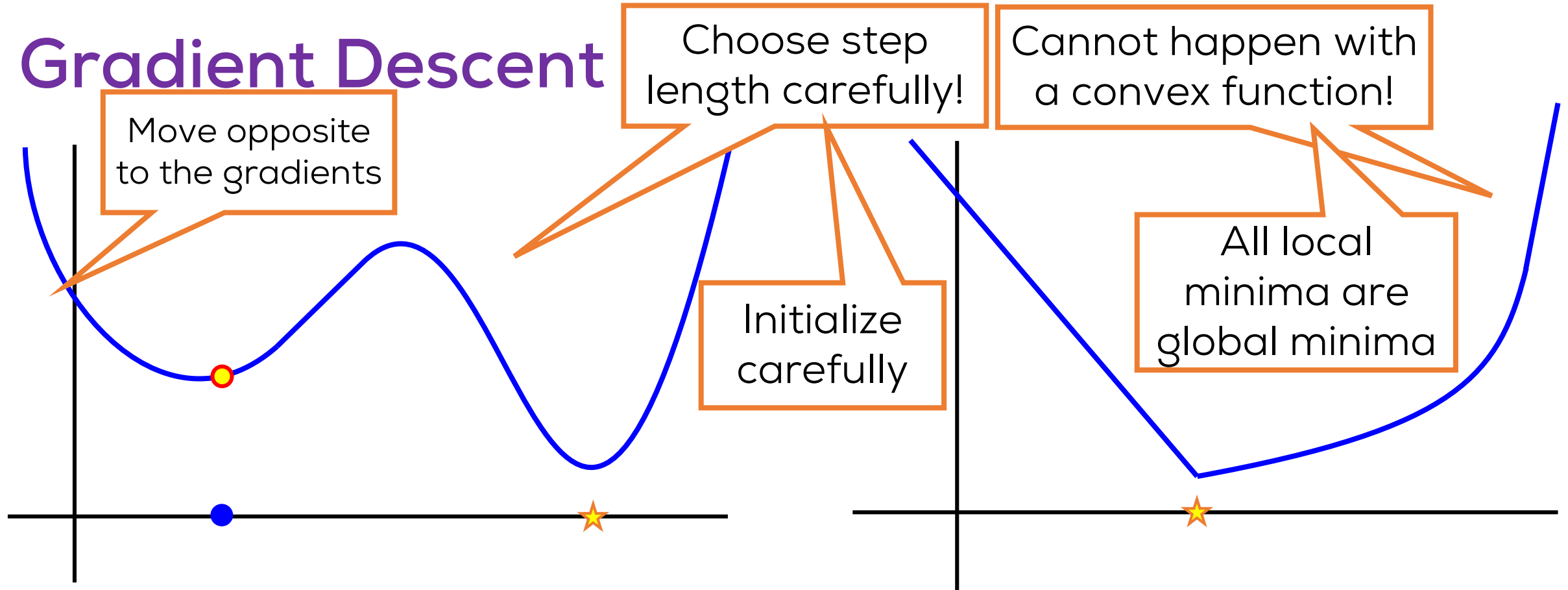
Gradient Descent



Gradient Descent



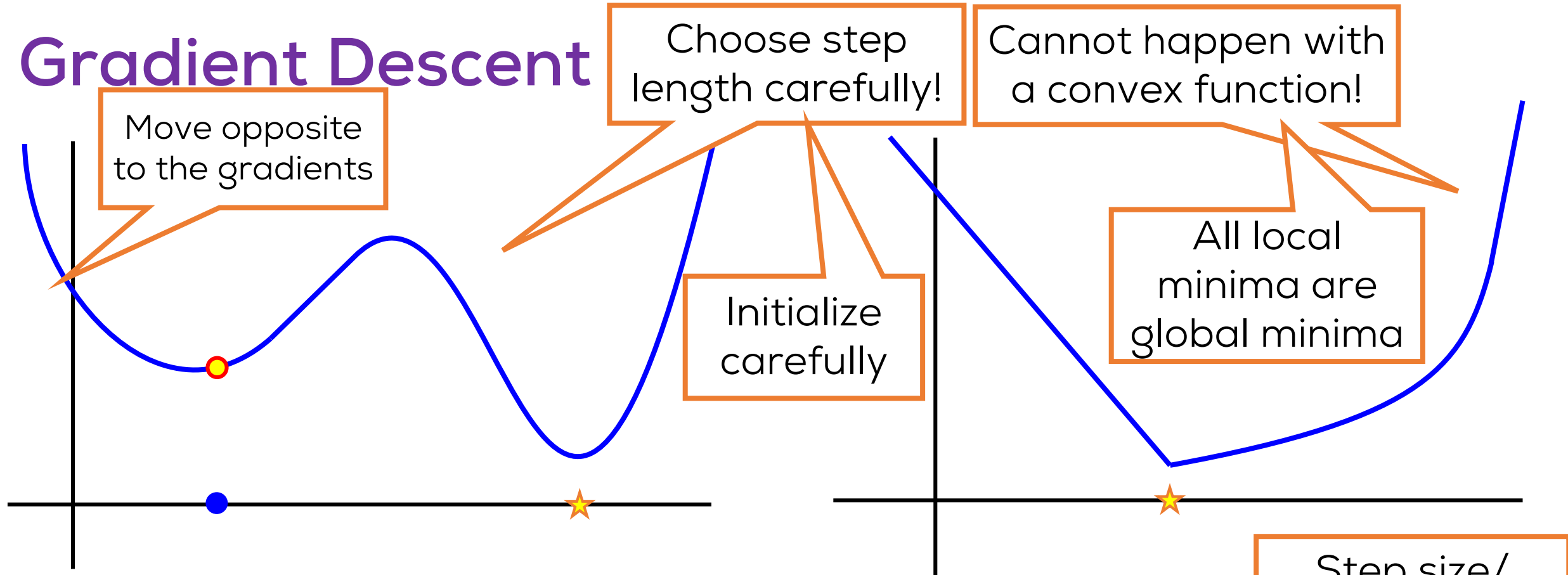
Gradient Descent



GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
3. Repeat until convergence

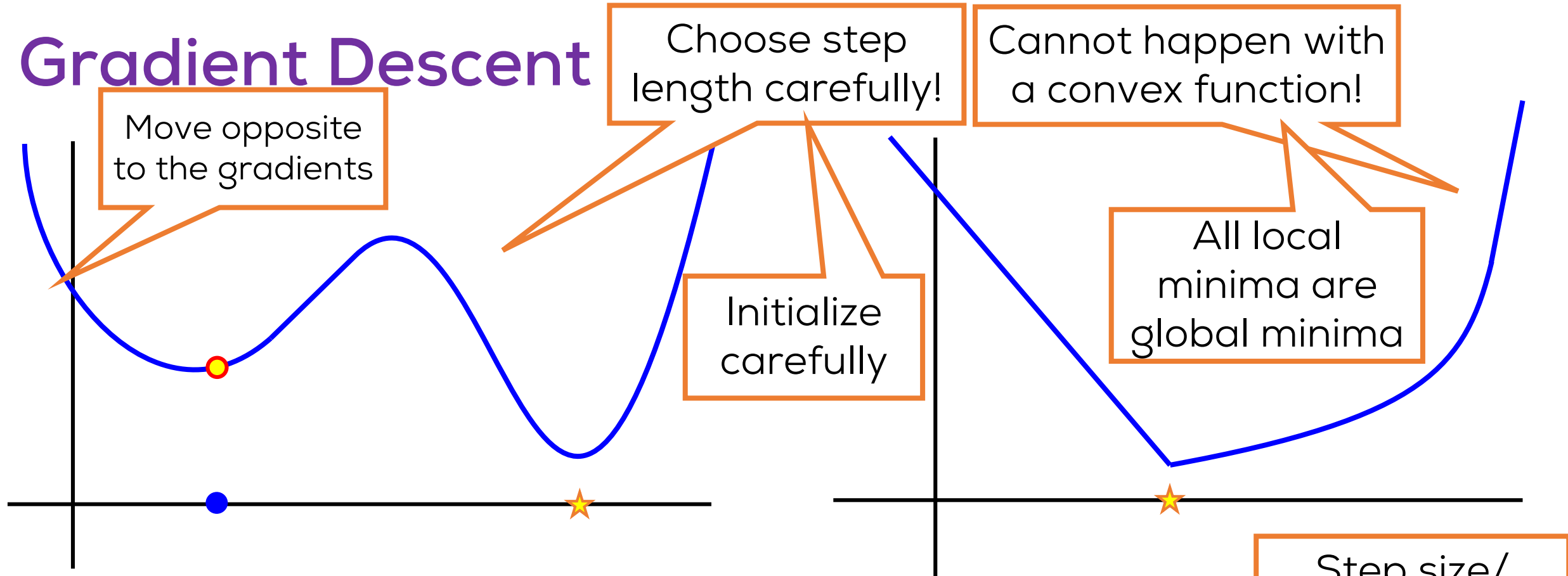
Gradient Descent



GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
3. Repeat until convergence

Gradient Descent



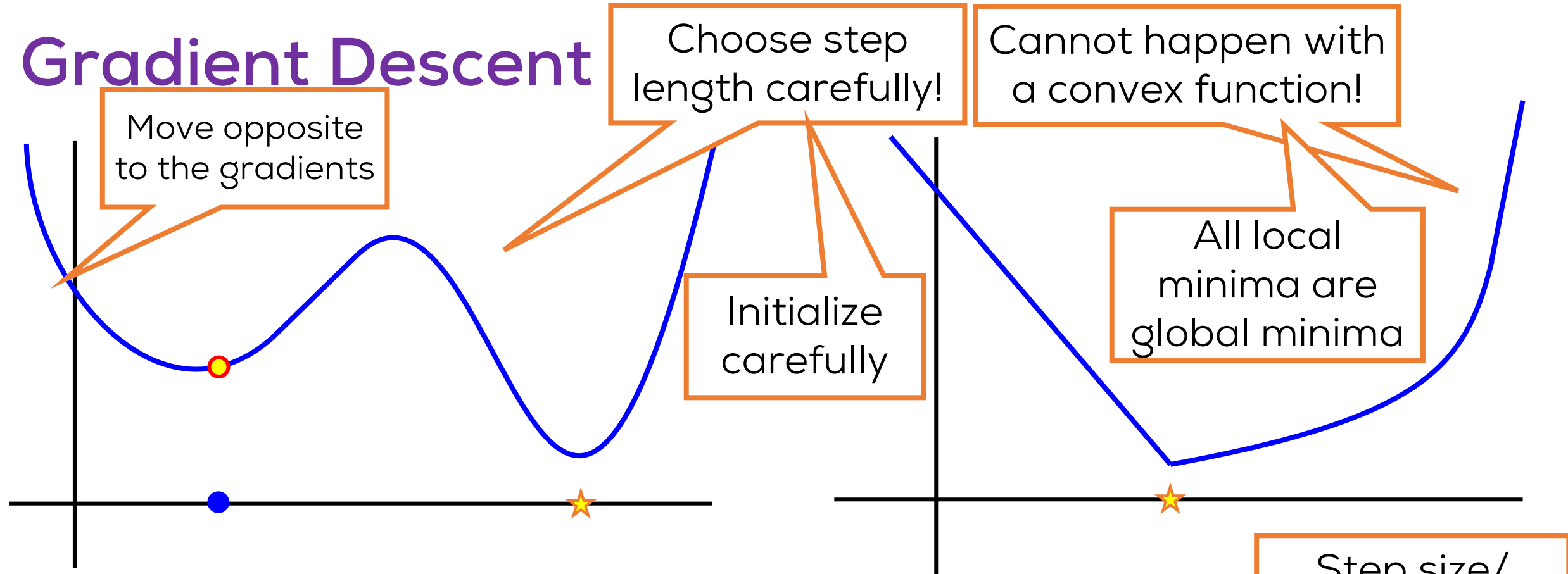
GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
3. Repeat until convergence

Step size/
learning rate

Tuned carefully
CV, Adam,
Adagrad

Gradient Descent



Many convergence criteria – length of gradient, performance threshold, dual criteria

GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
3. Repeat until convergence

Step size/
learning rate

Tuned carefully
CV, Adam,
Adagrad

Recap

First-order Optimality

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

$$\nabla f(\mathbf{w}) + \nabla r(\mathbf{w}) = \mathbf{0}$$

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

$$r(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{w} = (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$$

$$f(\mathbf{w}) = \sum_{i=1}^n \log(1 + \exp(-y^i \langle \mathbf{w}, \mathbf{x}^i \rangle))$$

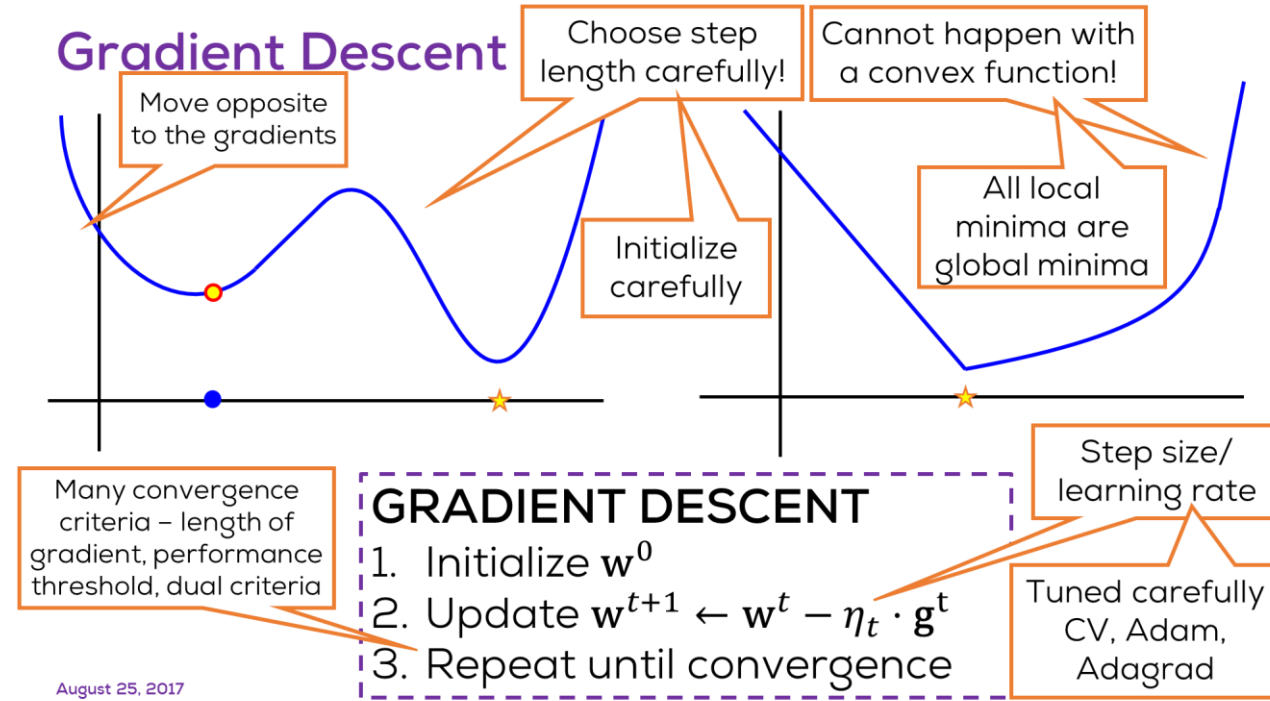
$$r(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\sum_{i=1}^n (1 - \sigma(y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)) y^i \cdot \mathbf{x}^i + 2\lambda \cdot \mathbf{w} = \mathbf{0}$$

August 25, 2017

5

Gradient Descent



August 25, 2017

August 30, 2017



CS771: Intro to ML

Recap

First-order Optimality

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

$$\nabla f(\mathbf{w}) + \nabla r(\mathbf{w}) = \mathbf{0}$$

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

$$r(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{w} = (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$$

$$f(\mathbf{w}) = \sum_{i=1}^n \log(1 + \exp(-y^i \langle \mathbf{w}, \mathbf{x}^i \rangle))$$

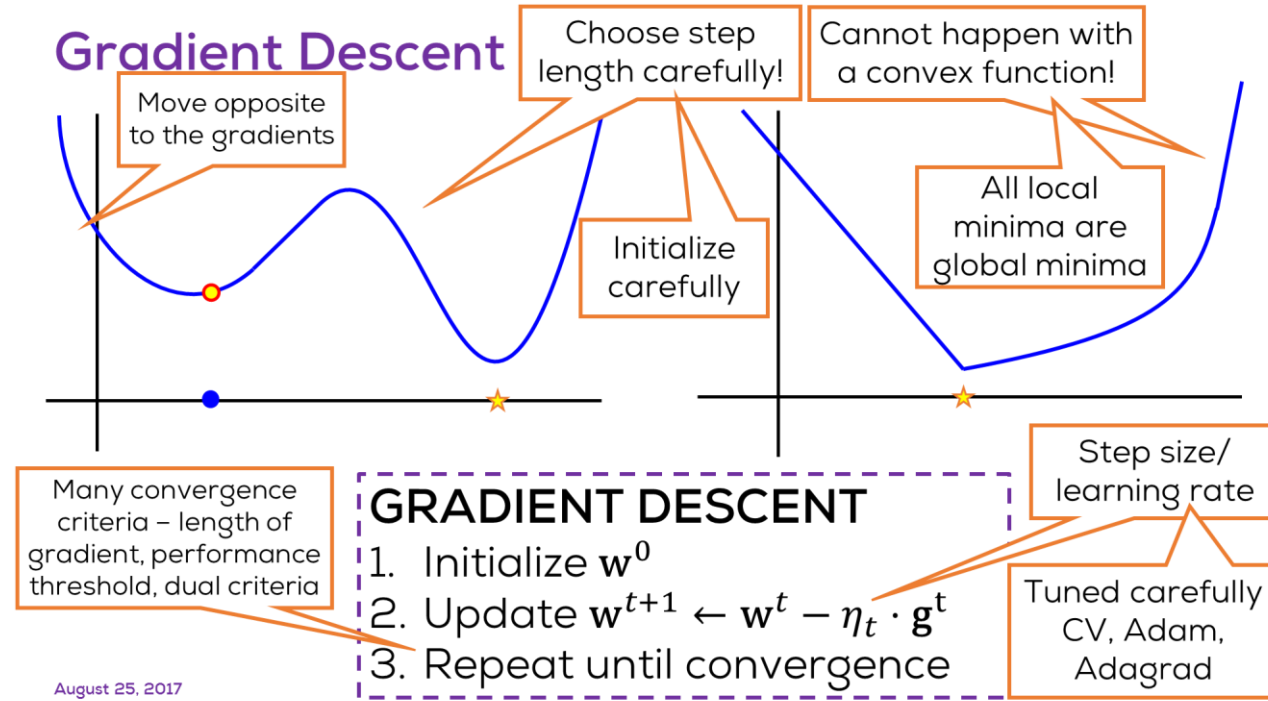
$$r(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\sum_{i=1}^n (1 - \sigma(y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)) y^i \cdot \mathbf{x}^i + 2\lambda \cdot \mathbf{w} = \mathbf{0}$$

August 25, 2017

5

Gradient Descent



August 25, 2017

GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
3. Repeat until convergence

What if function non-differentiable?

August 30, 2017



CS771: Intro to ML

Recap

First-order Optimality

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

$$\nabla f(\mathbf{w}) + \nabla r(\mathbf{w}) = \mathbf{0}$$

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

$$r(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{w} = (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$$

$$f(\mathbf{w}) = \sum_{i=1}^n \log(1 + \exp(-y^i \langle \mathbf{w}, \mathbf{x}^i \rangle))$$

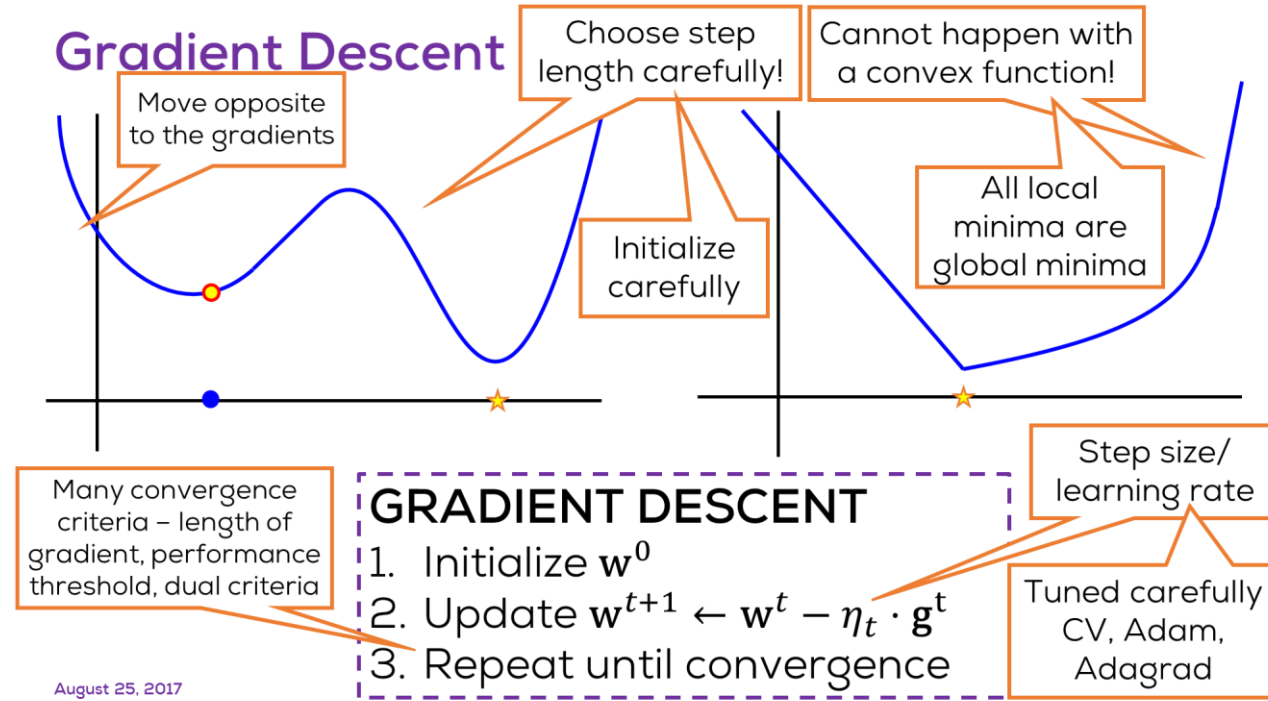
$$r(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\sum_{i=1}^n (1 - \sigma(y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)) y^i \cdot \mathbf{x}^i + 2\lambda \cdot \mathbf{w} = \mathbf{0}$$

August 25, 2017

5

Gradient Descent



August 25, 2017

What if function non-differentiable?

Use subgradients!

August 30, 2017



CS771: Intro to ML

Recap

First-order Optimality

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

$$\nabla f(\mathbf{w}) + \nabla r(\mathbf{w}) = \mathbf{0}$$

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

$$r(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{w} = (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$$

$$f(\mathbf{w}) = \sum_{i=1}^n \log(1 + \exp(-y^i \langle \mathbf{w}, \mathbf{x}^i \rangle))$$

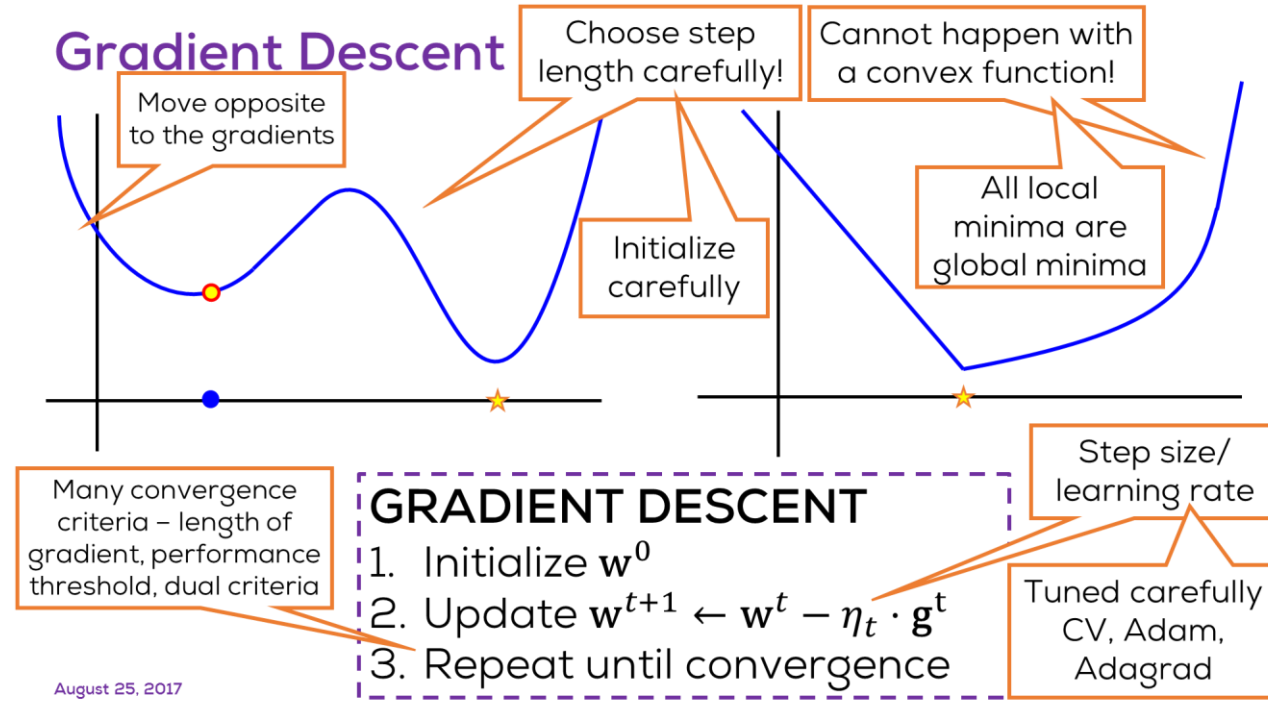
$$r(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\sum_{i=1}^n (1 - \sigma(y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)) y^i \cdot \mathbf{x}^i + 2\lambda \cdot \mathbf{w} = \mathbf{0}$$

August 25, 2017

5

Gradient Descent



August 25, 2017

What if function non-differentiable?

Use subgradients!

Defined for convex functions

August 30, 2017



CS771: Intro to ML

Recap

First-order Optimality

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

$$\nabla f(\mathbf{w}) + \nabla r(\mathbf{w}) = \mathbf{0}$$

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

$$r(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{w} = (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$$

$$f(\mathbf{w}) = \sum_{i=1}^n \log(1 + \exp(-y^i \langle \mathbf{w}, \mathbf{x}^i \rangle))$$

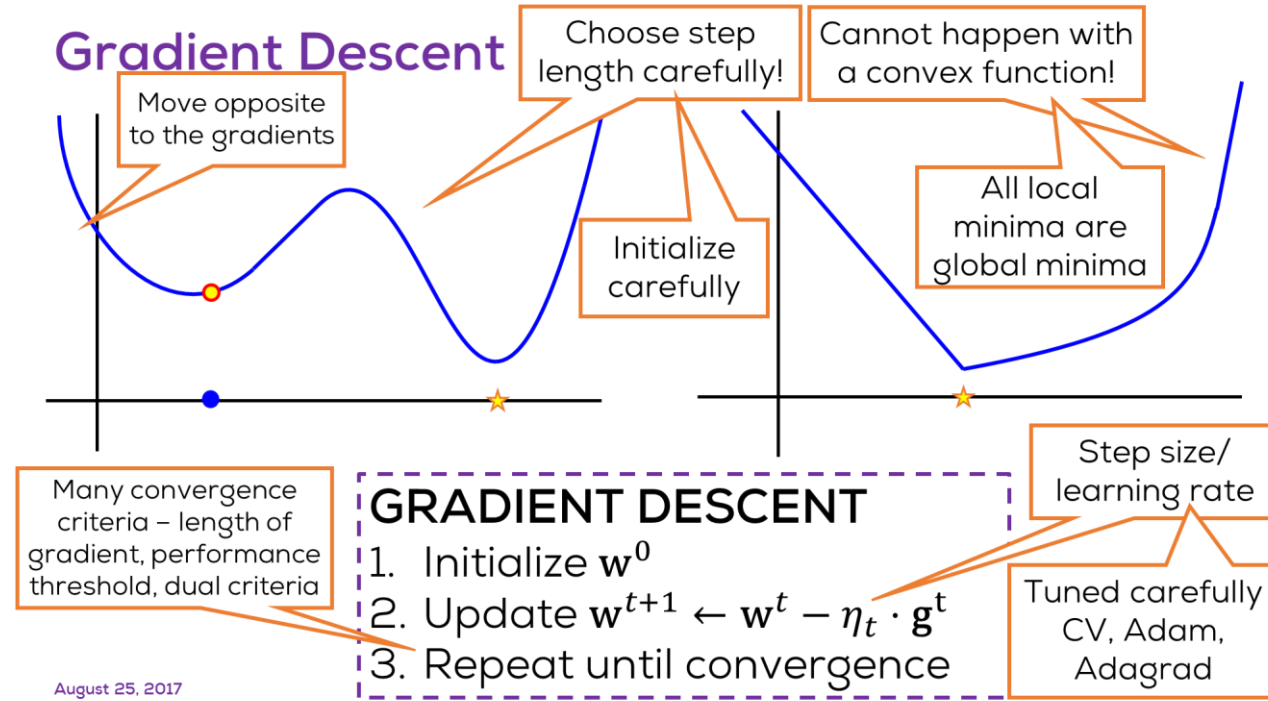
$$r(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\sum_{i=1}^n (1 - \sigma(y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)) y^i \cdot \mathbf{x}^i + 2\lambda \cdot \mathbf{w} = \mathbf{0}$$

August 25, 2017

5

Gradient Descent



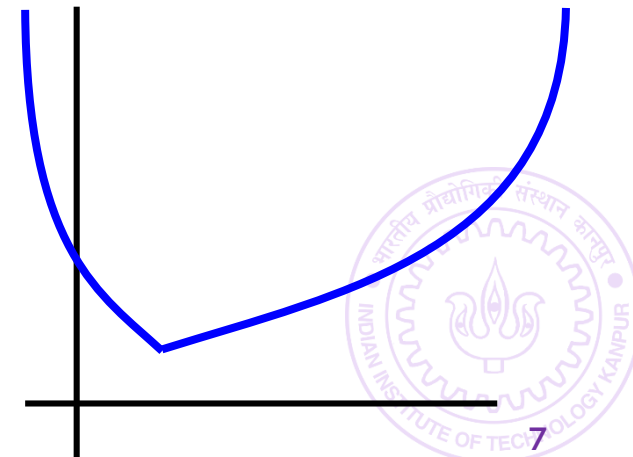
August 25, 2017

What if function non-differentiable?

Use subgradients!

Defined for convex functions

August 30, 2017



Recap

First-order Optimality

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

$$\nabla f(\mathbf{w}) + \nabla r(\mathbf{w}) = \mathbf{0}$$

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

$$r(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{w} = (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$$

$$f(\mathbf{w}) = \sum_{i=1}^n \log(1 + \exp(-y^i \langle \mathbf{w}, \mathbf{x}^i \rangle))$$

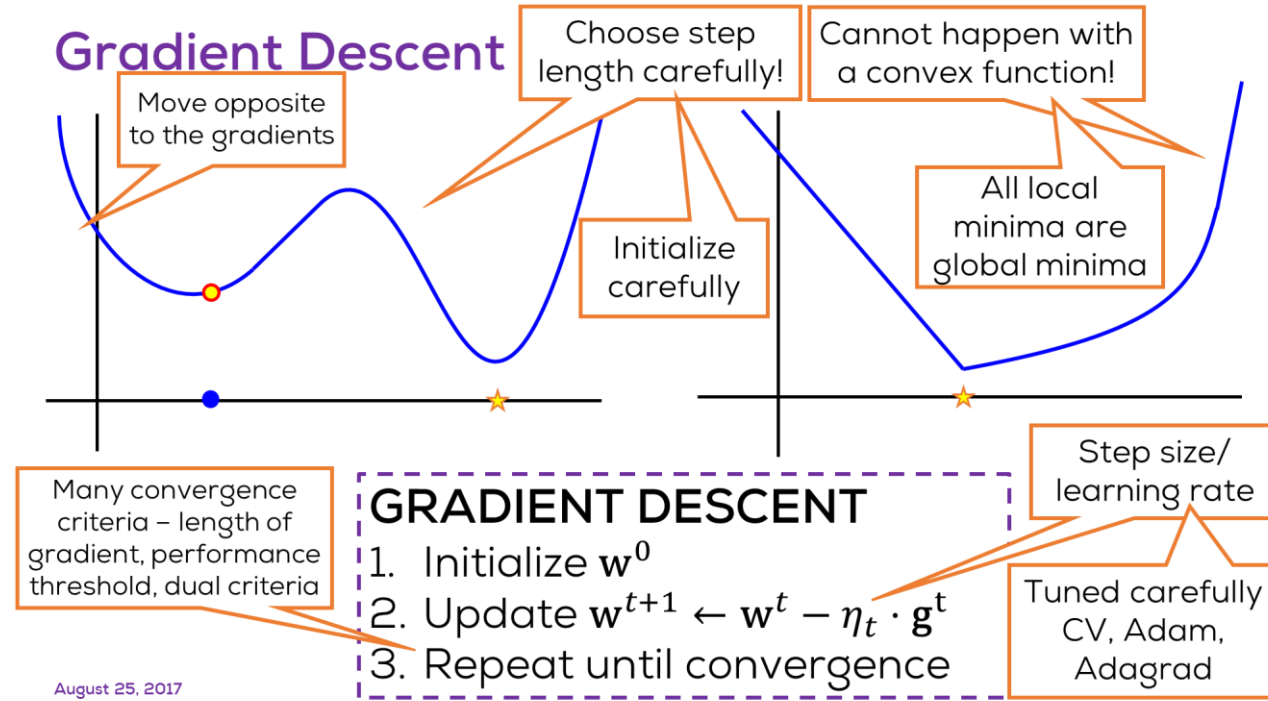
$$r(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\sum_{i=1}^n (1 - \sigma(y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)) y^i \cdot \mathbf{x}^i + 2\lambda \cdot \mathbf{w} = \mathbf{0}$$

August 25, 2017

5

Gradient Descent



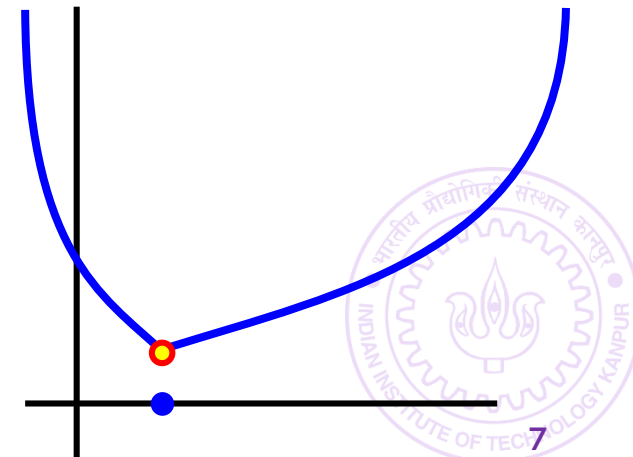
August 25, 2017

What if function non-differentiable?

Use subgradients!

Defined for convex functions

August 30, 2017



Recap

First-order Optimality

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

$$\nabla f(\mathbf{w}) + \nabla r(\mathbf{w}) = \mathbf{0}$$

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

$$r(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\mathbf{w} = (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$$

$$f(\mathbf{w}) = \sum_{i=1}^n \log(1 + \exp(-y^i \langle \mathbf{w}, \mathbf{x}^i \rangle))$$

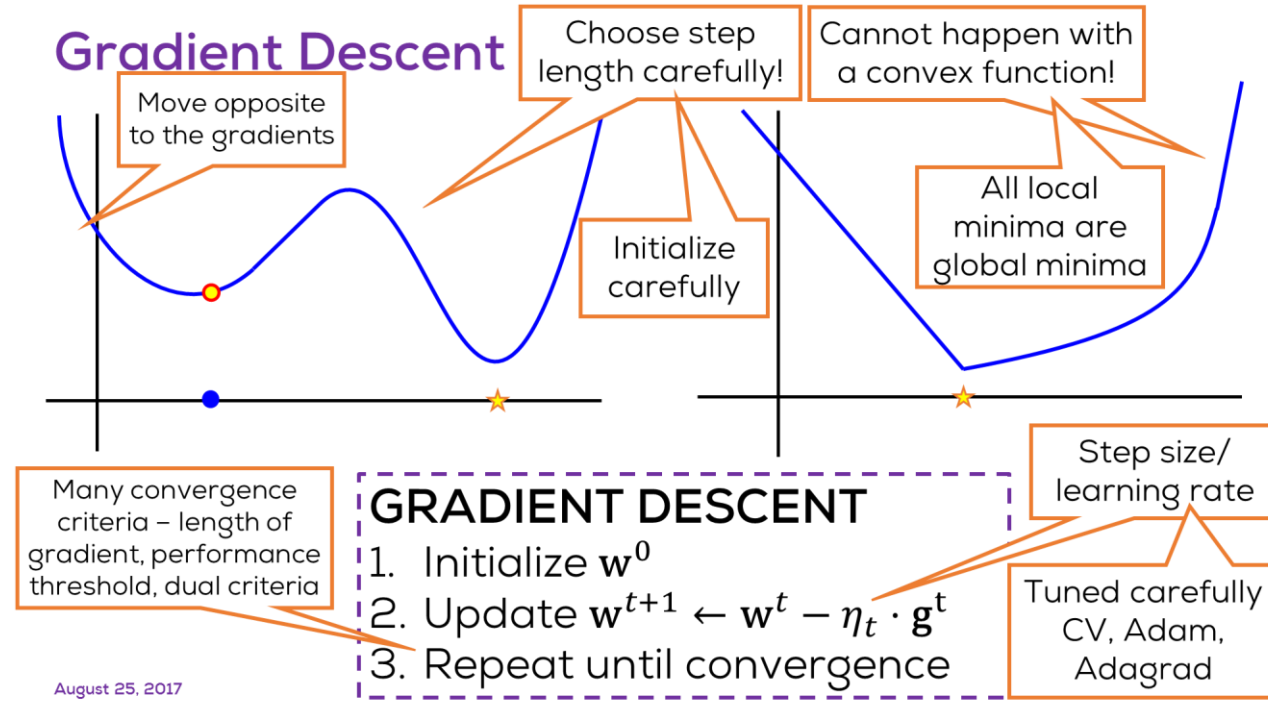
$$r(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$$

$$\sum_{i=1}^n (1 - \sigma(y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)) y^i \cdot \mathbf{x}^i + 2\lambda \cdot \mathbf{w} = \mathbf{0}$$

August 25, 2017

5

Gradient Descent



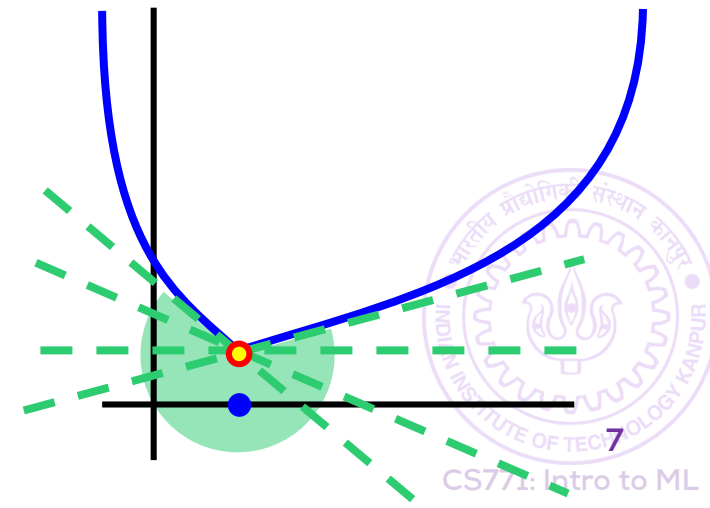
August 25, 2017

What if function non-differentiable?

Use subgradients!

Defined for convex functions

August 30, 2017



App: Linear Regression via GD

August 30, 2017



App: Linear Regression via GD

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

App: Linear Regression via

Convex?

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

App: Linear Regression via

Convex?

But but ...
 $(\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1}\mathbf{X}\mathbf{y}$

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



App: Linear Regression via

Convex?

But but ...
 $(\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$

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

Expensive!
 $O(d^3 + d^2n)$ time



App: Linear Regression via

Convex?

But but ...
 $(\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$

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

Expensive!
 $O(d^3 + d^2n)$ time

GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
3. Repeat until convergence



App: Linear Regression via

Convex?

But but ...
 $(\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$

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

Expensive!
 $O(d^3 + d^2n)$ time

GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
3. Repeat until convergence

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



App: Linear Regression via

Convex?

But but ...
 $(\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$

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

Expensive!
 $O(d^3 + d^2n)$ time

GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
3. Repeat until convergence

$$\begin{aligned} \mathbf{g}^t &= \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t) \\ &= 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}^i + 2\lambda \cdot \mathbf{w}^t \end{aligned}$$



App: Linear Regression via

Convex?

But but ...
 $(\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$

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

Expensive!
 $O(d^3 + d^2n)$ time

GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
3. Repeat until convergence

Only $O(nd)$
time per iter!

$$\begin{aligned} \mathbf{g}^t &= \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t) \\ &= 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}^i + 2\lambda \cdot \mathbf{w}^t \end{aligned}$$



App: Linear Regression via

Convex?

But but ...
 $(\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$

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

Expensive!
 $O(d^3 + d^2n)$ time

GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
3. Repeat until convergence

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

Only $O(nd)$
time per iter!

$$\begin{aligned} \mathbf{g}^t &= \nabla f(\mathbf{w}^t) + \nabla r(\mathbf{w}^t) \\ &= 2 \sum_{i=1}^n (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}^i + 2\lambda \cdot \mathbf{w}^t \end{aligned}$$



App: Linear Regression via

Convex?

But but ...
 $(\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$

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

Expensive!
 $O(d^3 + d^2n)$ time

GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
3. Repeat until convergence

Step length

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

Only $O(nd)$
time per iter!

$O\left(\frac{1}{\epsilon^2}\right)$ iterations
suffice to reach ϵ -
optimal solution

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

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



App: Linear Regression via

Convex?

But but ...
 $(\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$

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

Expensive!
 $O(d^3 + d^2n)$ time

GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
3. Repeat until convergence

Step length

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

Only $O(nd)$
time per iter!

$O\left(\frac{1}{\epsilon^2}\right)$ iterations
suffice to reach ϵ -
optimal solution

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

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

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



App: Linear Regression via

Convex?

But but ...
 $(\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$

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

Expensive!
 $O(d^3 + d^2n)$ time

GRADIENT DESCENT

1. Initialize \mathbf{w}^0
2. Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_t \cdot \mathbf{g}^t$
3. Repeat until convergence

Step length

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

Only $O(nd)$
time per iter!

$O\left(\frac{1}{\epsilon^2}\right)$ iterations
suffice to reach ϵ -
optimal solution

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

For convex
problems

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

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



A Closer look at the GD Update

August 30, 2017



A Closer look at the GD Update

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

A Closer look at the GD Update

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

A Closer look at the GD Update

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

Does not let model
change too much!

A Closer look at the GD Update

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

Does not let model
change too much!

What if $\lambda = 0$?
What if $\eta_t = 0$?
What if no $f(\cdot)$?

A Closer look at the GD Update

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

Does not let model
change too much!

What if $\lambda = 0$?
What if $\eta_t = 0$?
What if no $f(\cdot)$?

A Closer look at the GD Update

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

Does not let model
change too much!

What if $\lambda = 0$?
What if $\eta_t = 0$?
What if no $f(\cdot)$?

If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle > y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle = y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle < y^i$

A Closer look at the GD Update

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

Does not let model
change too much!

What if $\lambda = 0$?
What if $\eta_t = 0$?
What if no $f(\cdot)$?

Push \mathbf{w}^t "away"
from \mathbf{x}^i

If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle > y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle = y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle < y^i$

A Closer look at the GD Update

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

Does not let model
change too much!

What if $\lambda = 0$?
What if $\eta_t = 0$?
What if no $f(\cdot)$?

Push \mathbf{w}^t "away"
from \mathbf{x}^i

No need to
care about \mathbf{x}^i

If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle > y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle = y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle < y^i$

A Closer look at the GD Update

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

Does not let model
change too much!

What if $\lambda = 0$?
What if $\eta_t = 0$?
What if no $f(\cdot)$?

Push \mathbf{w}^t "away"
from \mathbf{x}^i

No need to
care about \mathbf{x}^i

Push \mathbf{w}^t
"towards" \mathbf{x}^i

If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle > y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle = y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle < y^i$

A Closer look at the GD Update

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

Does not let model
change too much!

What if $\lambda = 0$?
What if $\eta_t = 0$?
What if no $f(\cdot)$?

Push \mathbf{w}^t "away"
from \mathbf{x}^i

No need to
care about \mathbf{x}^i

Push \mathbf{w}^t
"towards" \mathbf{x}^i

If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle > y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle = y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle < y^i$

What if $\eta_t = 0$?
What if no $r(\cdot)$?

A Closer look at the GD Update

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

Does not let model
change too much!

What if $\lambda = 0$?
What if $\eta_t = 0$?
What if no $f(\cdot)$?

Push \mathbf{w}^t "away"
from \mathbf{x}^i

No need to
care about \mathbf{x}^i

Push \mathbf{w}^t
"towards" \mathbf{x}^i

If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle > y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle = y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle < y^i$

What if $\eta_t = 0$?
What if no $r(\cdot)$?

Exercise: for clarity, see what happens when $d = 1$

A Closer look at the GD Update

Lets see the case
 $n = 1, \lambda = 0$

$$\mathbf{w}^{t+1} =$$

$$\mathbf{w}^t - 2\eta_t \cdot (\langle \mathbf{w}^t, \mathbf{x}^i \rangle - y^i) \cdot \mathbf{x}^i$$

Does not let model
change too much!

What if $\lambda = 0$?
What if $\eta_t = 0$?
What if no $f(\cdot)$?

Push \mathbf{w}^t "away"
from \mathbf{x}^i

No need to
care about \mathbf{x}^i

Push \mathbf{w}^t
"towards" \mathbf{x}^i

If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle > y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle = y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle < y^i$

What if $\eta_t = 0$?
What if no $r(\cdot)$?

A Closer look at the GD Update

Lets see the case
 $n = 1, \lambda = 0$

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

Does not let model
change too much!

$$\langle \mathbf{w}^{t+1}, \mathbf{x}^i \rangle < \langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

Push \mathbf{w}^t "away"
from \mathbf{x}^i

If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle > y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle = y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle < y^i$

What if $\lambda = 0$?
What if $\eta_t = 0$?
What if no $f(\cdot)$?

No need to
care about \mathbf{x}^i

Push \mathbf{w}^t
"towards" \mathbf{x}^i

What if $\eta_t = 0$?
What if no $r(\cdot)$?

A Closer look at the GD Update

Lets see the case
 $n = 1, \lambda = 0$

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

Does not let model
change too much!

$$\langle \mathbf{w}^{t+1}, \mathbf{x}^i \rangle < \langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

Push \mathbf{w}^t "away"
from \mathbf{x}^i

If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle > y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle = y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle < y^i$

What if $\lambda = 0$?
What if $\eta_t = 0$?
What if no $f(\cdot)$?

$$\langle \mathbf{w}^{t+1}, \mathbf{x}^i \rangle = \langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

No need to
care about \mathbf{x}^i

Push \mathbf{w}^t
"towards" \mathbf{x}^i

What if $\eta_t = 0$?
What if no $r(\cdot)$?

A Closer look at the GD Update

Lets see the case
 $n = 1, \lambda = 0$

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

Does not let model
change too much!

$$\langle \mathbf{w}^{t+1}, \mathbf{x}^i \rangle < \langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

Push \mathbf{w}^t "away"
from \mathbf{x}^i

If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle > y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle = y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle < y^i$

What if $\lambda = 0$?
What if $\eta_t = 0$?
What if no $f(\cdot)$?

$$\langle \mathbf{w}^{t+1}, \mathbf{x}^i \rangle = \langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

No need to
care about \mathbf{x}^i

$$\langle \mathbf{w}^{t+1}, \mathbf{x}^i \rangle > \langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

Push \mathbf{w}^t
"towards" \mathbf{x}^i

What if $\eta_t = 0$?
What if no $r(\cdot)$?

A Closer look at the GD Update

Lets see the case
 $n = 1, \lambda = 0$

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

Does not let model
change too much!

$$\langle \mathbf{w}^{t+1}, \mathbf{x}^i \rangle < \langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

Push \mathbf{w}^t "away"
from \mathbf{x}^i

If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle > y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle = y^i$
If $\langle \mathbf{w}^t, \mathbf{x}^i \rangle < y^i$

What if $\lambda = 0$?
What if $\eta_t = 0$?
What if no $f(\cdot)$?

$$\langle \mathbf{w}^{t+1}, \mathbf{x}^i \rangle = \langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

No need to
care about \mathbf{x}^i

$$\langle \mathbf{w}^{t+1}, \mathbf{x}^i \rangle > \langle \mathbf{w}^t, \mathbf{x}^i \rangle$$

Push \mathbf{w}^t
"towards" \mathbf{x}^i

What if $\eta_t = 0$?
What if no $r(\cdot)$?

We are making corrective updates!

Subgradient Calculus

Do not always
hold ☹

Should hold in
nice settings ☺

MTH101 Calculus

Subgradient Calculus

$$\mathbf{x} \in \mathbb{R}^d, \mathbf{A} \in \mathbb{R}^{k \times d}, \mathbf{b} \in \mathbb{R}^k, c \in \mathbb{R}$$

Scaling	<ul style="list-style-type: none">• $\nabla(c \cdot f)(\mathbf{x}) = c \cdot \nabla f(\mathbf{x})$	<ul style="list-style-type: none">• $\partial(c \cdot f)(\mathbf{x}) = c \cdot \partial f(\mathbf{x})$ $= \{c \cdot \mathbf{v} : \mathbf{v} \in \partial f(\mathbf{x})\}$
Sum Rule	<ul style="list-style-type: none">• $\nabla(f + g)(x) = \nabla f(x) + \nabla g(x)$	<ul style="list-style-type: none">• $\partial(f + g)(\mathbf{x}) = \partial f(\mathbf{x}) + \partial g(\mathbf{x})$ $= \{\mathbf{u} + \mathbf{v} : \mathbf{u} \in \partial f(\mathbf{x}), \mathbf{v} \in \partial g(\mathbf{x})\}$
Chain Rule	<ul style="list-style-type: none">• $\nabla f(\mathbf{Ax} + \mathbf{b}) = \mathbf{A}^\top \nabla f(\mathbf{Ax} + \mathbf{b})$	<ul style="list-style-type: none">• $\partial f(\mathbf{Ax} + \mathbf{b}) = \mathbf{A}^\top \partial f(\mathbf{Ax} + \mathbf{b})$ $= \{\mathbf{A}^\top \mathbf{v} : \mathbf{v} \in \partial f(\mathbf{x})\}$
Max Rule	$\ell_{cs}(y, \{\eta_j\}) = [1 + \max_{k \neq y} \eta_k - \eta_y]_+$	$\circ \circ \circ \quad f(\mathbf{x}) = \max_i f_i(\mathbf{x})$ $\partial f(\mathbf{x}) = \{\sum \alpha_i \mathbf{v}^i : \mathbf{v}^i \in \partial f_i(\mathbf{x}), \sum \alpha_i = 1, \alpha_i \geq 0\}$

Subgradient Calculus

Do not always
hold ☹

Should hold in
nice settings ☺

MTH101 Calculus

Subgradient Calculus

$$\mathbf{x} \in \mathbb{R}^d, \mathbf{A} \in \mathbb{R}^{k \times d}, \mathbf{b} \in \mathbb{R}^k, c \in \mathbb{R}$$

Scaling

$$\bullet \nabla(c \cdot f)(\mathbf{x}) = c \cdot \nabla f(\mathbf{x})$$

$$\bullet \partial(c \cdot f)(\mathbf{x}) = c \cdot \partial f(\mathbf{x}) \\ = \{c \cdot \mathbf{v} : \mathbf{v} \in \partial f(\mathbf{x})\}$$

Sum Rule

$$\bullet \nabla(f + g)(x) = \nabla f(x) + \nabla g(x)$$

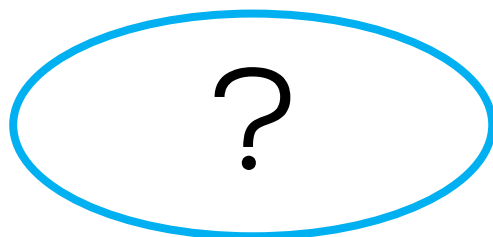
$$\bullet \partial(f + g)(\mathbf{x}) = \partial f(\mathbf{x}) + \partial g(\mathbf{x}) \\ = \{\mathbf{u} + \mathbf{v} : \mathbf{u} \in \partial f(\mathbf{x}), \mathbf{v} \in \partial g(\mathbf{x})\}$$

Chain Rule

$$\bullet \nabla f(\mathbf{Ax} + \mathbf{b}) = \mathbf{A}^\top \nabla f(\mathbf{Ax} + \mathbf{b})$$

$$\bullet \partial f(\mathbf{Ax} + \mathbf{b}) = \mathbf{A}^\top \partial f(\mathbf{Ax} + \mathbf{b}) \\ = \{\mathbf{A}^\top \mathbf{v} : \mathbf{v} \in \partial f(\mathbf{x})\}$$

Max Rule



$$f(\mathbf{x}) = \max_i f_i(\mathbf{x})$$

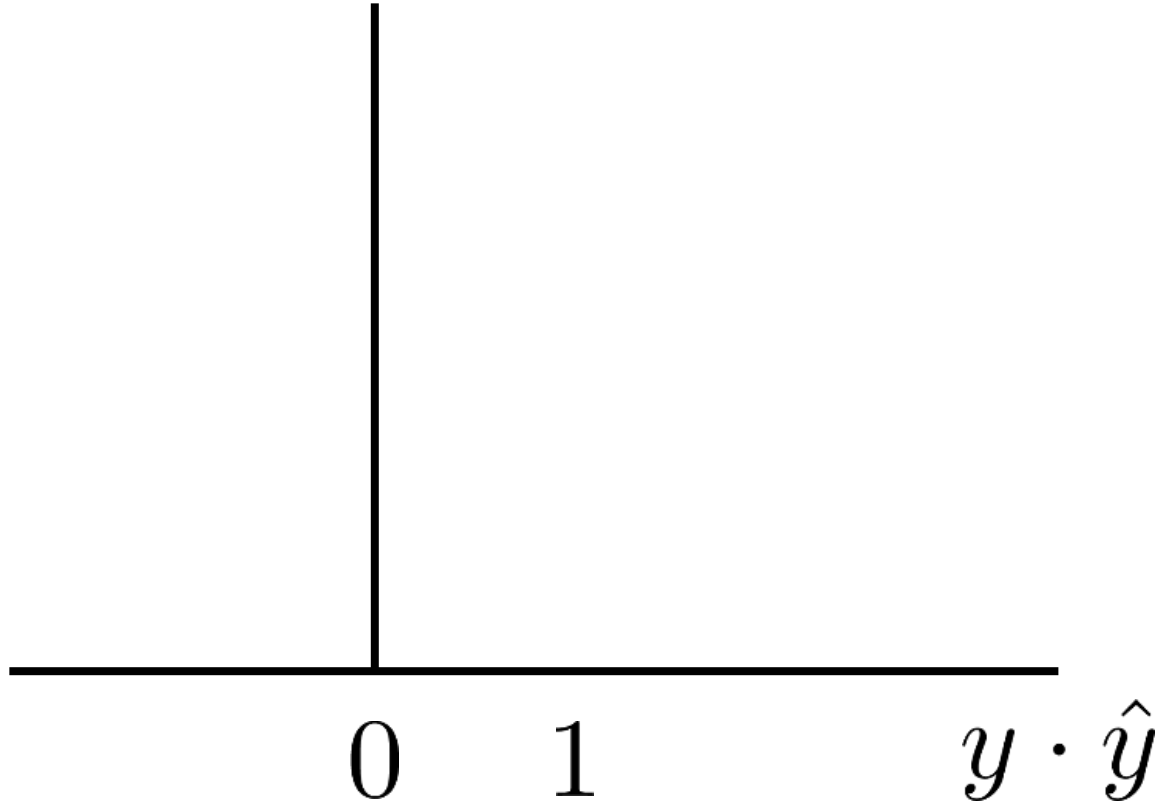
$$\partial f(\mathbf{x}) = \{\sum \alpha_i \mathbf{v}^i : \mathbf{v}^i \in \partial f_i(\mathbf{x}), \sum \alpha_i = 1, \alpha_i \geq 0\}$$

Application: SVM via subGD

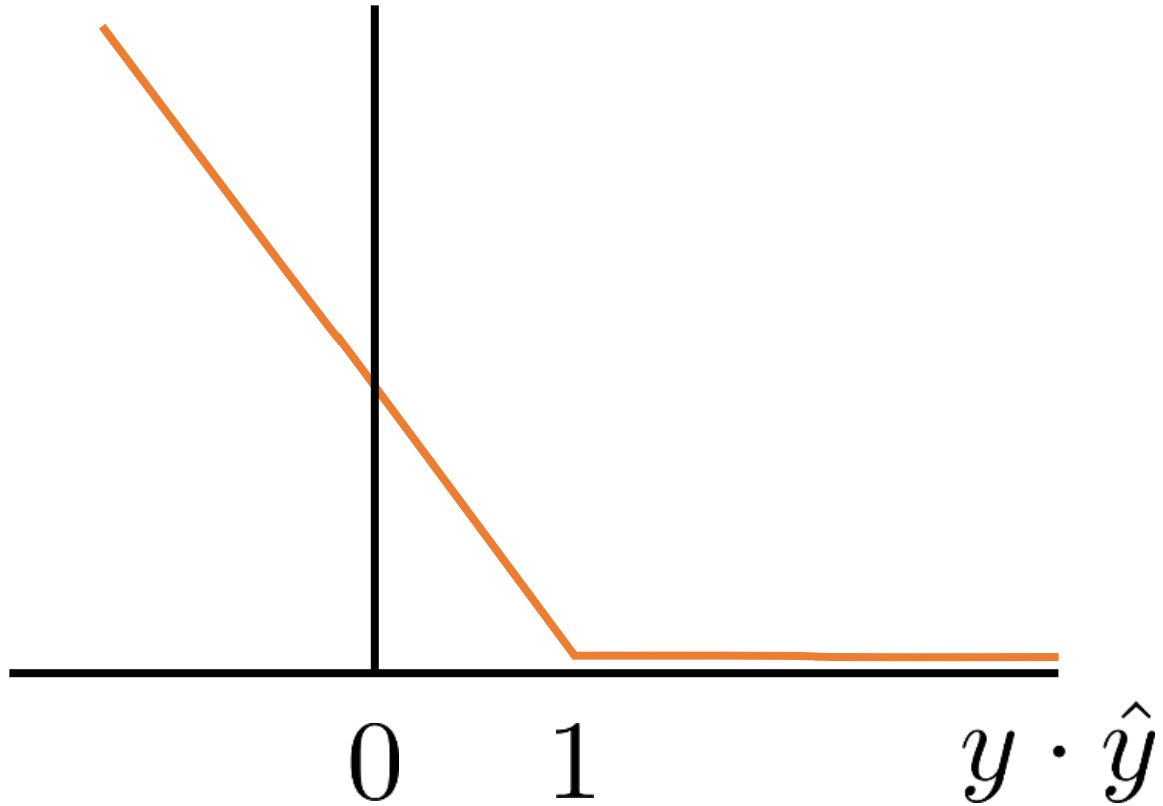
August 30, 2017



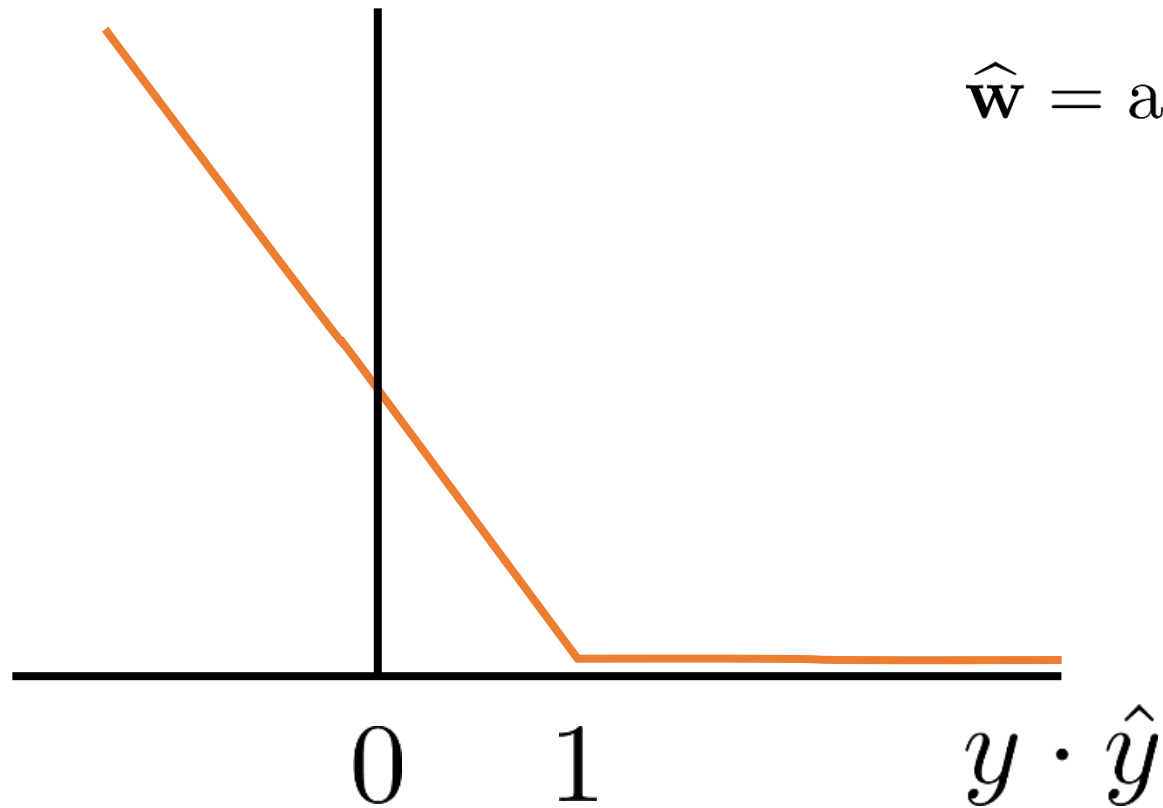
Application: SVM via subGD



Application: SVM via subGD

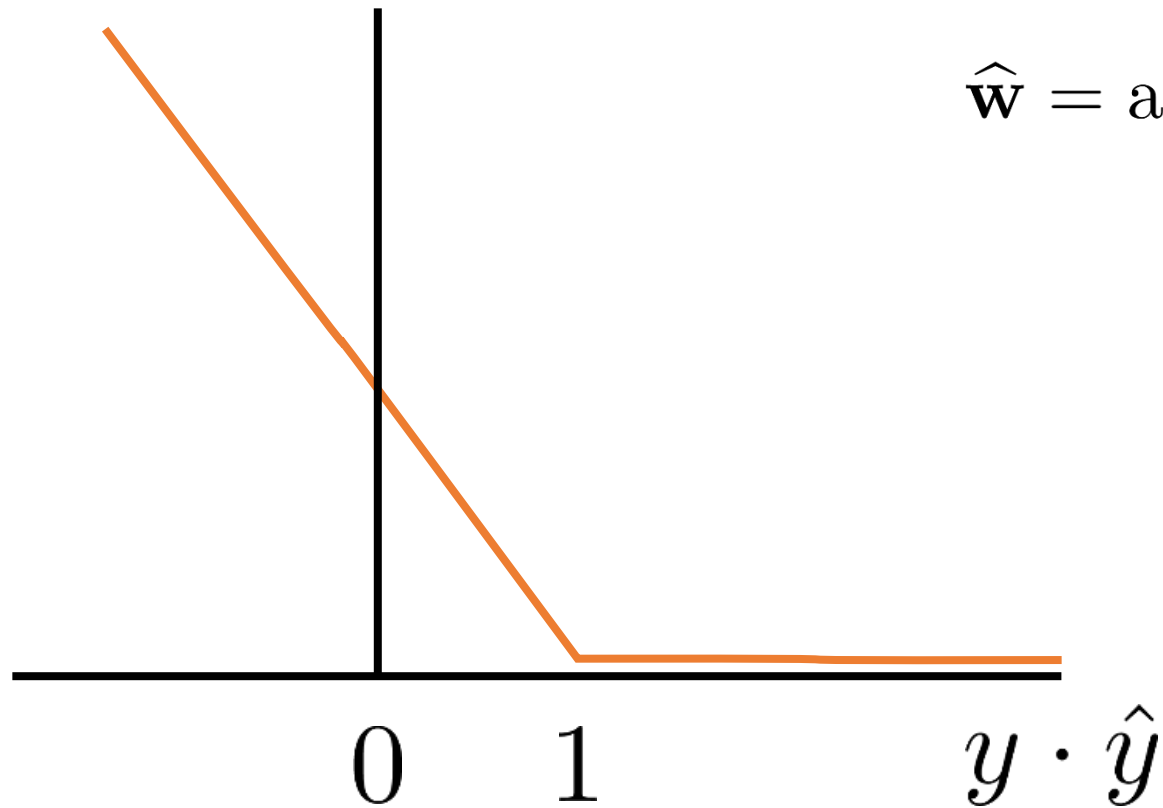


Application: SVM via subGD



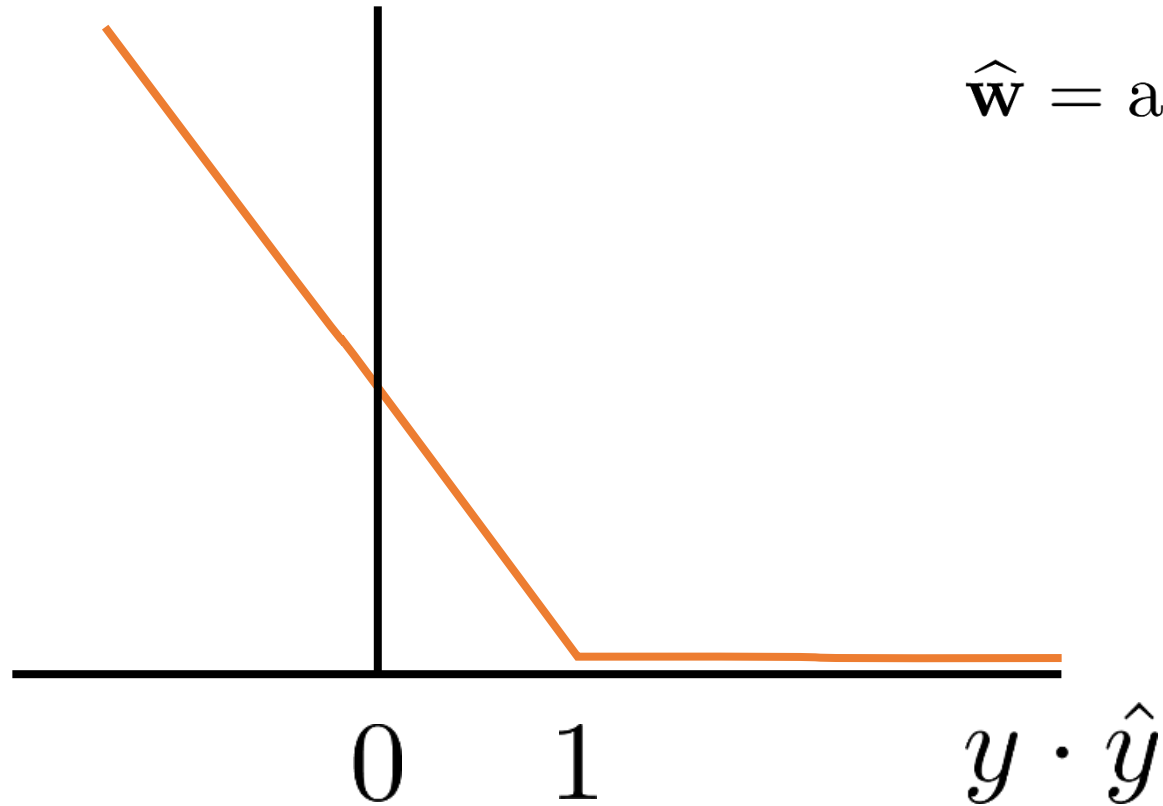
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{i=1}^n \ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Application: SVM via subGD



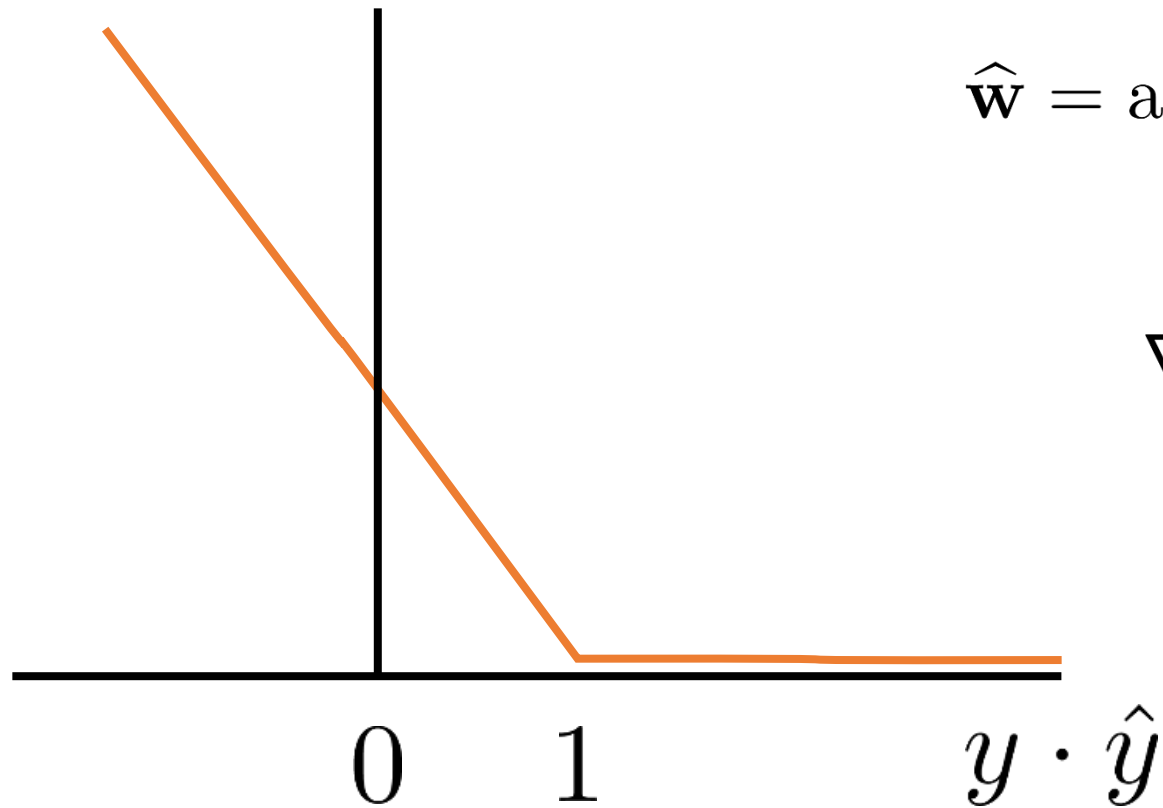
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \underbrace{\sum_{i=1}^n \ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)}_{f(\mathbf{w})} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Application: SVM via subGD



$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \underbrace{\sum_{i=1}^n \ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)}_{f(\mathbf{w})} + \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|_2^2}_{r(\mathbf{w})}$$

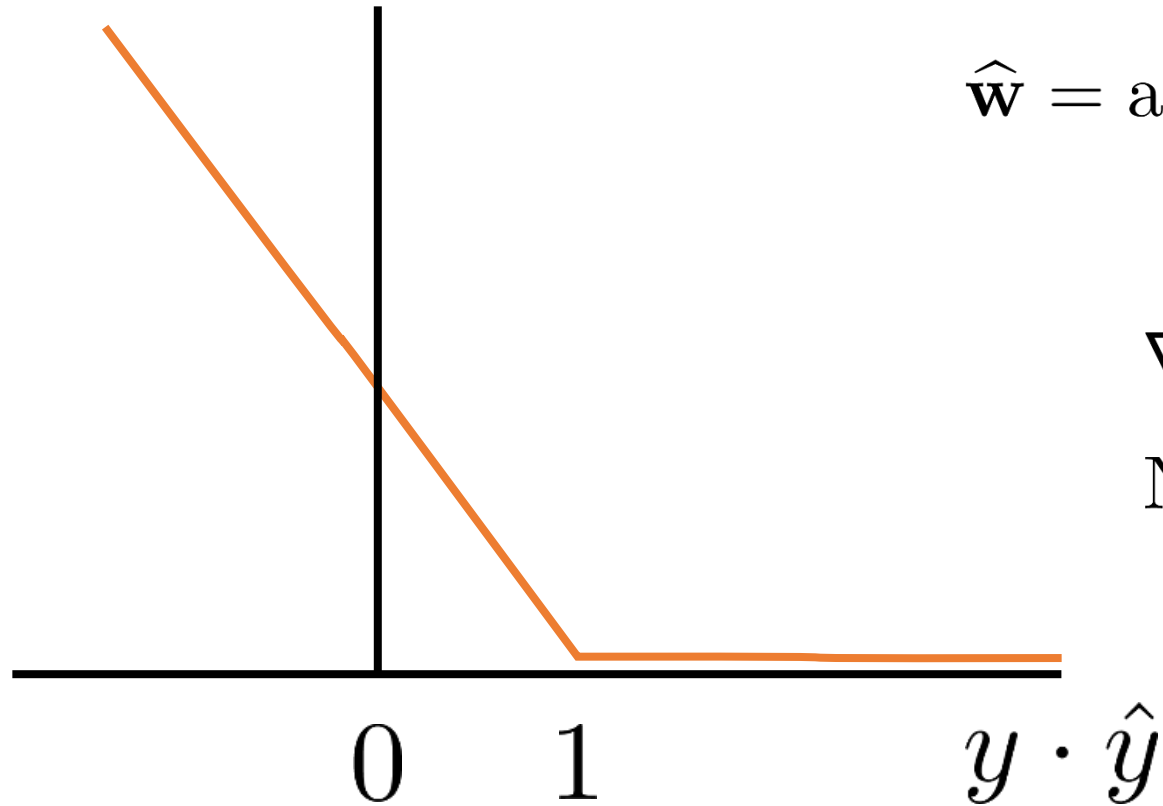
Application: SVM via subGD



$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \underbrace{\sum_{i=1}^n \ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)}_{f(\mathbf{w})} + \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|_2^2}_{r(\mathbf{w})}$$

$$\nabla r(\mathbf{w}) = \lambda \cdot \mathbf{w}$$

Application: SVM via subGD

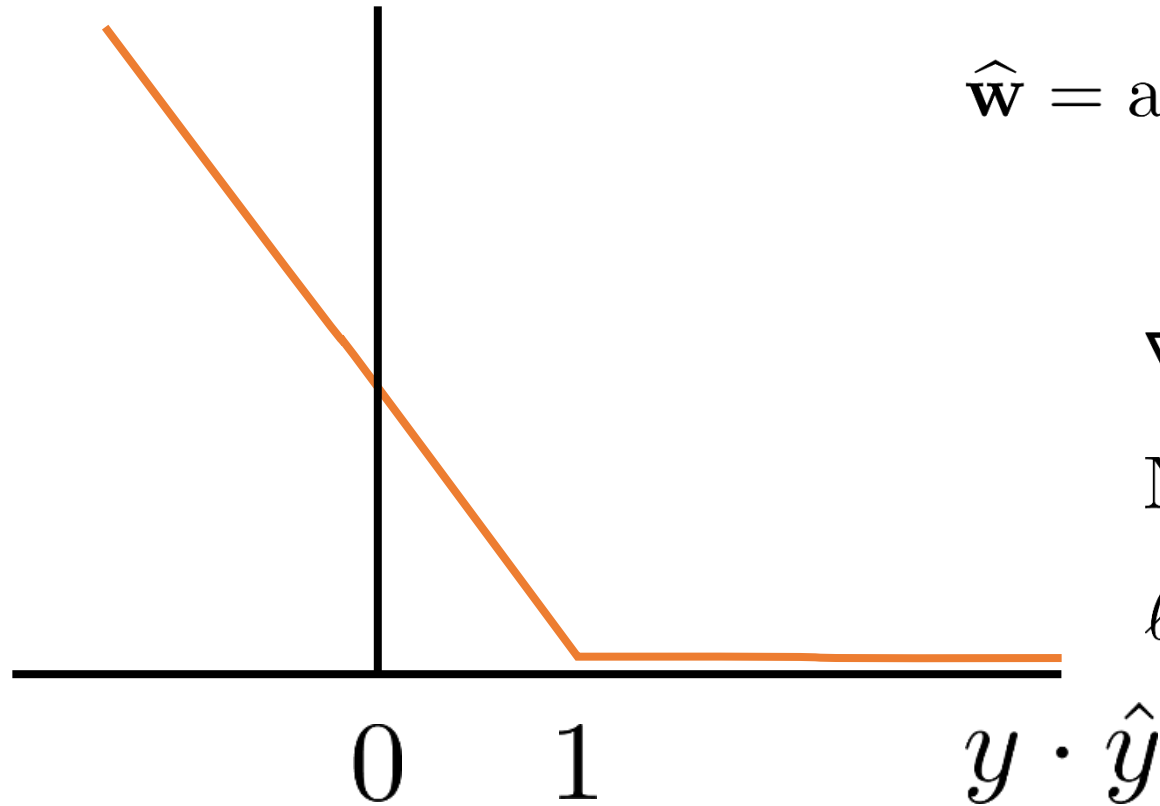


$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \underbrace{\sum_{i=1}^n \ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)}_{f(\mathbf{w})} + \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|_2^2}_{r(\mathbf{w})}$$

$$\nabla r(\mathbf{w}) = \lambda \cdot \mathbf{w}$$

Need $\mathbf{v}^i \in \partial \ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$

Application: SVM via subGD



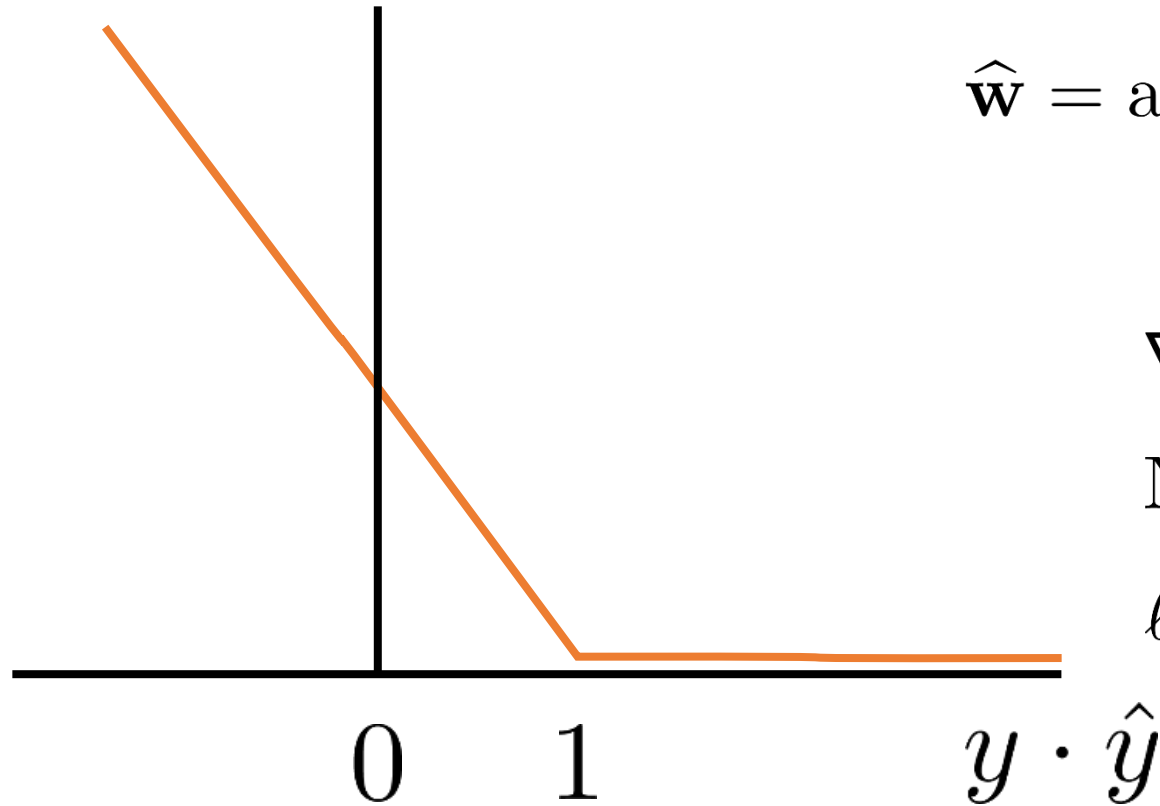
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \underbrace{\sum_{i=1}^n \ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)}_{f(\mathbf{w})} + \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|_2^2}_{r(\mathbf{w})}$$

$$\nabla r(\mathbf{w}) = \lambda \cdot \mathbf{w}$$

Need $\mathbf{v}^i \in \partial \ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$

$$\ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle) = [1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle]_+$$

Application: SVM via subGD



$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \underbrace{\sum_{i=1}^n \ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)}_{f(\mathbf{w})} + \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|_2^2}_{r(\mathbf{w})}$$

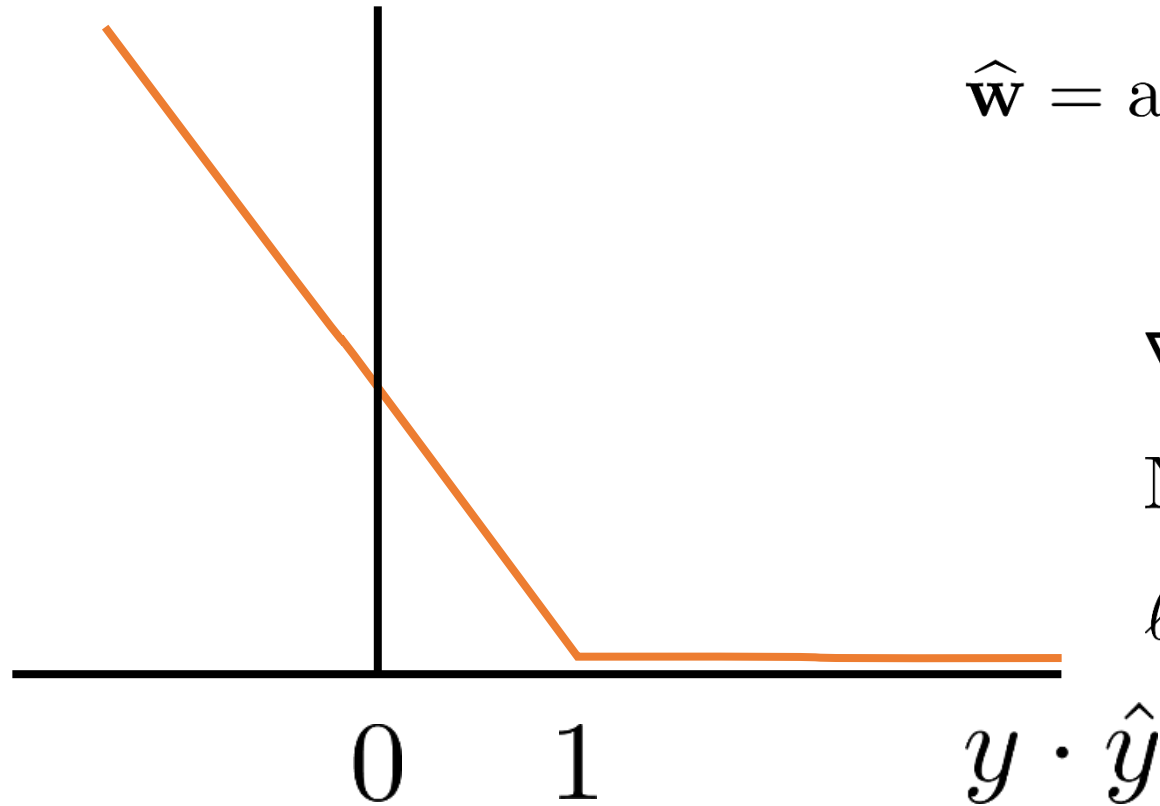
$$\nabla r(\mathbf{w}) = \lambda \cdot \mathbf{w}$$

Need $\mathbf{v}^i \in \partial \ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$

$$\ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle) = [1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle]_+$$

$$\mathbf{v}^i = \begin{cases} \mathbf{0} & \text{if } y^i \langle \mathbf{w}, \mathbf{x}^i \rangle > 1 \\ -y^i \cdot \mathbf{x}^i & \text{if } y^i \langle \mathbf{w}, \mathbf{x}^i \rangle < 1 \\ c \cdot y^i \cdot \mathbf{x}^i & \text{if } y^i \langle \mathbf{w}, \mathbf{x}^i \rangle = 1 \\ c \in [-1, 0] \end{cases}$$

Application: SVM via subGD



$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \underbrace{\sum_{i=1}^n \ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)}_{f(\mathbf{w})} + \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|_2^2}_{r(\mathbf{w})}$$

$$\nabla r(\mathbf{w}) = \lambda \cdot \mathbf{w}$$

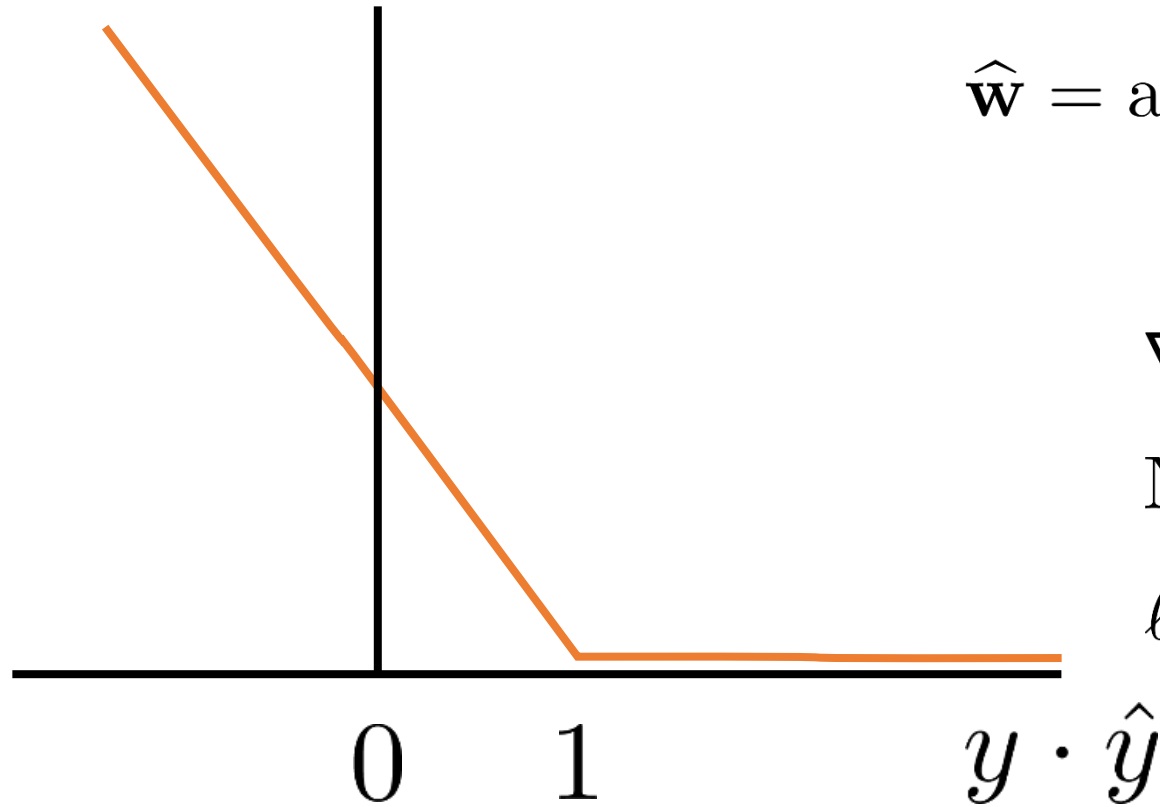
Need $\mathbf{v}^i \in \partial \ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$

$$\ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle) = [1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle]_+$$

$O(d)$ time per
data point!

$$\mathbf{v}^i = \begin{cases} \mathbf{0} & \text{if } y^i \langle \mathbf{w}, \mathbf{x}^i \rangle > 1 \\ -y^i \cdot \mathbf{x}^i & \text{if } y^i \langle \mathbf{w}, \mathbf{x}^i \rangle < 1 \\ c \cdot y^i \cdot \mathbf{x}^i & \text{if } y^i \langle \mathbf{w}, \mathbf{x}^i \rangle = 1 \\ c \in [-1, 0] \end{cases}$$

Application: SVM via subGD



$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \underbrace{\sum_{i=1}^n \ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)}_{f(\mathbf{w})} + \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|_2^2}_{r(\mathbf{w})}$$

$$\nabla r(\mathbf{w}) = \lambda \cdot \mathbf{w}$$

Need $\mathbf{v}^i \in \partial \ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle)$

$$\ell_{\text{hinge}}(y^i, \langle \mathbf{w}, \mathbf{x}^i \rangle) = [1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle]_+$$

$$\mathbf{v}^i = \begin{cases} \mathbf{0} & \text{if } y^i \langle \mathbf{w}, \mathbf{x}^i \rangle > 1 \\ -y^i \cdot \mathbf{x}^i & \text{if } y^i \langle \mathbf{w}, \mathbf{x}^i \rangle < 1 \\ c \cdot y^i \cdot \mathbf{x}^i & \text{if } y^i \langle \mathbf{w}, \mathbf{x}^i \rangle = 1 \\ c \in [-1, 0] \end{cases}$$

$O(nd)$ time
per iter

$O(d)$ time per
data point!

App: Sparse Regression via subGD

August 30, 2017



App: Sparse Regression via subGD

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \lambda \cdot \|\mathbf{w}\|_1$$

App: Sparse Regression via ℓ_1 Convex?

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_1$$

App: Sparse Regression via sul

Convex?

Non-
differentiable

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_1$$

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

App: Sparse Regression via sul

Convex?

Non-
differentiable

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_1$$

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

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

App: Sparse Regression via sub

Convex?

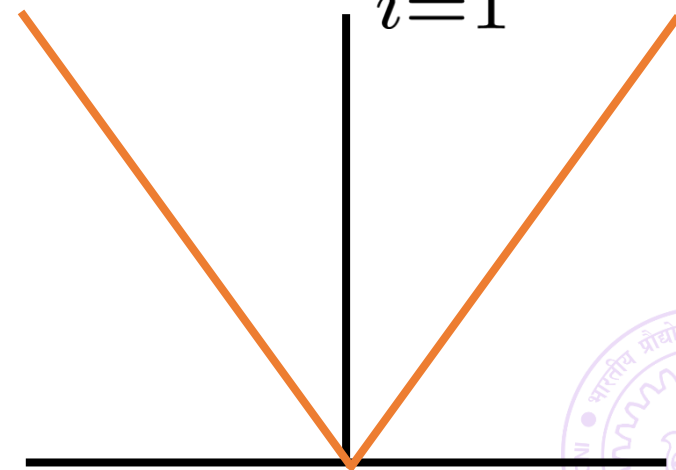
Non-differentiable

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_1$$

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

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

$$\|\mathbf{w}^t\|_1 = \sum_{i=1}^d |\mathbf{w}_i^t|$$



App: Sparse Regression via sub

Convex?

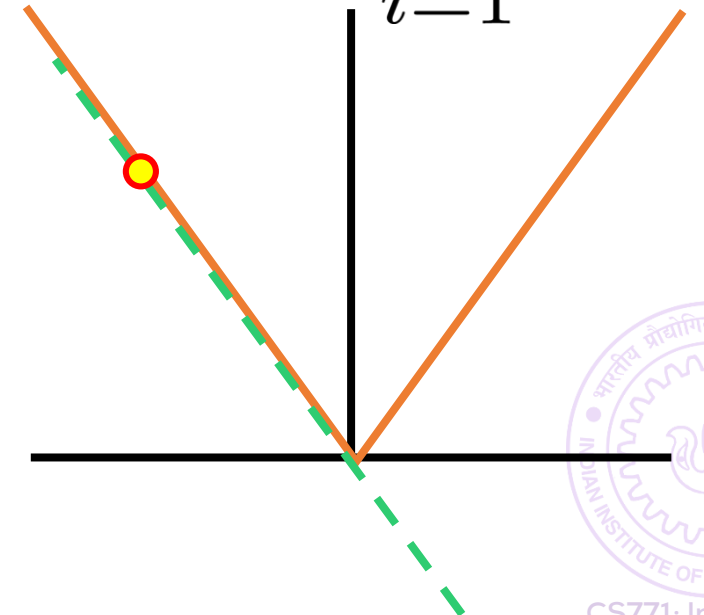
Non-differentiable

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_1$$

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

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

$$\|\mathbf{w}^t\|_1 = \sum_{i=1}^d |\mathbf{w}_i^t|$$



App: Sparse Regression via sul

Convex?

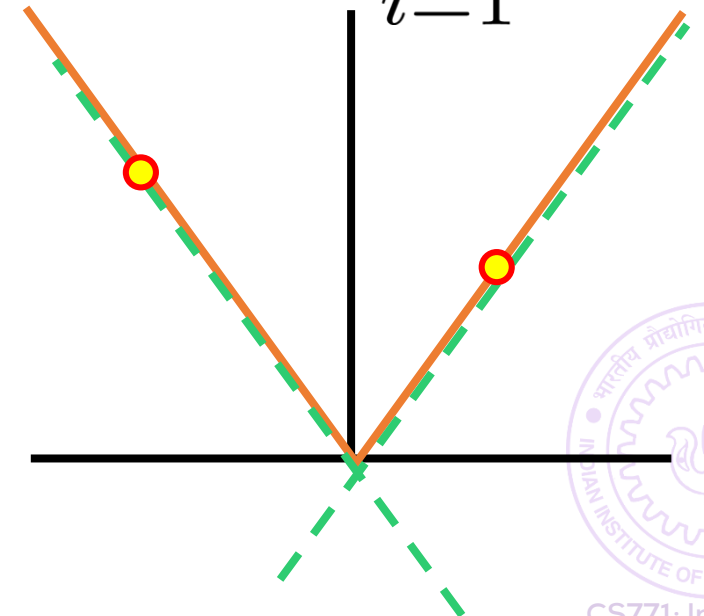
Non-differentiable

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_1$$

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

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

$$\|\mathbf{w}^t\|_1 = \sum_{i=1}^d |\mathbf{w}_i^t|$$



App: Sparse Regression via sub

Convex?

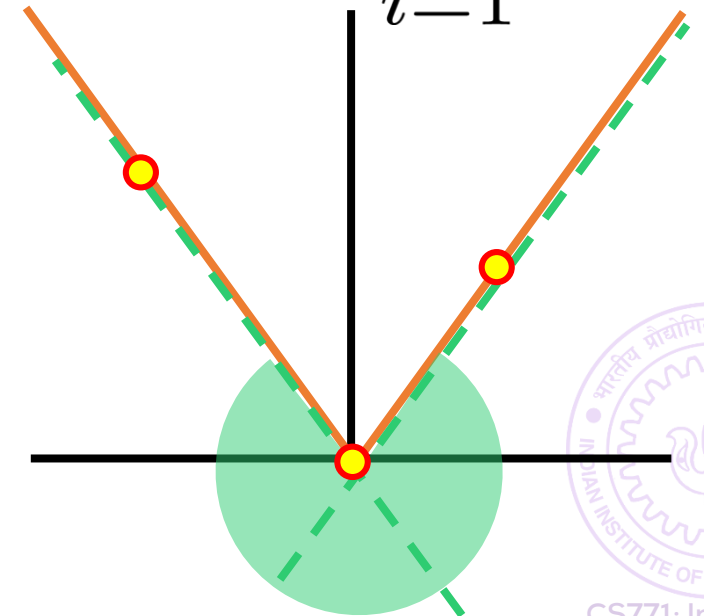
Non-differentiable

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_1$$

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

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

$$\|\mathbf{w}^t\|_1 = \sum_{i=1}^d |\mathbf{w}_i^t|$$



App: Sparse Regression via sub

Convex?

Non-differentiable

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_1$$

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

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

$$\|\mathbf{w}^t\|_1 = \sum_{i=1}^d |\mathbf{w}_i^t|$$

$$\rho_i^t = \begin{cases} +1 & \text{if } \mathbf{w}_i^t > 0 \\ [-1, 1] & \text{if } \mathbf{w}_i^t = 0 \\ -1 & \text{if } \mathbf{w}_i^t < 0 \end{cases}$$

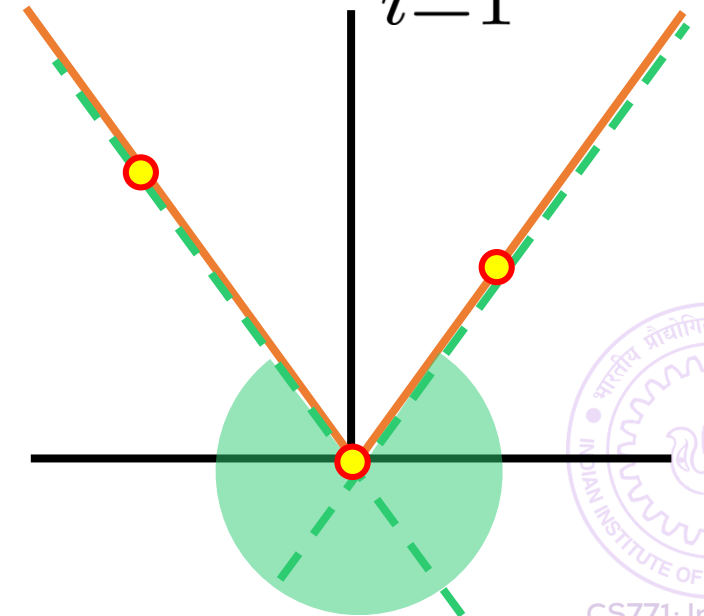
$\partial |\mathbf{w}_i|$

+1

-1

\mathbf{w}_i

e.g. $\rho_i^t = \text{sign}(\mathbf{w}_i^t)$



App: Sparse Regression via sub

Convex?

Non-differentiable

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \lambda \cdot \|\mathbf{w}\|_1$$

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

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

$O(nd)$ time per iter!

$\|\mathbf{w}\|_1$

$$\sum_{i=1}^d |\mathbf{w}_i^t|$$

$$\rho_i^t = \begin{cases} +1 & \text{if } \mathbf{w}_i^t > 0 \\ [-1, 1] & \text{if } \mathbf{w}_i^t = 0 \\ -1 & \text{if } \mathbf{w}_i^t < 0 \end{cases}$$

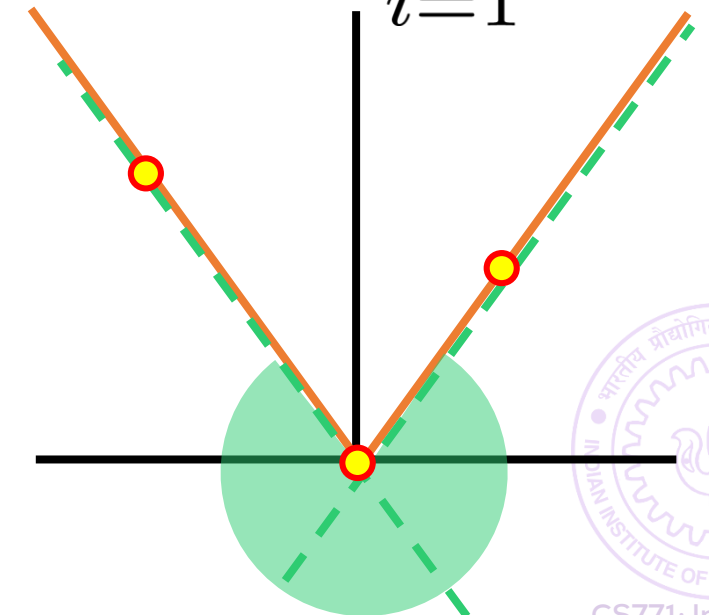
$\partial |\mathbf{w}_i|$

+1

-1

\mathbf{w}_i

e.g. $\rho_i^t = \text{sign}(\mathbf{w}_i^t)$



App: Logistic Regression via GD

August 30, 2017



App: Logistic Regression via GD

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

App: Logistic Regression via GD

Convex?

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

App: Logistic Regression via GD

Convex?

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

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

App: Logistic Regression via GD

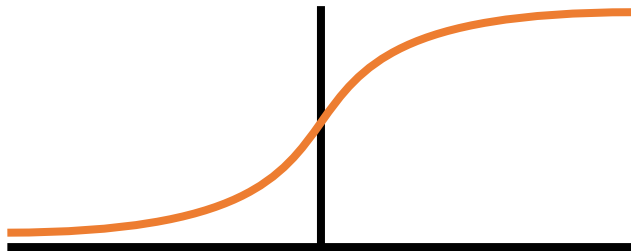
Convex?

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

$O(nd)$ time
per iter!

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

Whats going on
here?



App: Logistic Regression via GD

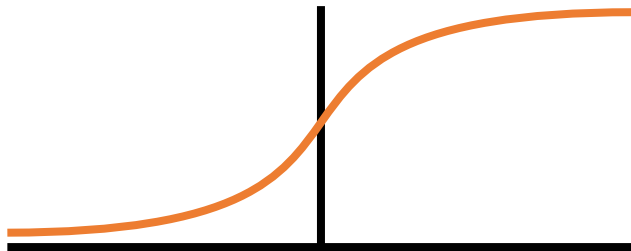
Convex?

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

$O(nd)$ time
per iter!

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

Whats going on
here?



Stochastic Gradient Method

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

Stochastic Gradient Method

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

Stochastic Gradient Method

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

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

Stochastic Gradient Method

$O(nd)$ 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 Method

$O(nd)$ 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 Method

$O(d)$ 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 Method

$O(d)$ 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 Method

$O(d)$ 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 Method

$O(d)$ 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 Method

$O(d)$ 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 Method

$O(d)$ 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 Method

$O(d)$ 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 Method

$O(d)$ 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 Method

$O(d)$ 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. 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 Method

$O(d)$ 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 $\mathbf{g}^t \leftarrow \nabla f_{I_t}(\mathbf{w}^t) + \nabla r(\mathbf{w}^t)$
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 Method

$O(d)$ 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 $\mathbf{g}^t \leftarrow \nabla f(\mathbf{w}^t, \mathbf{z}^t) + \nabla r(\mathbf{w}^t)$
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

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

Please give your Feedback

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