

COS234: INTRODUCTION TO MACHINE LEARNING

Prof. Yoram Singer

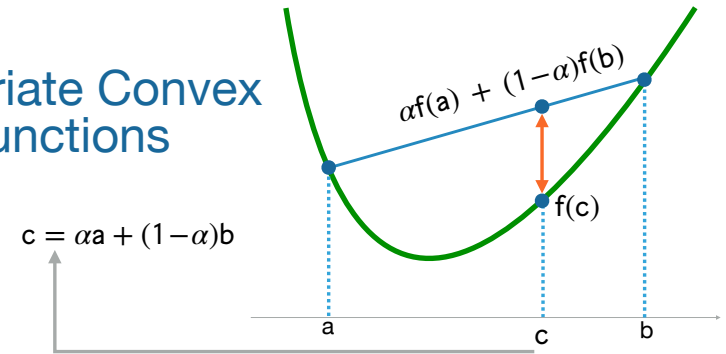


Topic: Gradient-Based Learning - Part II

© 2020 YORAM SINGER

$$\alpha \in [0, 1] : f(\alpha a + (1 - \alpha)b) \leq \alpha f(a) + (1 - \alpha)f(b)$$

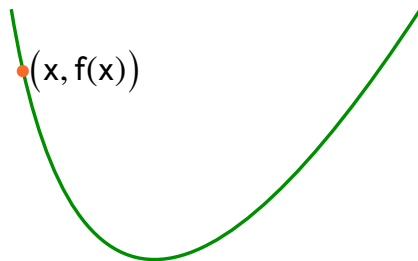
Univariate Convex Functions



© 2020 YORAM SINGER

2

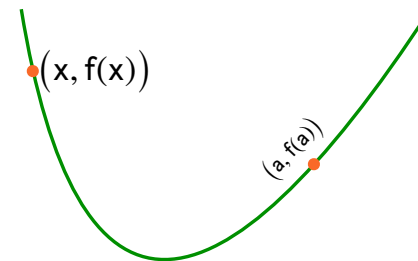
Univariate Convex Functions



© 2020 YORAM SINGER

3

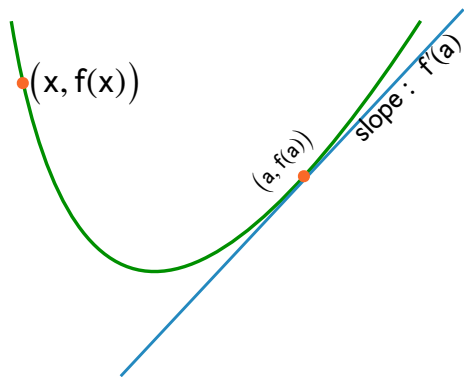
Univariate Convex Functions



© 2020 YORAM SINGER

3

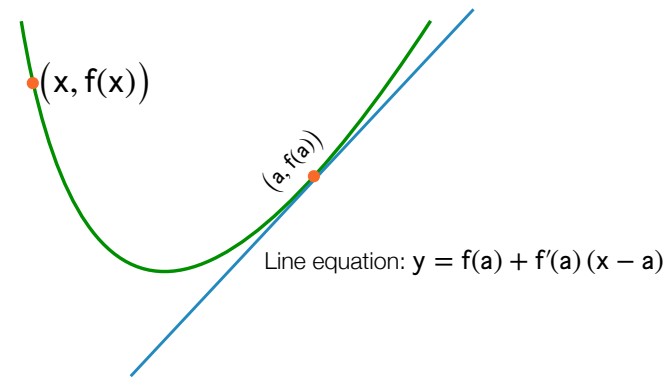
Univariate Convex Functions



© 2020 YORAM SINGER

3

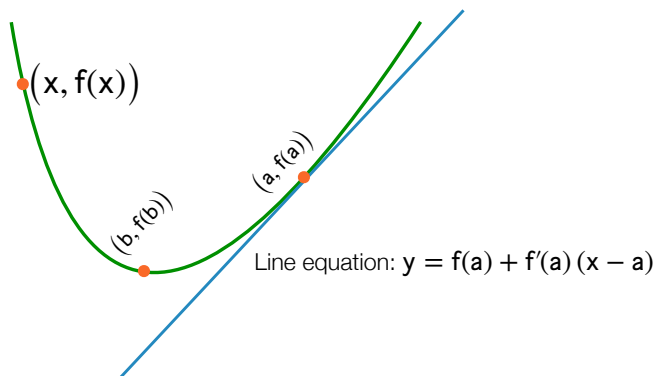
Univariate Convex Functions



© 2020 YORAM SINGER

3

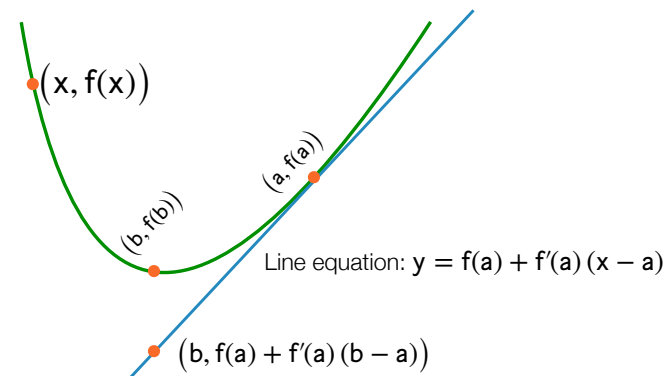
Univariate Convex Functions



© 2020 YORAM SINGER

3

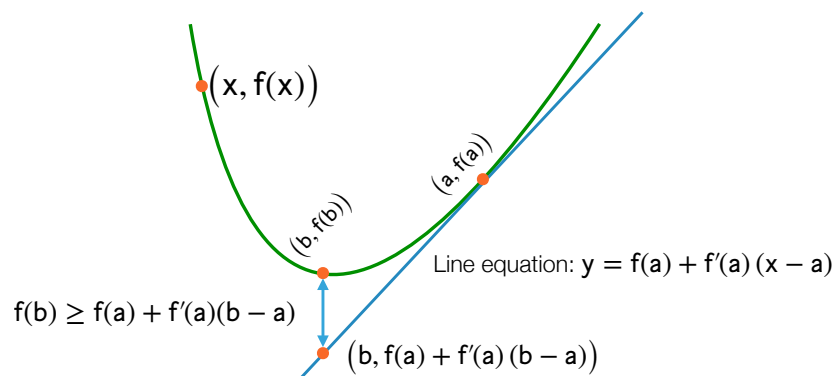
Univariate Convex Functions



© 2020 YORAM SINGER

3

Univariate Convex Functions



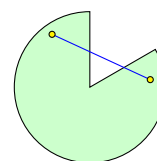
© 2020 YORAM SINGER

3

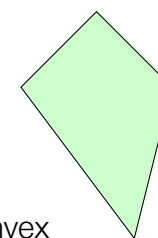
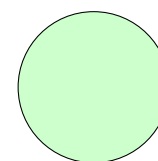
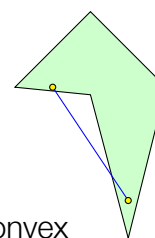
Convex Sets

Set Ω is convex if for any $\mathbf{u}, \mathbf{v} \in \Omega$ line segment between \mathbf{u} and \mathbf{v} is in Ω :

$$\forall \alpha \in [0, 1] : \alpha \mathbf{u} + (1 - \alpha) \mathbf{v} \in \Omega$$



Non-Convex



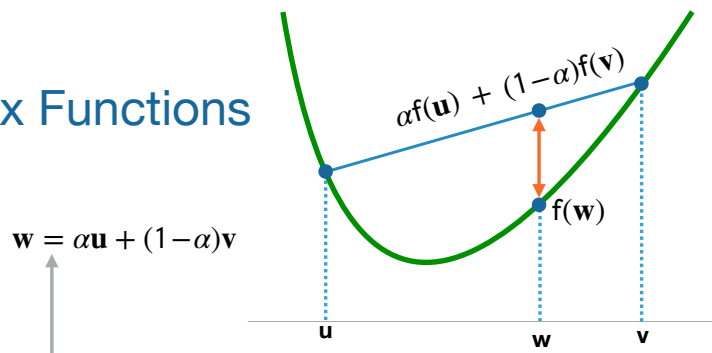
Convex

© 2020 YORAM SINGER

4

$$\alpha \in [0, 1] : f(\alpha \mathbf{u} + (1 - \alpha) \mathbf{v}) \leq \alpha f(\mathbf{u}) + (1 - \alpha) f(\mathbf{v})$$

Convex Functions



© 2020 YORAM SINGER

5

Gradients

$\mathcal{L}(\mathbf{w})$ is a function from $\mathbf{R}^d \rightarrow \mathbf{R}_+$

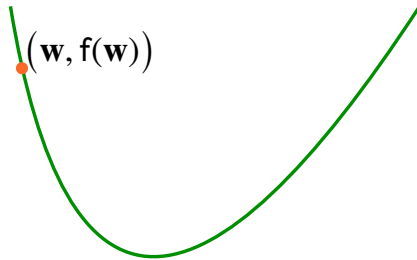
We need switch from *derivatives* to *gradients*:

$$\nabla \mathcal{L}(\mathbf{w}) \equiv \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = \left(\frac{\partial \mathcal{L}}{\partial w_1}, \frac{\partial \mathcal{L}}{\partial w_2}, \dots, \frac{\partial \mathcal{L}}{\partial w_d} \right)$$

© 2020 YORAM SINGER

6

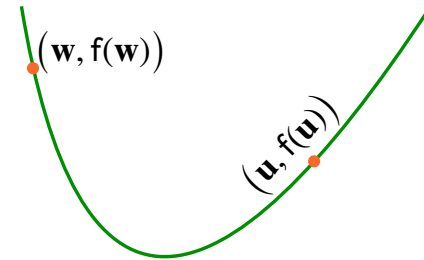
Multivariate Convex Functions



© 2020 YORAM SINGER

7

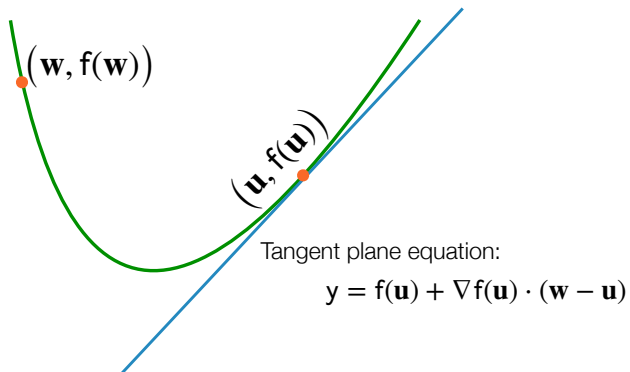
Multivariate Convex Functions



© 2020 YORAM SINGER

7

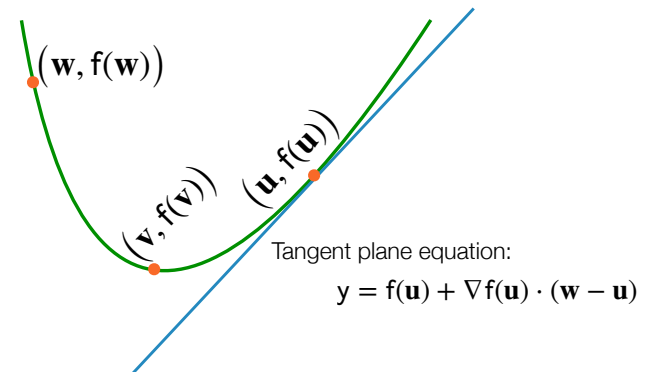
Multivariate Convex Functions



© 2020 YORAM SINGER

7

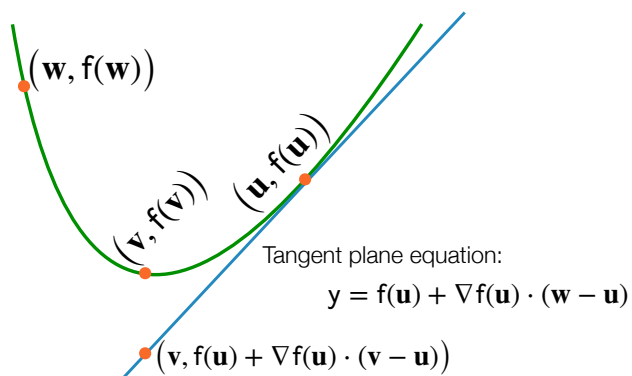
Multivariate Convex Functions



© 2020 YORAM SINGER

7

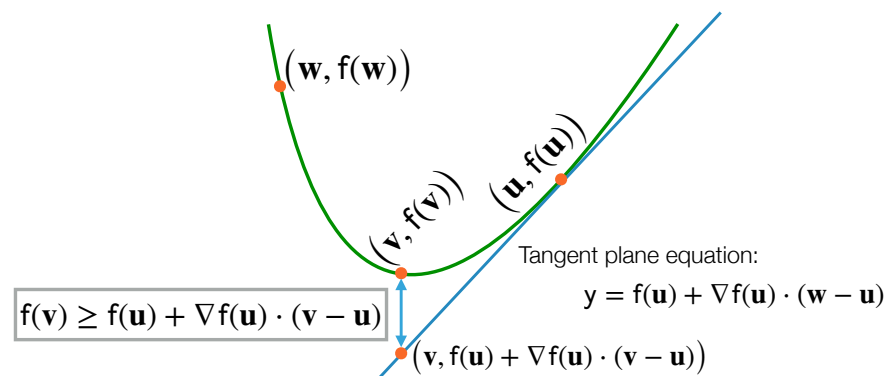
Multivariate Convex Functions



© 2020 YORAM SINGER

7

Multivariate Convex Functions



© 2020 YORAM SINGER

7

First Order Conditions

- Function $f : \mathbf{R}^d \rightarrow \mathbf{R}$ is convex iff:
 - For $\forall \alpha \in [0, 1], \mathbf{u} \in \mathbf{R}^d, \mathbf{v} \in \mathbf{R}^d$:

$$f(\alpha \mathbf{u} + (1 - \alpha) \mathbf{v}) \leq \alpha f(\mathbf{u}) + (1 - \alpha) f(\mathbf{v})$$
 - $f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u}) \cdot (\mathbf{v} - \mathbf{u})$
 - Fix $\mathbf{u} \in \mathbf{R}^d, \mathbf{v} \in \mathbf{R}^d$ and define $h : \mathbf{R} \rightarrow \mathbf{R}$ as $h(t) = f(\mathbf{u} + t\mathbf{v})$
 then $h(t)$ is a convex univariate function ($h''(t) \geq 0$)

© 2020 YORAM SINGER

8

Second Order Conditions

Hessian of $f : \mathbf{R}^d \rightarrow \mathbf{R}$ is $d \times d$ matrix of second order derivatives $\nabla^2 f$:

$$H_{ij} = \frac{\partial^2 f}{\partial w_i \partial w_j}$$

f is convex iff:

$$\forall \mathbf{u} : \mathbf{u}^T \mathbf{H} \mathbf{u} \geq 0$$

This means that the smallest Eigen value of H is non-negative

© 2020 YORAM SINGER

9

Second Order Conditions

- Function $f : \mathbf{R}^d \rightarrow \mathbf{R}$ is convex iff:

- For $\forall \alpha \in [0, 1], \mathbf{u} \in \mathbf{R}^d, \mathbf{v} \in \mathbf{R}^d$:

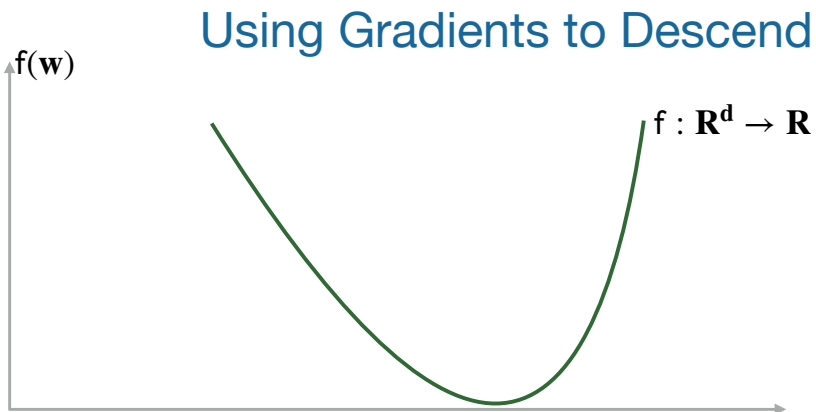
$$f(\alpha \mathbf{u} + (1 - \alpha)\mathbf{v}) \leq \alpha f(\mathbf{u}) + (1 - \alpha)f(\mathbf{v})$$

- $f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u}) \cdot (\mathbf{v} - \mathbf{u})$

- Fix $\mathbf{u} \in \mathbf{R}^d, \mathbf{v} \in \mathbf{R}^d$ and define $h : \mathbf{R} \rightarrow \mathbf{R}$ as $h(t) = f(\mathbf{u} + t\mathbf{v})$

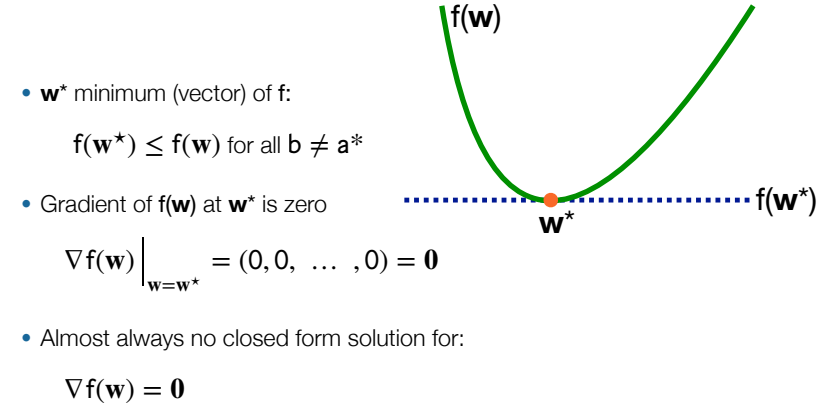
© 2020 YORAM SINGER

10



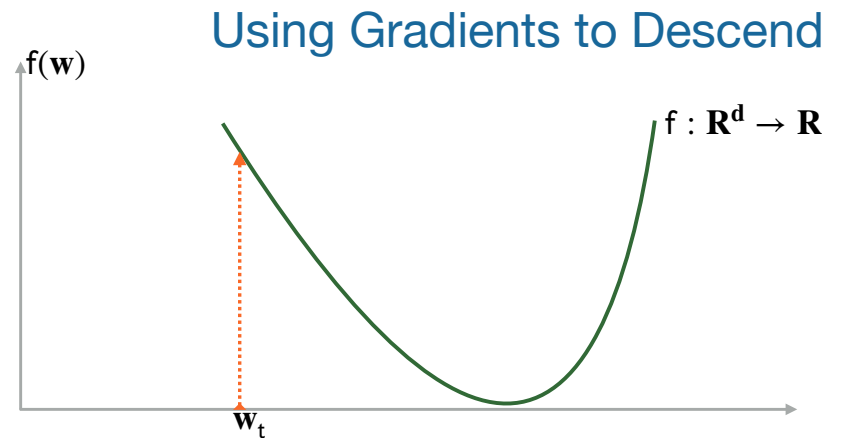
© 2020 YORAM SINGER

12



© 2020 YORAM SINGER

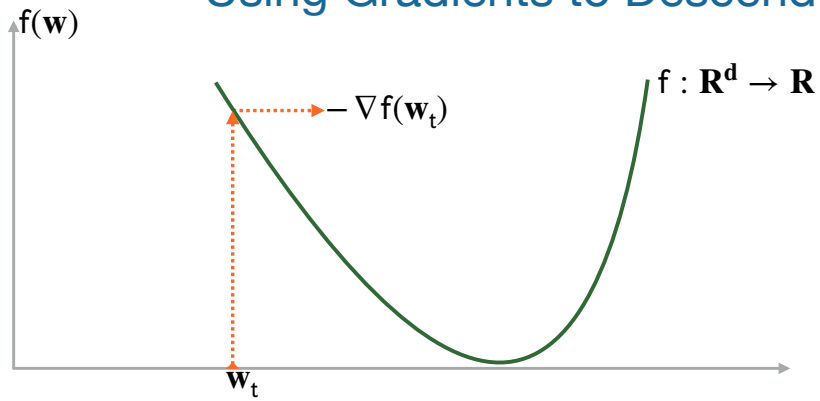
11



© 2020 YORAM SINGER

12

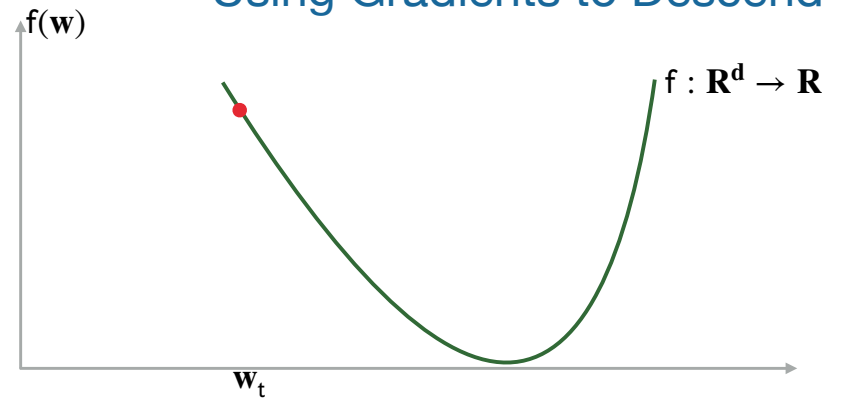
Using Gradients to Descend



© 2020 YORAM SINGER

12

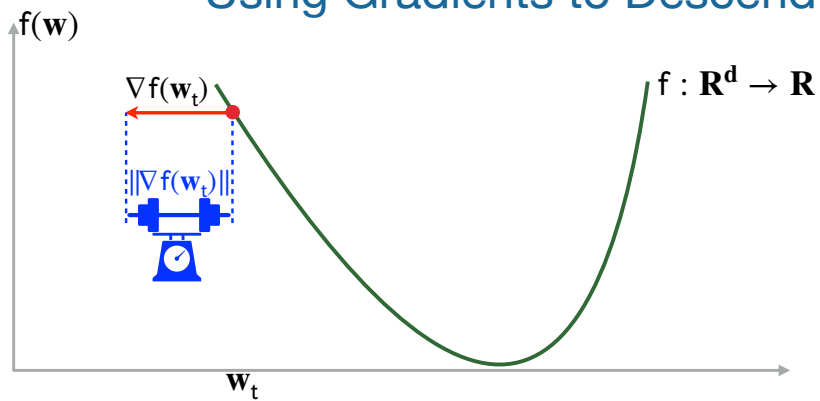
Using Gradients to Descend



© 2020 YORAM SINGER

12

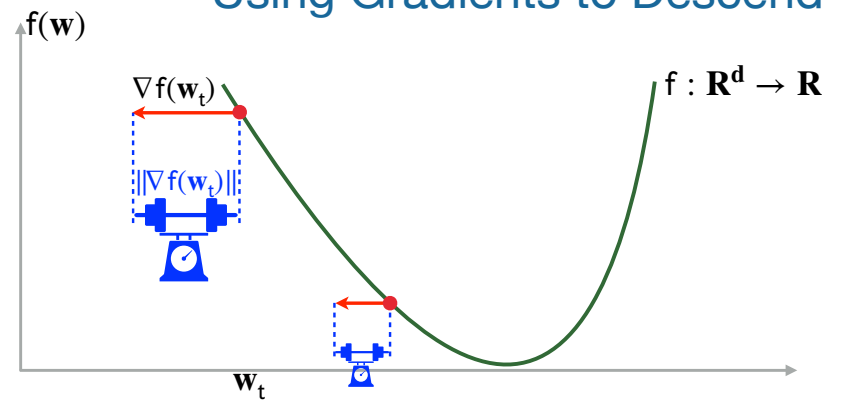
Using Gradients to Descend



© 2020 YORAM SINGER

12

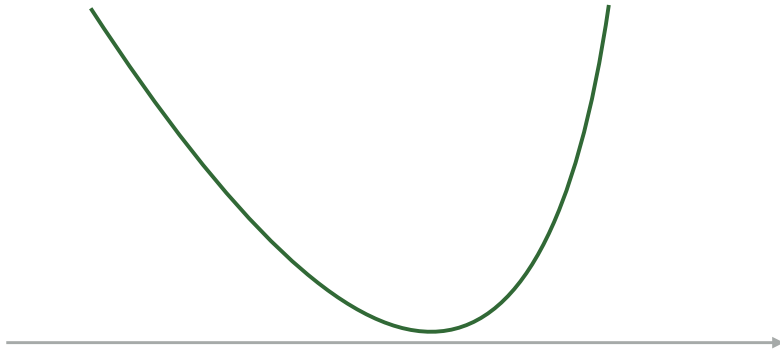
Using Gradients to Descend



© 2020 YORAM SINGER

12

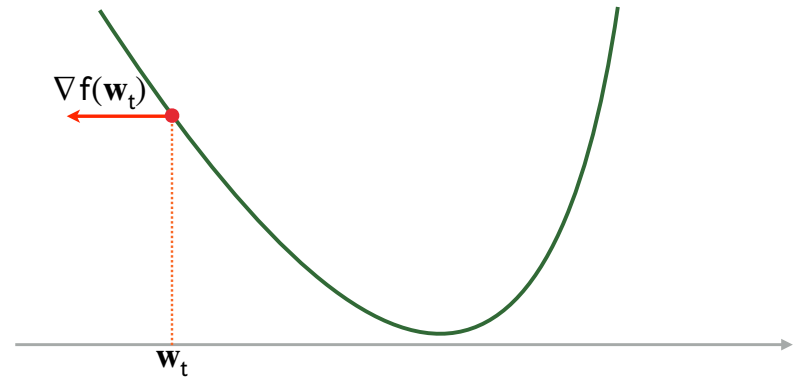
Gradient Descent



© 2020 YORAM SINGER

13

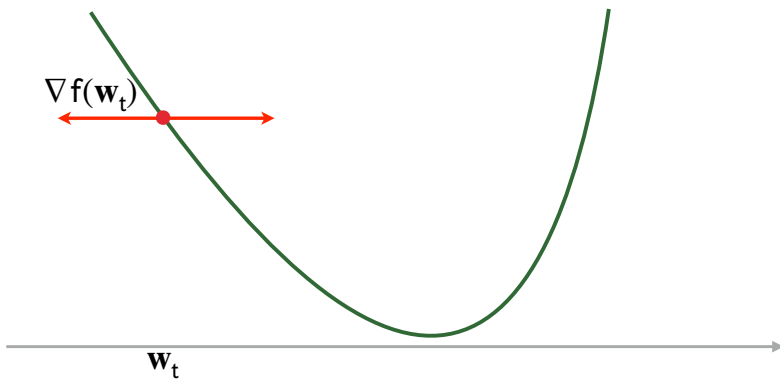
Gradient Descent



© 2020 YORAM SINGER

13

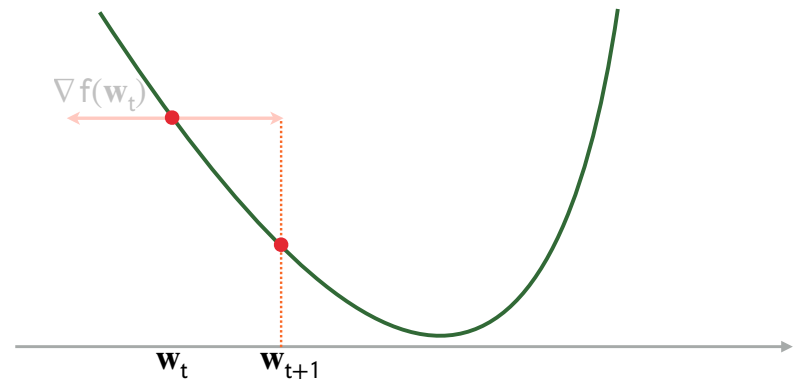
Gradient Descent



© 2020 YORAM SINGER

13

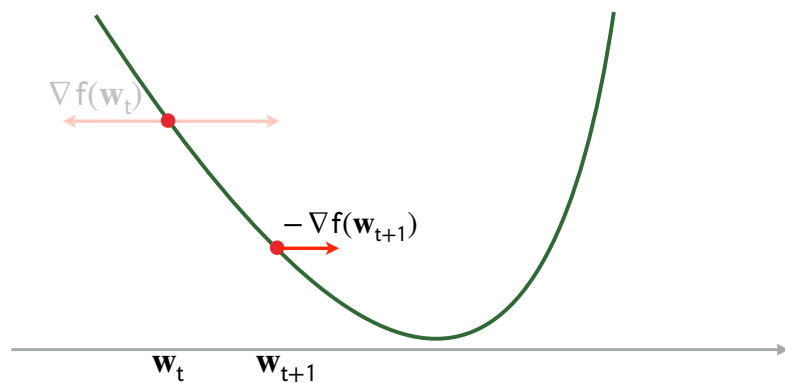
Gradient Descent



© 2020 YORAM SINGER

13

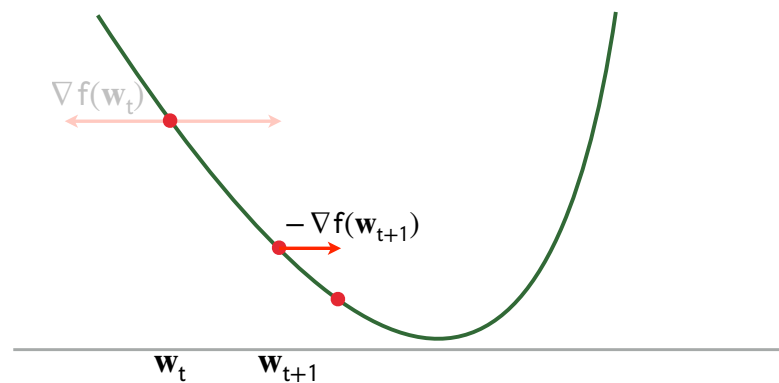
Gradient Descent



© 2020 YORAM SINGER

13

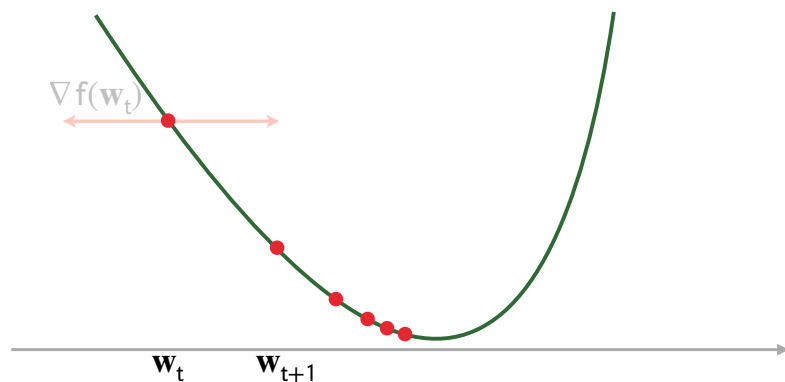
Gradient Descent



© 2020 YORAM SINGER

13

Gradient Descent



© 2020 YORAM SINGER

13

Gradient Descent Procedure

- Input: function $f : \mathbf{R} \rightarrow \mathbf{R}_+$
- Goal: find $\hat{\mathbf{w}}$ such that $f(\hat{\mathbf{w}}) - f(\mathbf{w}^*) \leq \epsilon$
- Choose initial value \mathbf{w}_1
- Loop:
 - $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t)$
- Until ...

© 2020 YORAM SINGER

14

Gradient Descent Procedure

- Input: function $f : \mathbf{R} \rightarrow \mathbf{R}_+$
- Goal: find $\hat{\mathbf{w}}$ such that $f(\hat{\mathbf{w}}) - f(\mathbf{w}^*) \leq \epsilon$
- Choose initial value \mathbf{w}_1
- Loop:
 - $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t)$
- Until ...

★ learning rate
★ step size

© 2020 YORAM SINGER

14

Gradient Descent

```
def derivative_descent(w0, gradient_func, eta):
    T, d = len(eta) + 1, len(w0)
    w = w0
    for t in range(1, T):
        w = w - eta[t-1] * gradient_func(w)
    return w
```

© 2020 YORAM SINGER

15

Learning Rate

- Crucial in many learning problems
- Fixed learning-rate can be used in restricted circumstances
- Self-tuning procedure of learning-rate exist, notably AdaGrad
- In many applications:
 - Linear decrease $\eta_t = \frac{\eta_0}{b + st}$ where $\eta_0 \in [0.1, 1]$ (typically)
 - Sub-linear decrease $\eta_t \sim \frac{\eta_0}{\sqrt{t}}$

© 2020 YORAM SINGER

16

Generalized Linear Models

Loss for linear predictors $\mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, \mathbf{h}(\mathbf{w} \cdot \mathbf{x}_i))$

\mathbf{h} : transfer (activation) function $\mathbf{h} : \mathbf{R} \rightarrow \mathbf{R}$ and $\hat{y}_i = \mathbf{h}(\mathbf{w} \cdot \mathbf{x}_i)$

$\ell(y, \hat{y})$ binary function from $\mathbf{R} \times \mathbf{R}$ to \mathbf{R}_+

We can find \mathbf{w} which (approximately) minimizes $\mathcal{L}(\mathbf{w})$ using GD and computing the gradient using chain rule

© 2020 YORAM SINGER

17

GLM: Examples

Linear regression:

$$y \in \mathbf{R} \quad h(z) = z \quad \ell(y, \hat{y}) = (y - \hat{y})^2$$

Classification with hinge-loss:

$$y \in \{-1, +1\} \quad h(z) = z \quad \ell(y, \hat{y}) = [1 - y\hat{y}]_+ \text{ where } [z]_+ = \max\{0, z\}$$

Logistic regression:

$$y \in \{0, 1\} \quad h(z) = \frac{1}{1 + e^{-z}} \quad \ell(y, \hat{y}) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

© 2020 YORAM SINGER

18

Gradients for GLM

$$\nabla \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \nabla_{\mathbf{w}} \ell(y_i, h(\mathbf{w} \cdot \mathbf{x}_i))$$

$$\text{Given } y_i \text{ derivative w.r.t } \hat{y}_i \text{ is } \ell'(y_i, \hat{y}_i) = \frac{d \ell(y_i, \hat{y}_i)}{d \hat{y}_i}$$

Define $\mathbf{z}_i = \mathbf{w} \cdot \mathbf{x}_i$ then derivative of h w.r.t. \mathbf{z}_i is $h'(\mathbf{z}_i)$

Finally $\nabla_{\mathbf{w}} \mathbf{z}_i = \mathbf{x}_i$

$$\text{Thus } \nabla \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell'(y_i, \hat{y}_i) h'(\mathbf{z}_i) \mathbf{x}_i$$

© 2020 YORAM SINGER

19

Gradients for GLM

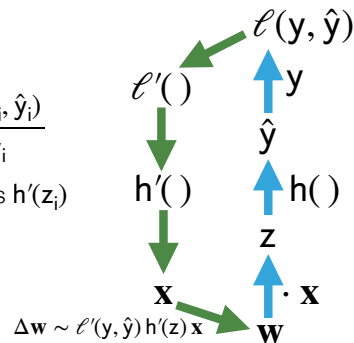
$$\nabla \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \nabla_{\mathbf{w}} \ell(y_i, h(\mathbf{w} \cdot \mathbf{x}_i))$$

$$\text{Given } y_i \text{ derivative w.r.t } \hat{y}_i \text{ is } \ell'(y_i, \hat{y}_i) = \frac{d \ell(y_i, \hat{y}_i)}{d \hat{y}_i}$$

Define $\mathbf{z}_i = \mathbf{w} \cdot \mathbf{x}_i$ then derivative of h w.r.t. \mathbf{z}_i is $h'(\mathbf{z}_i)$

Finally $\nabla_{\mathbf{w}} \mathbf{z}_i = \mathbf{x}_i$

$$\text{Thus } \nabla \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell'(y_i, \hat{y}_i) h'(\mathbf{z}_i) \mathbf{x}_i$$



© 2020 YORAM SINGER

19

Gradient for Logistic Regression

© 2020 YORAM SINGER

20

Gradient for Logistic Regression

$$y \in \{0, 1\} \quad h(z) = \frac{1}{1 + e^{-z}} \quad \ell(y, \hat{y}) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

Gradient for Logistic Regression

$$y \in \{0, 1\} \quad h(z) = \frac{1}{1 + e^{-z}} \quad \ell(y, \hat{y}) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

$$\ell'(y, \hat{y}) = \frac{1 - y}{1 - \hat{y}} - \frac{y}{\hat{y}} = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})}$$

Gradient for Logistic Regression

$$y \in \{0, 1\} \quad h(z) = \frac{1}{1 + e^{-z}} \quad \ell(y, \hat{y}) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

$$\ell'(y, \hat{y}) = \frac{1 - y}{1 - \hat{y}} - \frac{y}{\hat{y}} = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})}$$

$$h'(z) = h(z)(1 - h(z)) = \hat{y}(1 - \hat{y})$$

Gradient for Logistic Regression

$$y \in \{0, 1\} \quad h(z) = \frac{1}{1 + e^{-z}} \quad \ell(y, \hat{y}) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

$$\ell'(y, \hat{y}) = \frac{1 - y}{1 - \hat{y}} - \frac{y}{\hat{y}} = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})}$$

$$h'(z) = h(z)(1 - h(z)) = \hat{y}(1 - \hat{y}) \quad \longrightarrow \quad \ell'(y, \hat{y}) h'(z) = \hat{y} - y$$

Gradient for Logistic Regression

$$y \in \{0, 1\} \quad h(z) = \frac{1}{1 + e^{-z}} \quad \ell(y, \hat{y}) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

$$\ell'(y, \hat{y}) = \frac{1 - y}{1 - \hat{y}} - \frac{y}{\hat{y}} = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})}$$



$$\ell'(y, \hat{y}) h'(z) = \hat{y} - y$$

$$h'(z) = h(z)(1 - h(z)) = \hat{y}(1 - \hat{y})$$

$$\nabla \mathcal{L}(\mathbf{w}) = -\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i) \mathbf{x}_i$$

© 2020 YORAM SINGER

20

Gradient for Logistic Regression

$$y \in \{0, 1\} \quad h(z) = \frac{1}{1 + e^{-z}} \quad \ell(y, \hat{y}) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

$$\ell'(y, \hat{y}) = \frac{1 - y}{1 - \hat{y}} - \frac{y}{\hat{y}} = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})}$$



$$\ell'(y, \hat{y}) h'(z) = \hat{y} - y$$

$$h'(z) = h(z)(1 - h(z)) = \hat{y}(1 - \hat{y})$$

$$\nabla \mathcal{L}(\mathbf{w}) = -\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i) \mathbf{x}_i$$

Interpretation?

© 2020 YORAM SINGER

20

Properties of GD

Assume that $\|\nabla \mathcal{L}(\mathbf{w})\| \leq \gamma$

This typically amounts to $\forall i : \|\mathbf{x}_i\| \leq c\gamma$

Assume that $\forall t : \|\mathbf{w}_t\| \leq r$ and $\|\mathbf{w}^*\| \leq r$

Recall that from convexity $\mathbf{f}(\mathbf{v}) \geq \mathbf{f}(\mathbf{u}) + \nabla \mathbf{f}(\mathbf{u}) \cdot (\mathbf{v} - \mathbf{u})$

Using $\mathbf{v} = \mathbf{w}^*$ and $\mathbf{u} = \mathbf{w}_t$ we get

$$\mathcal{L}(\mathbf{w}_t) - \mathcal{L}(\mathbf{w}^*) \leq \nabla \mathcal{L}(\mathbf{w}_t) \cdot (\mathbf{w}_t - \mathbf{w}^*)$$

© 2020 YORAM SINGER

21

Properties of GD (cont)

Going forward we abbreviate $\mathbf{g}_t = \nabla \mathcal{L}(\mathbf{w}_t)$ and thus we can write

$$\mathcal{L}(\mathbf{w}_t) - \mathcal{L}(\mathbf{w}^*) \leq \mathbf{g}_t \cdot (\mathbf{w}_t - \mathbf{w}^*)$$

Define progress towards (unknown) optimum as $\Delta_t = \|\mathbf{w}_t - \mathbf{w}^*\|^2$

From GD update & boundedness of gradients:

$$\Delta_{t+1} \leq \Delta_t - 2\eta \mathbf{g}_t \cdot (\mathbf{w}_t - \mathbf{w}^*) + \eta^2 \gamma^2$$

Rearranging terms and dividing by 2η gives

$$\mathbf{g}_t \cdot (\mathbf{w}_t - \mathbf{w}^*) \leq \frac{1}{2\eta} (\Delta_t - \Delta_{t+1}) + \frac{\eta \gamma^2}{2}$$

© 2020 YORAM SINGER

22

Properties of GD (cont)

Going forward we abbreviate $\mathbf{g}_t = \nabla \mathcal{L}(\mathbf{w}_t)$ and thus we can write

$$\mathcal{L}(\mathbf{w}_t) - \mathcal{L}(\mathbf{w}^*) \leq \mathbf{g}_t \cdot (\mathbf{w}_t - \mathbf{w}^*) \leq \frac{1}{2\eta}(\Delta_t - \Delta_{t+1}) + \frac{\eta\gamma^2}{2}$$

Define progress towards (unknown) optimum as $\Delta_t = \|\mathbf{w}_t - \mathbf{w}^*\|^2$

From GD update & boundedness of gradients:

$$\Delta_{t+1} \leq \Delta_t - 2\eta \mathbf{g}_t \cdot (\mathbf{w}_t - \mathbf{w}^*) + \eta^2 \gamma^2$$

Rearranging terms and dividing by 2η gives

$$\mathbf{g}_t \cdot (\mathbf{w}_t - \mathbf{w}^*) \leq \frac{1}{2\eta}(\Delta_t - \Delta_{t+1}) + \frac{\eta\gamma^2}{2}$$

© 2020 YORAM SINGER

22

We established that $\mathcal{L}(\mathbf{w}_t) - \mathcal{L}(\mathbf{w}^*) \leq \frac{1}{2\eta}(\Delta_t - \Delta_{t+1}) + \frac{\eta\gamma^2}{2}$

Taking the average from $t=1$ through T

$$\frac{1}{T} \sum_{t=1}^T \mathcal{L}(\mathbf{w}_t) - \mathcal{L}(\mathbf{w}^*) \leq \frac{1}{2\eta T}(\Delta_1 - \Delta_{T+1}) + \frac{\eta\gamma^2}{2} \leq \frac{r^2}{\eta T} + \frac{\eta\gamma^2}{2}$$

Assume $\forall t: \mathcal{L}(\mathbf{w}_t) \leq \mathcal{L}(\mathbf{w}_{t-1})$ then $\mathcal{L}(\mathbf{w}_T) - \mathcal{L}(\mathbf{w}^*) \leq \frac{r^2}{\eta T} + \frac{\eta\gamma^2}{2}$

Choosing $\eta = r/(\gamma T^{1/2})$ gives $\mathcal{L}(\mathbf{w}_T) - \mathcal{L}(\mathbf{w}^*) \leq \frac{2r\gamma}{\sqrt{T}}$

© 2020 YORAM SINGER

23

Convergence of GD

Assumptions:

$$\|\nabla \mathcal{L}(\mathbf{w})\| \leq \gamma, \forall t: \|\mathbf{w}_t\| \leq r, \|\mathbf{w}^*\| \leq r$$

$\mathcal{L}(\mathbf{w}_t)$ is monotonically decreasing

Using learning rate of $\eta = r/(\gamma\sqrt{t})$

Then $\mathcal{L}(\mathbf{w}_t)$ converges to $\mathcal{L}(\mathbf{w}^*)$ at a rate of $O(\gamma r t^{-1/2})$

© 2020 YORAM SINGER

24

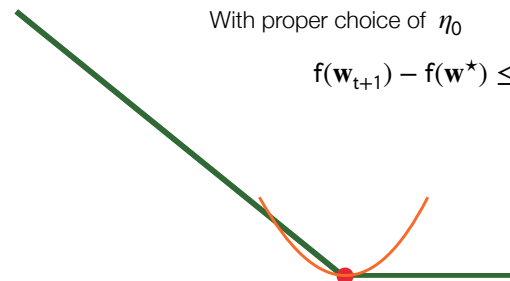
Summary: GD for Convex Losses

Learning rate goes to zero $\eta_0 t^{-1/2}$

With proper choice of η_0

$$f(\mathbf{w}_{t+1}) - f(\mathbf{w}^*) \leq \frac{\kappa}{\sqrt{t}}$$

problem dependent



© 2020 YORAM SINGER

25

Next

When data is prevalent:

- Stochastic Gradient Descent (SGD)

Solving multiclass problems with SGD

Beyond GLIM: deep learning