



Deep Learning

CS60010

Abir Das

Computer Science and Engineering Department
Indian Institute of Technology Kharagpur

<http://cse.iitkgp.ac.in/~adas/>



Agenda

- Understand basics of Matrix/Vector Calculus and Optimization concepts to be used in the course.
- Understand different types or errors in learning



Resources

- "Deep Learning", I. Goodfellow, Y. Bengio, A. Courville. (Chapter 8)



Optimization and Deep Learning- Connections

- Deep learning (machine learning, in general) involves optimization in many contexts
- The goal is to find parameters θ of a neural network that significantly reduce a cost function or objective function $J(\theta)$.
- Gradient based optimization is the most popular way for training Deep Neural Networks.
- There are other ways too, *e.g.*, evolutionary or derivative free optimization, but they come with issues particularly crucial for neural network training.
- Its easy to spend a semester on optimization. Thus, these few lectures will only be a scratch on a very small part of the surface.



Optimization and Deep Learning- Differences

- In learning we care about some performance measure P (e.g., image classification accuracy, language translation accuracy etc.) on **test set**, but we minimize a different cost function $J(\theta)$ on **training set**, with the hope that doing so will improve P
- This is in contrast to pure optimization where minimizing $J(\theta)$ is a goal in itself
- Lets see what types of errors creep in as a result

Expected and Empirical Risk

- Training data is $\{\mathbf{x}^{(i)}, y^{(i)} : i = 1, \dots, N\}$ coming from probability distribution $\mathcal{D}(\mathbf{x}, y)$
- Let the neural network learns the output functions as $g^{\mathcal{D}}(\mathbf{x})$ and the loss is denoted as $l(g^{\mathcal{D}}(\mathbf{x}), y)$
- What we want to minimize is the expected risk

$$E(g^{\mathcal{D}}) = \int l(g^{\mathcal{D}}(\mathbf{x}), y) d\mathcal{D}(\mathbf{x}, y) = \mathbb{E} \left(l(g^{\mathcal{D}}(\mathbf{x}), y) \right)$$

- If we could minimize this risk we would have got the true optimal function

$$f = \arg \min_{g^{\mathcal{D}}} E(g^{\mathcal{D}})$$

- But we don't know the actual distribution that generates the data. So, what we actually minimize is the empirical risk

$$E_n(g^{\mathcal{D}}) = \frac{1}{N} \sum_{i=1}^N l(g^{\mathcal{D}}(\mathbf{x}_i), y_i) = \mathbb{E}_n \left(l(g^{\mathcal{D}}(\mathbf{x}), y) \right)$$

- We choose a family \mathcal{H} of candidate prediction functions and find the function that minimizes the empirical risk

$$g_n^{\mathcal{D}} = \arg \min_{g^{\mathcal{D}} \in \mathcal{H}} E_n(g^{\mathcal{D}})$$

- Since f may not be found in the family \mathcal{H} , we also define

$$g_{\mathcal{H}^*}^{\mathcal{D}} = \arg \min_{g^{\mathcal{D}} \in \mathcal{H}} E(g^{\mathcal{D}})$$



Sources of Error

	Minimizes	Staying within	
f	expected risk	No constraints on function family	<ul style="list-style-type: none">• True data distribution known• Family of functions exhaustive
$g_{\mathcal{H}*}^D$	expected risk	family of functions \mathcal{H}	<ul style="list-style-type: none">• True data distribution known• Family of functions not exhaustive
g_n^D	empirical risk	family of functions \mathcal{H}	<ul style="list-style-type: none">• True data distribution not known• Family of functions not exhaustive

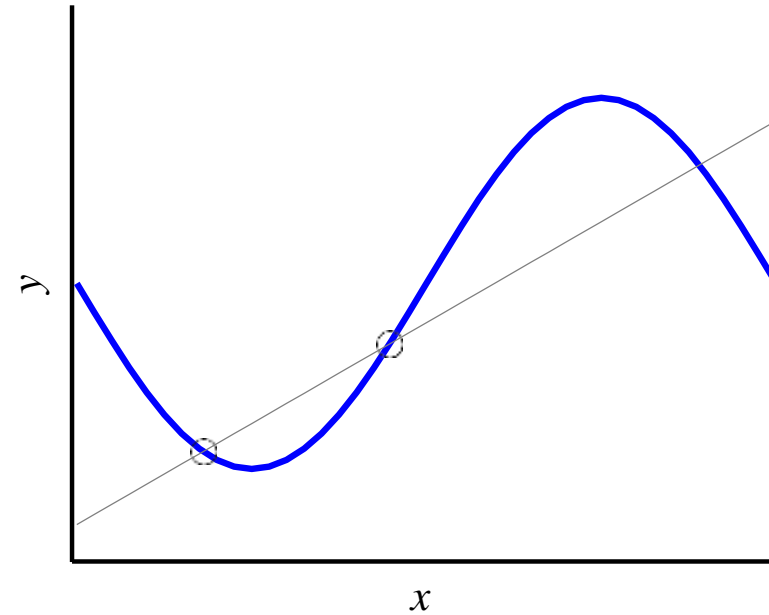
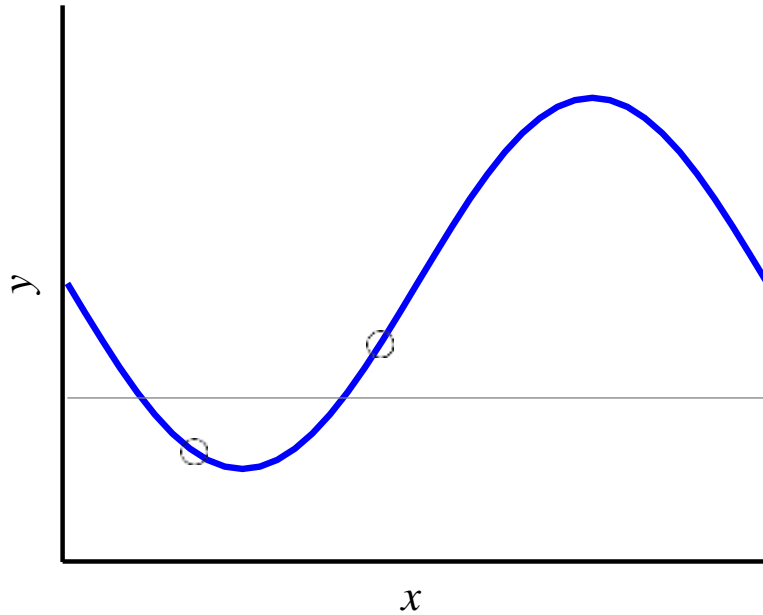
- True data generation procedure not known
- Family of functions [or hypothesis] to try is not exhaustive
- (And not to forget) we rely on a surrogate loss in place of the true classification error rate

A Simple Learning Problem

2 Data Points. 2 hypothesis sets:

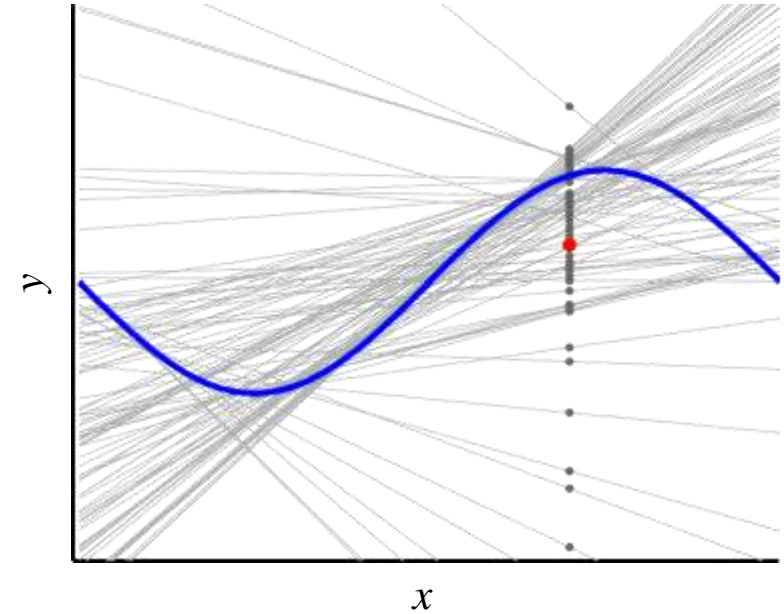
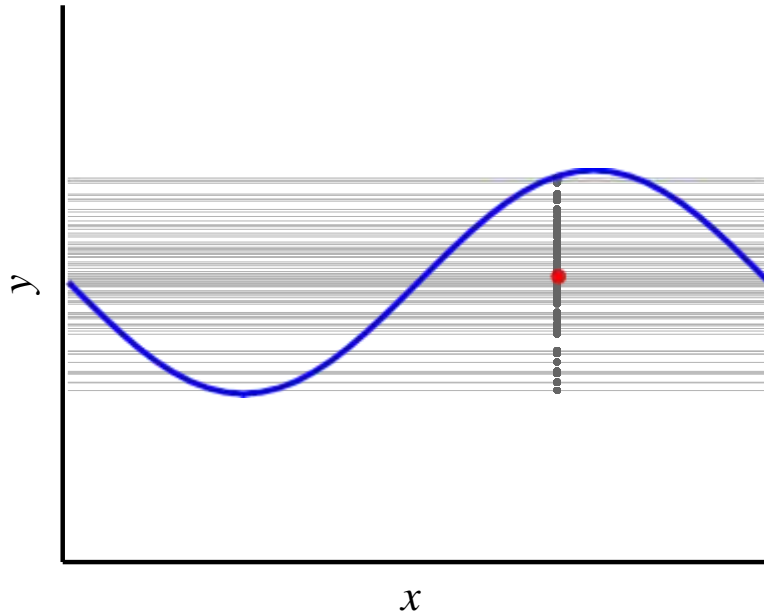
$$\mathcal{H}_0: h(x) = b$$

$$\mathcal{H}_1: h(x) = ax + b$$



Slide courtesy: Malik Magdon-Ismail

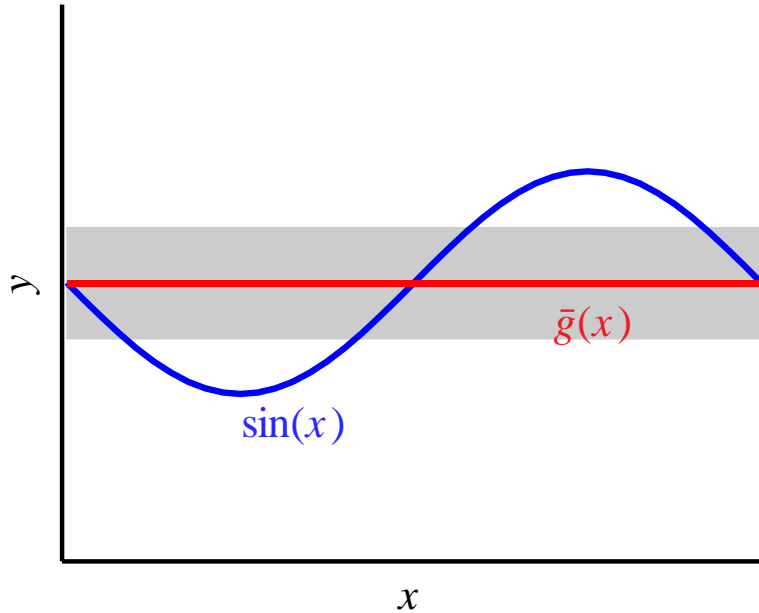
Let's Repeat the Experiment Many Times



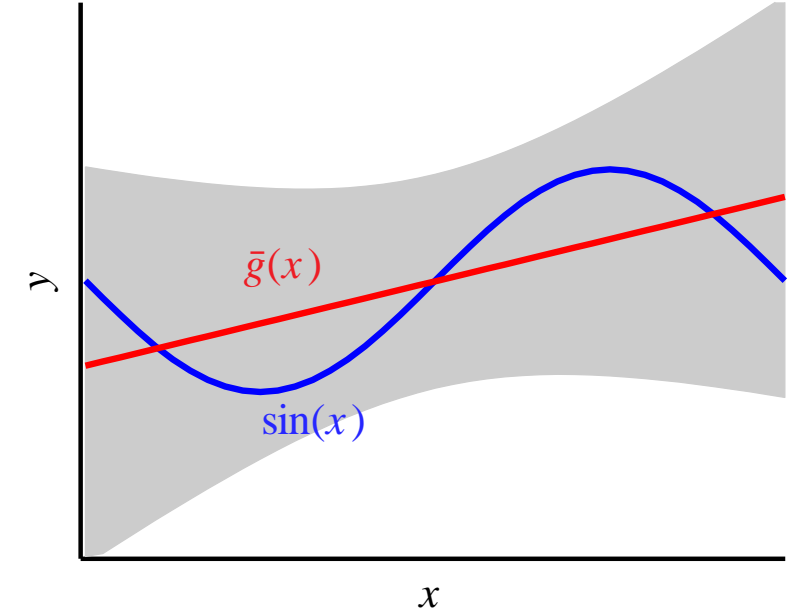
- For each data set \mathcal{D} you get a different $g_n^{\mathcal{D}}$
- For a fixed x , $g_n^{\mathcal{D}}(x)$ is random value, depending on \mathcal{D}

Slide courtesy: Malik Magdon-Ismail

What's Happening on Average



We can define



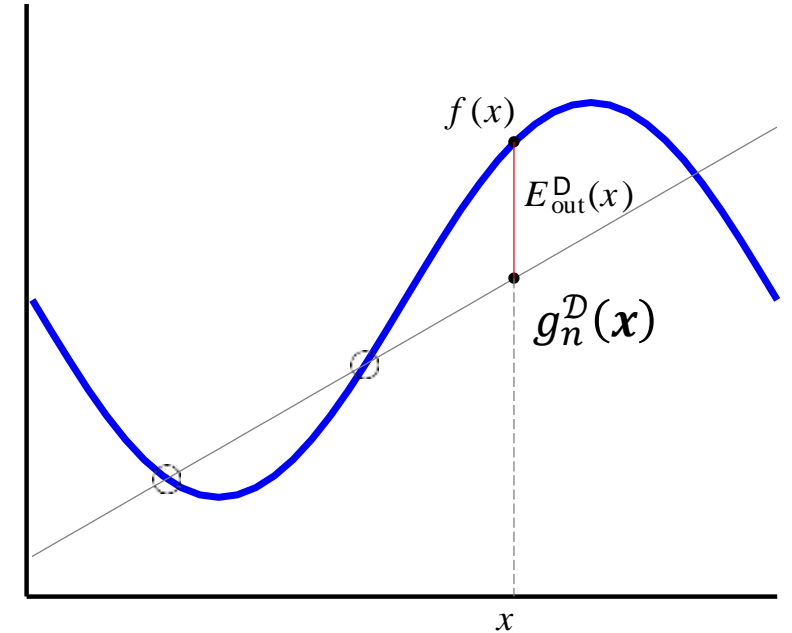
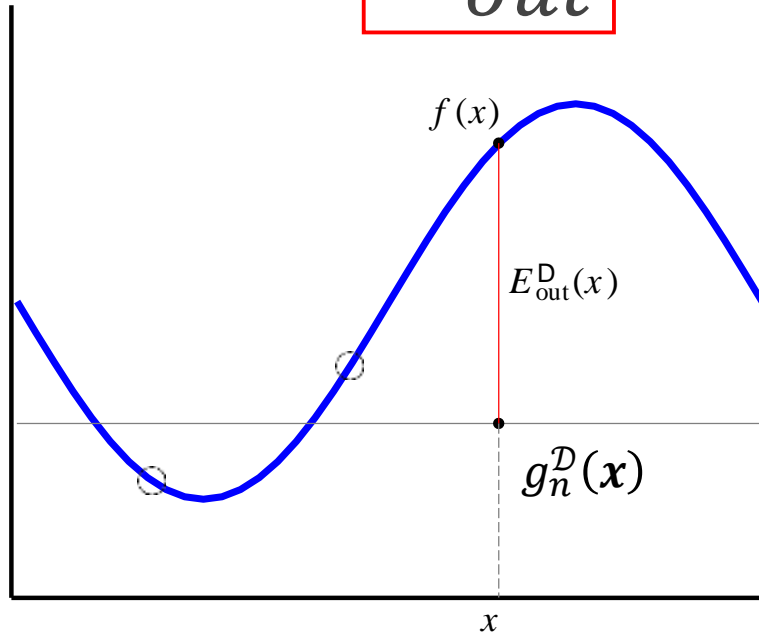
$g_n^{\mathcal{D}}(x)$ \longleftarrow Random value, depending on \mathcal{D}

$\bar{g}(x) = \mathbb{E}_{\mathcal{D}}[g_n^{\mathcal{D}}(x)] = \frac{1}{K} (g_n^{\mathcal{D}_1}(x) + g_n^{\mathcal{D}_2}(x) + \dots + g_n^{\mathcal{D}_K}(x))$ \longleftarrow Your average prediction on x

$\text{var}(x) = \mathbb{E}_{\mathcal{D}} \left[\left(g_n^{\mathcal{D}}(x) - \bar{g}(x) \right)^2 \right] = \mathbb{E}_{\mathcal{D}} [g_n^{\mathcal{D}}(x)^2] - \bar{g}(x)^2$ \longleftarrow How variable is your prediction

Slide courtesy: Malik Magdon-Ismail

E_{out} on Test Point x for Data \mathcal{D}



$$E_{out}^{\mathcal{D}}(x) = \left(g_n^{\mathcal{D}}(x) - f(x)\right)^2 \leftarrow \text{Squared error, a random value depending on } \mathcal{D}$$

$$E_{out}(x) = \mathbb{E}_{\mathcal{D}}[E_{out}^{\mathcal{D}}(x)] = \mathbb{E}_{\mathcal{D}}\left[\left(g_n^{\mathcal{D}}(x) - f(x)\right)^2\right] \leftarrow \text{Expected value of the above random variable}$$

Slide motivation: Malik Magdon-Ismail



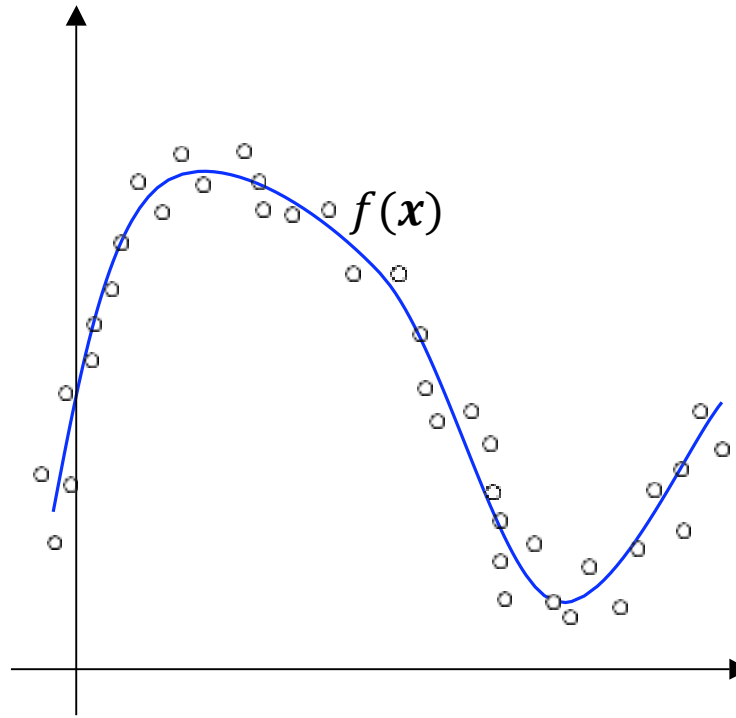
The Bias-Variance Decomposition

$$\begin{aligned} E_{out}(\mathbf{x}) &= \mathbb{E}_{\mathcal{D}} \left[\left(g_n^{\mathcal{D}}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right] \\ &= \mathbb{E}_{\mathcal{D}} [f(\mathbf{x})^2 - 2f(\mathbf{x})g_n^{\mathcal{D}}(\mathbf{x}) + g_n^{\mathcal{D}}(\mathbf{x})^2] \\ &= f(\mathbf{x})^2 - 2f(\mathbf{x})\mathbb{E}_{\mathcal{D}}[g_n^{\mathcal{D}}(\mathbf{x})] + \mathbb{E}_{\mathcal{D}}[g_n^{\mathcal{D}}(\mathbf{x})^2] \\ &= f(\mathbf{x})^2 - 2f(\mathbf{x})\bar{g}(\mathbf{x}) + \mathbb{E}_{\mathcal{D}}[g_n^{\mathcal{D}}(\mathbf{x})^2] \\ &= f(\mathbf{x})^2 - \bar{g}(\mathbf{x})^2 + \bar{g}(\mathbf{x})^2 - 2f(\mathbf{x})\bar{g}(\mathbf{x}) + \mathbb{E}_{\mathcal{D}}[g_n^{\mathcal{D}}(\mathbf{x})^2] \\ &= \underbrace{(f(\mathbf{x}) - \bar{g}(\mathbf{x}))^2}_{\text{Bias}} + \underbrace{\mathbb{E}_{\mathcal{D}}[g_n^{\mathcal{D}}(\mathbf{x})^2] - \bar{g}(\mathbf{x})^2}_{\text{Variance}} \end{aligned}$$

$$E_{out}(\mathbf{x}) = \text{Bias} + \text{Variance}$$

Bias-Variance to Overfitting-Underfitting

- Suppose the true underlying function is $f(x)$.
- But the observations (data points) are noisy i.e., $f(x) + \varepsilon$



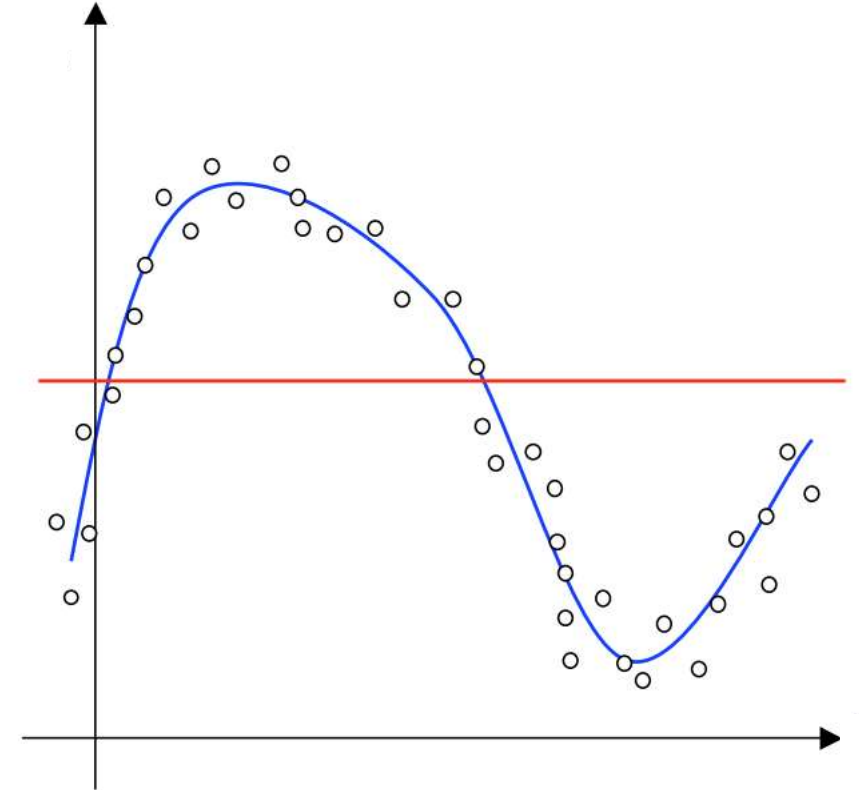
The Extreme Cases of Bias and Variance - Under-fitting

A good way to understand the concepts of bias and variance is by considering the two extreme cases of what a neural network might learn.

Suppose the neural network is lazy and just produces the same constant output whatever training data we give it, i.e. $g_n^{\mathcal{D}}(\mathbf{x}) = c$, then

$$\begin{aligned} E_{out}(\mathbf{x}) &= (f(\mathbf{x}) - \bar{g}(\mathbf{x}))^2 + \mathbb{E}_{\mathcal{D}}[g_n^{\mathcal{D}}(\mathbf{x})^2] - \bar{g}(\mathbf{x})^2 \\ &= (f(\mathbf{x}) - c)^2 + \mathbb{E}_{\mathcal{D}}[c^2] - c^2 \\ &= (f(\mathbf{x}) - c)^2 + 0 \end{aligned}$$

In this case the variance term will be zero, but the bias will be large, because the network has made no attempt to fit the data. We say we have extreme under-fitting



Ignore the data \Rightarrow

Big approximation error (high bias)

No variation between data sets (no variance)

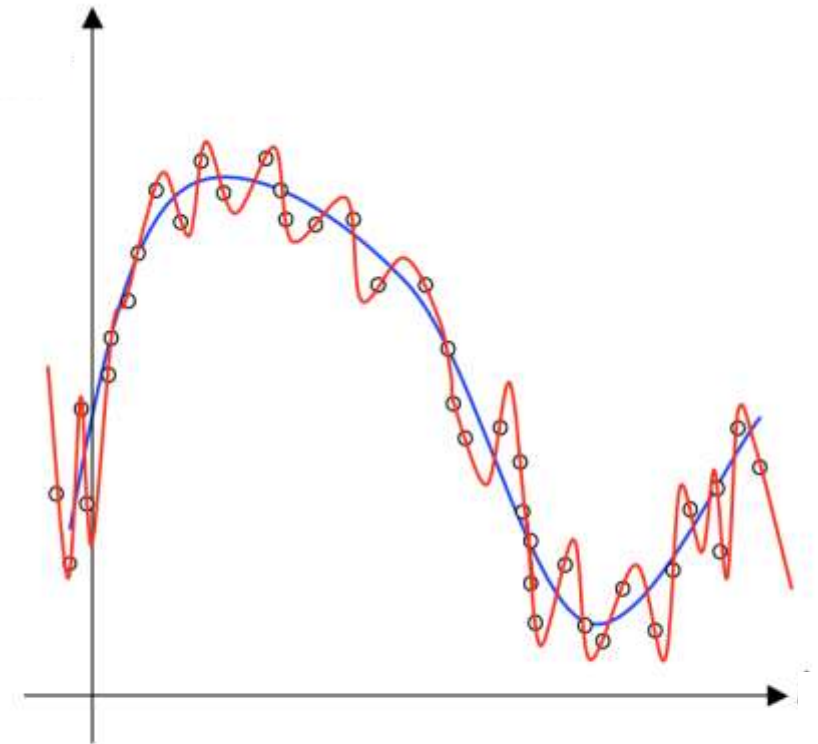
Slide motivation: John A. Bullinaria

The Extreme Cases of Bias and Variance - Over-fitting

On the other hand, suppose the neural network is very hard working and makes sure that it exactly fits every data point, i.e. $g_n^{\mathcal{D}}(\mathbf{x}) = f(\mathbf{x}) + \varepsilon$, then, $\bar{g}(\mathbf{x}) = \mathbb{E}_{\mathcal{D}}[g_n^{\mathcal{D}}(\mathbf{x})] = \mathbb{E}_{\mathcal{D}}[f(\mathbf{x}) + \varepsilon] = \mathbb{E}_{\mathcal{D}}[f(\mathbf{x})] + 0 = f(\mathbf{x})$

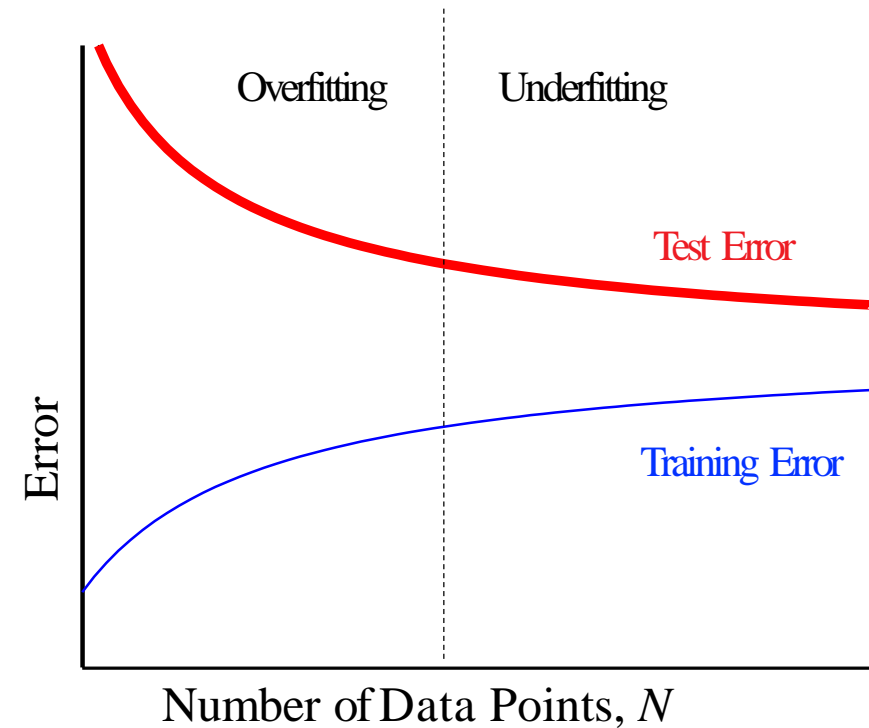
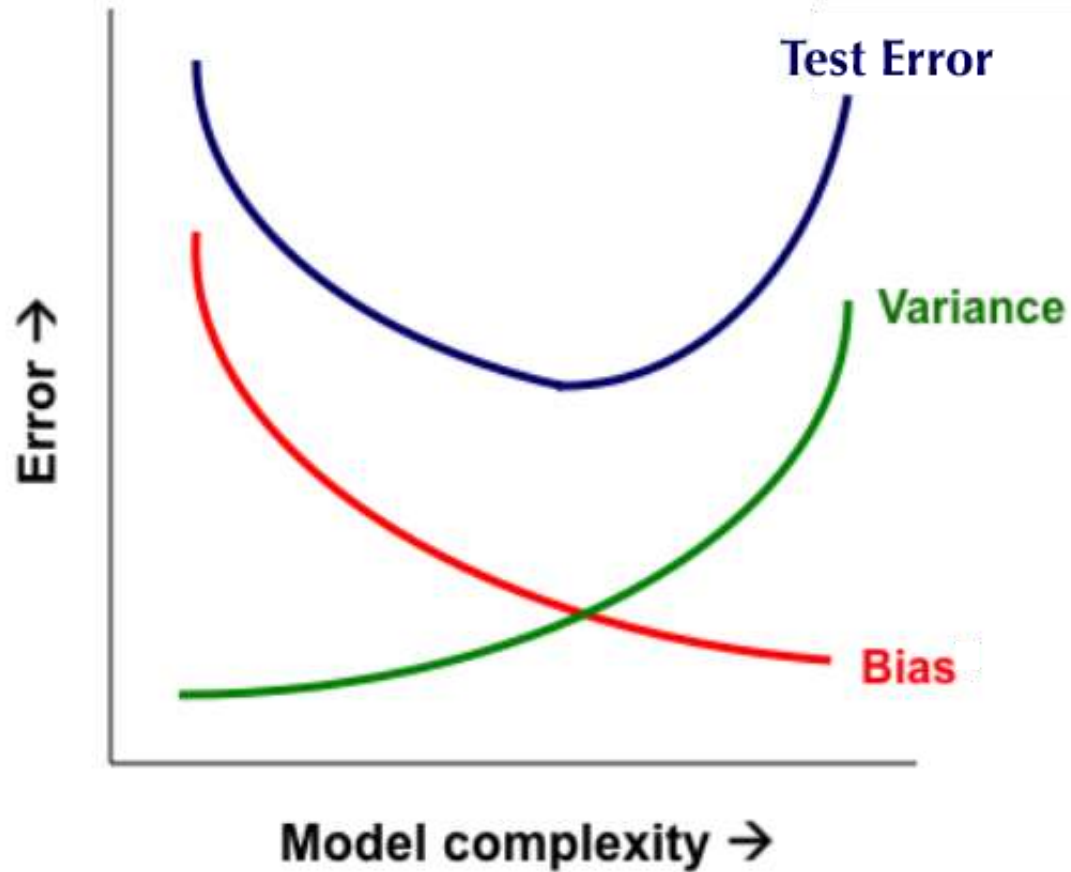
$$\begin{aligned} E_{out}(\mathbf{x}) &= (f(\mathbf{x}) - \bar{g}(\mathbf{x}))^2 + \mathbb{E}_{\mathcal{D}}[g_n^{\mathcal{D}}(\mathbf{x})^2] - \bar{g}(\mathbf{x})^2 \\ &= (f(\mathbf{x}) - f(\mathbf{x}))^2 + \mathbb{E}_{\mathcal{D}}[(g_n^{\mathcal{D}}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2] \\ &= 0 + \mathbb{E}_{\mathcal{D}}[(f(\mathbf{x}) + \varepsilon - f(\mathbf{x}))^2] \\ &= 0 + \mathbb{E}_{\mathcal{D}}[\varepsilon^2] \end{aligned}$$

In this case the bias term will be zero, but the variance is the square of the noise on the data, which could be substantial. In this case we say we have extreme over-fitting.



Fit every data point \Rightarrow
No approximation error (zero bias)
Variation between data sets (high variance)

Bias-Variance Trade-off





Vector/Matrix Calculus

- If $f(\mathbf{x}) \in \mathbb{R}$ and $\mathbf{x} \in \mathbb{R}^n$, then $\nabla_{\mathbf{x}} f \triangleq \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \cdots \quad \frac{\partial f}{\partial x_n} \right]^T$ is called the gradient of $f(\mathbf{x})$
- If $f(\mathbf{x}) = [f_1(x) \quad f_2(x) \quad \cdots \quad f_m(x)]^T \in \mathbb{R}^m$ and $\mathbf{x} \in \mathbb{R}^n$, then

$$\nabla_{\mathbf{x}} f = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}_{m \times n}$$

- is called "Jacobian matrix" of $f(\mathbf{x})$ w.r.t. \mathbf{x}



Vector/Matrix Calculus

- If $f(\mathbf{x}) \in \mathbb{R}$ and $\mathbf{x} \in \mathbb{R}^n$, then

- $$\nabla_{\mathbf{x}}^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 x_n} \\ \frac{\partial^2 f}{\partial x_1 x_2} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 x_n} & \frac{\partial^2 f}{\partial x_2 x_n} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

- is called the "Hessian matrix" of $f(\mathbf{x})$ w.r.t. \mathbf{x}



Vector/Matrix Calculus

- Some standard results:

- $\frac{\partial}{\partial \mathbf{x}} \mathbf{b}^T \mathbf{x} = \frac{\partial}{\partial \mathbf{x}} \mathbf{x}^T \mathbf{b} = \mathbf{b}$

- $\frac{\partial}{\partial \mathbf{x}} \mathbf{A} \mathbf{x} = \mathbf{A}$

- $\frac{\partial}{\partial \mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$; if $\mathbf{A} = \mathbf{A}^T$, $\frac{\partial}{\partial \mathbf{x}} \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{A} \mathbf{x}$

- Product rule: $\mathbf{u} \in \mathbb{R}^m$, $\mathbf{v} \in \mathbb{R}^m$, $\mathbf{x} \in \mathbb{R}^n$

- $$\left[\frac{\partial \mathbf{u}^T \mathbf{v}}{\partial \mathbf{x}} \right]_{1 \times n}^T = [\mathbf{u}^T]_{1 \times m} \left[\frac{\partial \mathbf{v}}{\partial \mathbf{x}} \right]_{m \times n} + [\mathbf{v}^T]_{1 \times m} \left[\frac{\partial \mathbf{u}}{\partial \mathbf{x}} \right]_{m \times n}$$



Vector/Matrix Calculus

- Some standard results:

- If $F(f(\mathbf{x})) \in \mathbb{R}$, $f(\mathbf{x}) \in \mathbb{R}$, $\mathbf{x} \in \mathbb{R}^n$, then

- $$\left[\frac{\partial F}{\partial \mathbf{x}} \right]_{n \times 1} = \left[\frac{\partial f}{\partial \mathbf{x}} \right]_{n \times 1} \left[\frac{\partial F}{\partial f} \right]_{1 \times 1}$$

- If $F(\mathbf{f}(\mathbf{x})) \in \mathbb{R}$, $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^m$, $\mathbf{x} \in \mathbb{R}^n$, then

- $$\left[\frac{\partial F}{\partial \mathbf{x}} \right]_{n \times 1} = \left[\frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right]_{n \times m}^T \left[\frac{\partial F}{\partial \mathbf{f}} \right]_{m \times 1}$$



Vector/Matrix Calculus

- Derivatives of norms:

- For two norm of a vector

$$\frac{\partial}{\partial \mathbf{x}} \|\mathbf{x}\|_2^2 = \frac{\partial}{\partial \mathbf{x}} \mathbf{x}^T \mathbf{x} = 2\mathbf{x}$$

- For Frobenius norm of a matrix

$$\frac{\partial}{\partial \mathbf{X}} \|\mathbf{X}\|_F^2 = \frac{\partial}{\partial \mathbf{X}} \text{tr}(\mathbf{X}\mathbf{X}^T) = 2\mathbf{X}$$



Optimization Problem

- Problem Statement:

$$\min_{s.t. \mathbf{x} \in \chi} f(\mathbf{x})$$

- Problem statement of convex optimization:

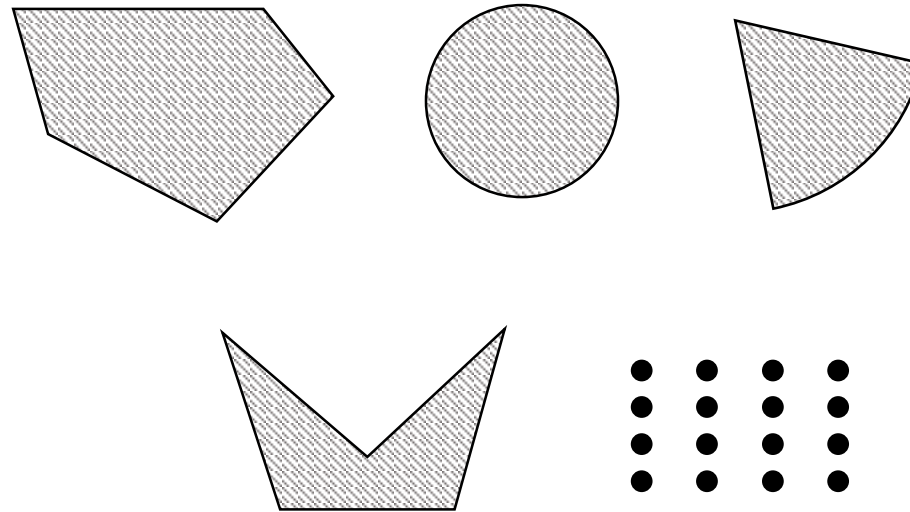
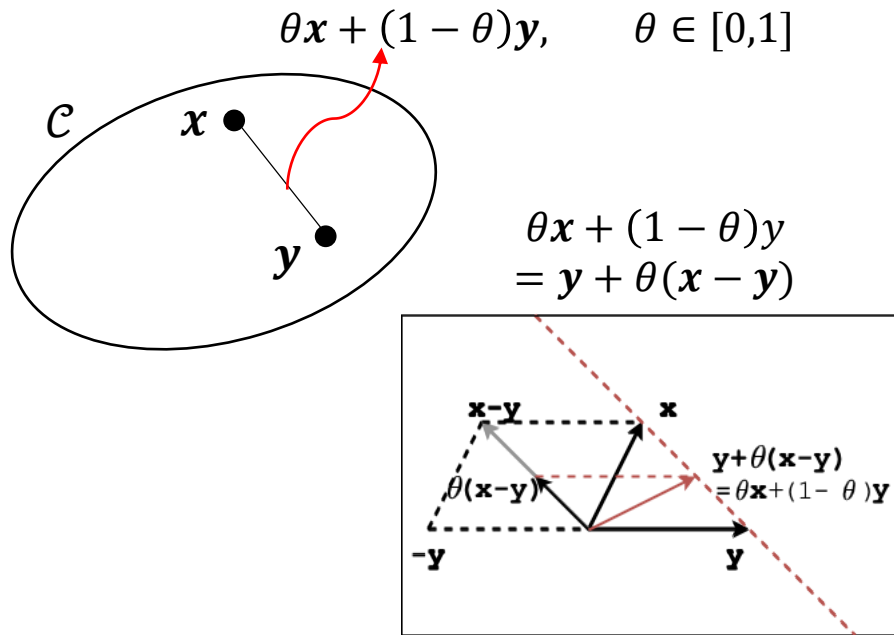
$$\min_{s.t. \mathbf{x} \in \chi} f(\mathbf{x})$$

- with $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a **convex function** and
- χ is a **convex set**

Convex Sets and Functions

- **Convex Set:** A set $\mathcal{C} \subseteq \mathbb{R}^n$ is a convex set if for all $x, y \in \mathcal{C}$, the line segment connecting x and y is in \mathcal{C} , i.e.,

$$\forall x, y \in \mathcal{C}, \theta \in [0, 1], \theta x + (1 - \theta)y \in \mathcal{C}$$





Convex Sets and Functions

- **Convex Function:** A function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex function if its domain $dom(f)$ is a convex set and for all $x, y \in dom(f)$, and $\theta \in [0,1]$, we have

$$f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y)$$

All norms are convex functions

$$\begin{aligned} \|\theta x + (1 - \theta)y\| &\leq \|\theta x\| + \|(1 - \theta)y\| \\ &= \theta \|x\| + (1 - \theta)\|y\| \end{aligned}$$





Alternative Definition of Convexity for differentiable functions

Theorem: (first order characterization) Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ be a differentiable function whose $\text{dom}(f)$ is convex. Then, f is convex *iff*

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f) \quad \dots(1)$$



Proof (part 1): $Eq^n(1) \Rightarrow \text{Convex}$

Given $Eq^n(1)$: $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$

Let us consider $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$ and $\theta \in [0,1]$

Let $\mathbf{z} = (1 - \theta)\mathbf{x} + \theta\mathbf{y}$

Now, By $eq^n(1)$,

$$f(\mathbf{x}) \geq f(\mathbf{z}) + \langle \nabla f(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle$$

... (2), taking \mathbf{x}, \mathbf{z}

$$f(\mathbf{y}) \geq f(\mathbf{z}) + \langle \nabla f(\mathbf{z}), \mathbf{y} - \mathbf{z} \rangle$$

... (3), taking \mathbf{y}, \mathbf{z}



Multiplying Eq^n (2) with $(1 - \theta)$, we get,

$$(1 - \theta)f(\mathbf{x}) \geq (1 - \theta)f(\mathbf{z}) + (1 - \theta)\langle \nabla f(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle \quad \dots(4)$$

and Eq^n (3) with θ

$$\theta f(\mathbf{y}) \geq \theta f(\mathbf{z}) + \theta \langle \nabla f(\mathbf{z}), \mathbf{y} - \mathbf{z} \rangle \quad \dots (5)$$

Now combining Eq^n s (4) and (5), we have,

$$(1 - \theta)f(\mathbf{x}) + \theta f(\mathbf{y}) \geq (1 - \theta)f(\mathbf{z}) + \theta f(\mathbf{z}) + (1 - \theta)\langle \nabla f(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle + \theta \langle \nabla f(\mathbf{z}), \mathbf{y} - \mathbf{z} \rangle$$



$$\begin{aligned}(1 - \theta)f(\mathbf{x}) + \theta f(\mathbf{y}) &\geq (1 - \theta)f(\mathbf{z}) + \theta f(\mathbf{z}) + (1 - \theta)\langle \nabla f(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle \\ &\quad + \theta \langle \nabla f(\mathbf{z}), \mathbf{y} - \mathbf{z} \rangle \\ \Rightarrow (1 - \theta)f(\mathbf{x}) + \theta f(\mathbf{y}) &\geq f(\mathbf{z}) - \cancel{\theta f(\mathbf{z})} + \cancel{\theta f(\mathbf{z})} + \langle \nabla f(\mathbf{z}), (1 - \theta)(\mathbf{x} - \mathbf{z}) \rangle \\ &\quad + \langle \nabla f(\mathbf{z}), \theta(\mathbf{y} - \mathbf{z}) \rangle \\ \Rightarrow (1 - \theta)f(\mathbf{x}) + \theta f(\mathbf{y}) &\geq f(\mathbf{z}) + \underbrace{\langle \nabla f(\mathbf{z}), (1 - \theta)(\mathbf{x} - \mathbf{z}) + \theta(\mathbf{y} - \mathbf{z}) \rangle}_{=0 \text{ (Why?)}}\end{aligned}$$

$$\begin{aligned}&\left[\begin{aligned}(1 - \theta)(\mathbf{x} - \mathbf{z}) + \theta(\mathbf{y} - \mathbf{z}) &= (1 - \theta)\mathbf{x} - (1 - \theta)\mathbf{z} + \theta\mathbf{y} - \theta\mathbf{z} \\ &= (1 - \theta)\mathbf{x} - \mathbf{z} + \theta\mathbf{z} + \theta\mathbf{y} - \theta\mathbf{z} \\ &= (1 - \theta)\mathbf{x} + \theta\mathbf{y} - \mathbf{z} = \mathbf{z} - \mathbf{z} = \mathbf{0}\end{aligned} \right] \\ &= f(\mathbf{z}) = f((1 - \theta)\mathbf{x} + \theta\mathbf{y})\end{aligned}$$



Proof (part 2): *Convex* \Rightarrow Eqⁿ (1)

Suppose f is convex, i.e., $\forall \mathbf{x}, \mathbf{y} \in \text{dom}(f)$ and $\forall \theta \in [0,1]$,

$$f((1 - \theta)\mathbf{x} + \theta\mathbf{y}) \leq (1 - \theta)f(\mathbf{x}) + \theta f(\mathbf{y})$$

Equivalently,

$$f(\mathbf{x} + \theta(\mathbf{y} - \mathbf{x})) \leq f(\mathbf{x}) + \theta(f(\mathbf{y}) - f(\mathbf{x}))$$

$$\Rightarrow f(\mathbf{y}) - f(\mathbf{x}) \geq \frac{f(\mathbf{x} + \theta(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\theta}$$



By taking the limit as $\theta \rightarrow 0$ on both sides and using the definition of derivative, we obtain,

$$f(\mathbf{y}) - f(\mathbf{x}) \geq \lim_{\theta \rightarrow 0} \frac{f(\mathbf{x} + \theta(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\theta} = \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$$

(How?)

$$g = \lim_{\theta \rightarrow 0} \frac{f(x + \theta(y - x)) - f(x)}{\theta} \quad [\text{Let us see in case of 1-D } x, y.]$$

$$g = \lim_{\theta \rightarrow 0} \frac{f(x + \theta(y - x)) - f(x)}{\theta(y - x)} (y - x) \quad [\text{Multiplying numerator and denominator by } (y - x)]$$

Let, $\theta(y - x) = h$, As $\theta \rightarrow 0$, $h \rightarrow 0$

$$\therefore g = \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h} (y - x) = f'(x)(y - x)$$

Gradient Descent

min loss function
 w, b

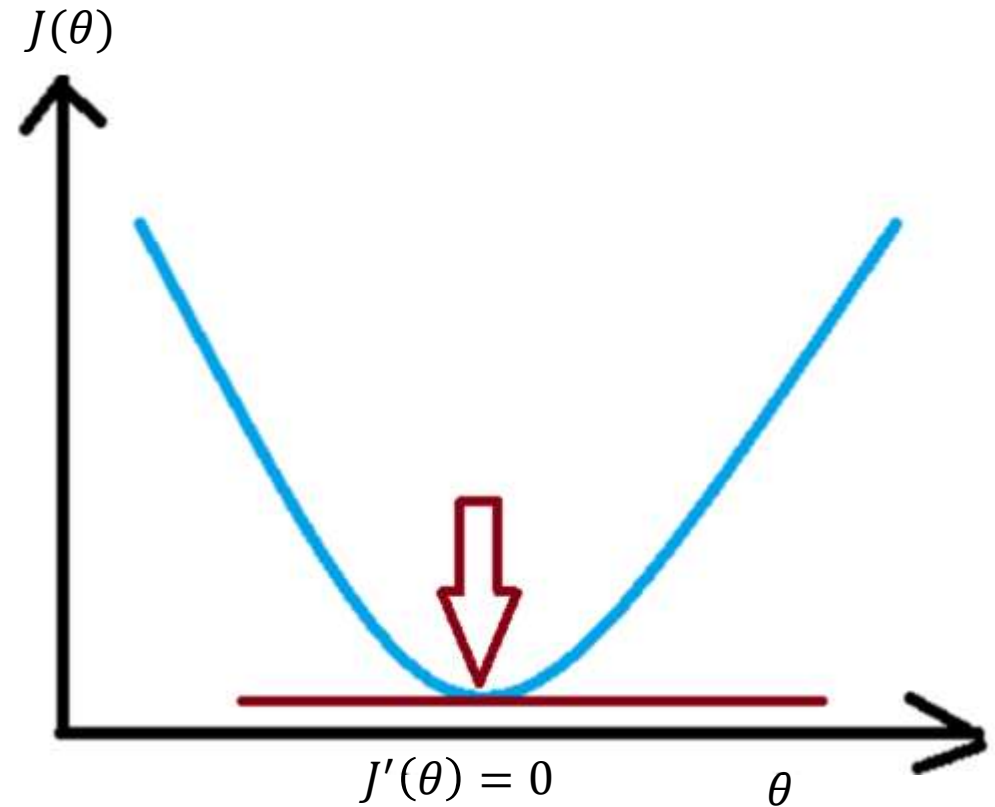
More formally, $\min_{\theta} J(\theta)$

For scalar θ , the condition is $J'(\theta) = \frac{\partial J}{\partial \theta} = 0$

For higher dimensional θ , the condition boils down to
 $J'(\theta) = \nabla_{\theta} J = \mathbf{0}$

$$\Rightarrow \left[\frac{\partial J}{\partial \theta_1}, \frac{\partial J}{\partial \theta_2}, \dots, \frac{\partial J}{\partial \theta_N} \right]^T = \mathbf{0}$$

$$\Rightarrow \begin{bmatrix} \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \\ \vdots \\ \frac{\partial J}{\partial \theta_N} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$



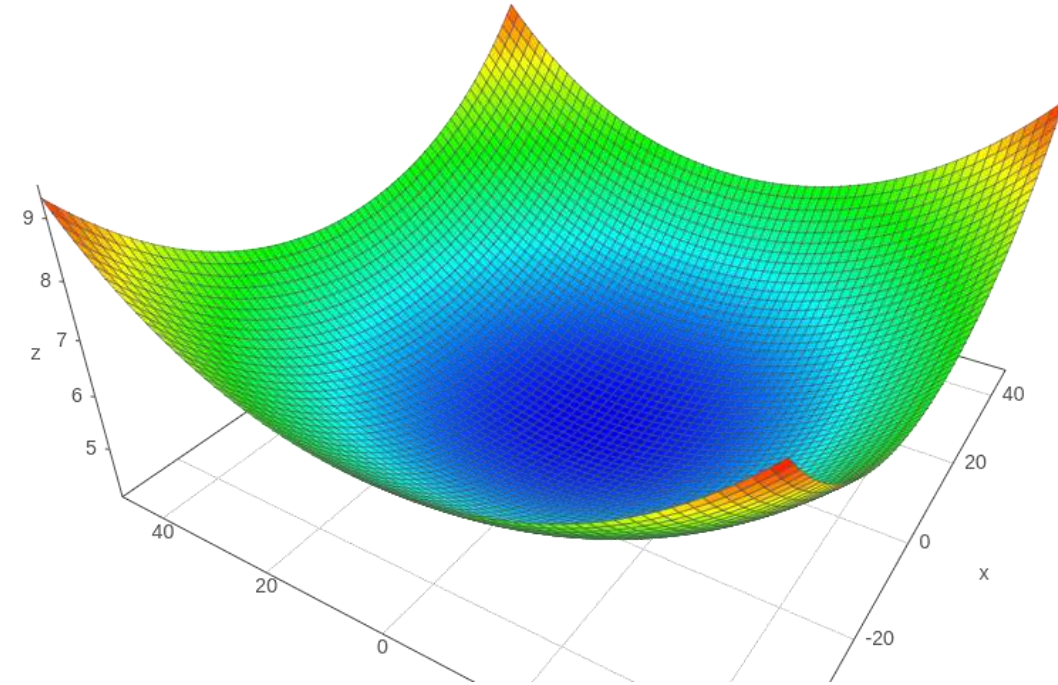
One Way to Find Minima – Gradient Descent

This is helpful but not always useful.

For example $J(\theta) = \log \sum_{i=1}^m e^{(\mathbf{a}_i^T \theta + b_i)}$ is a convex function with clear minima, but finding analytical solution is not easy.

$$\begin{aligned} \nabla_{\theta} J &= \mathbf{0} \\ \Rightarrow \frac{1}{\sum_{i=1}^m e^{(\mathbf{a}_i^T \theta + b_i)}} \sum_{i=1}^m e^{(\mathbf{a}_i^T \theta + b_i)} \mathbf{a}_i &= \mathbf{0} \\ \Rightarrow \sum_{i=1}^m e^{(\mathbf{a}_i^T \theta + b_i)} \mathbf{a}_i &= \mathbf{0} \end{aligned}$$

So, a numerical iterative solution is sought for.





One Way to Find Minima – Gradient Descent

Start with an initial guess θ^0

Repeatedly update θ by taking a small step: $\theta^k = \theta^{k-1} + \eta \Delta \theta \dots (1)$

so that $J(\theta)$ gets smaller with each update i.e., $J(\theta^k) \leq J(\theta^{k-1}) \dots (2)$

$$\begin{aligned} (1) \text{ implies, } J(\theta^k) &= J(\theta^{k-1} + \eta \Delta \theta) \\ &= J(\theta^{k-1}) + \eta (\Delta \theta)^T J'(\theta^{k-1}) + h. o. t. \text{ [Using Taylor series expansion]} \\ &\approx J(\theta^{k-1}) + \eta (\Delta \theta)^T J'(\theta^{k-1}) \text{ [Neglecting h. o. t.] } \dots (3) \end{aligned}$$

Combining (2) and (3),

$$\eta (\Delta \theta)^T J'(\theta^{k-1}) \leq 0 \text{ i. e., } (\Delta \theta)^T J'(\theta^{k-1}) \leq 0 \text{ [as } \eta \text{ is positive]}$$

So, for θ to minimize $J(\theta)$, i. e., to satisfy (2) we have to choose some $\Delta \theta$ that is negative when a dot product of it is done with the gradient $J'(\theta^{k-1})$.

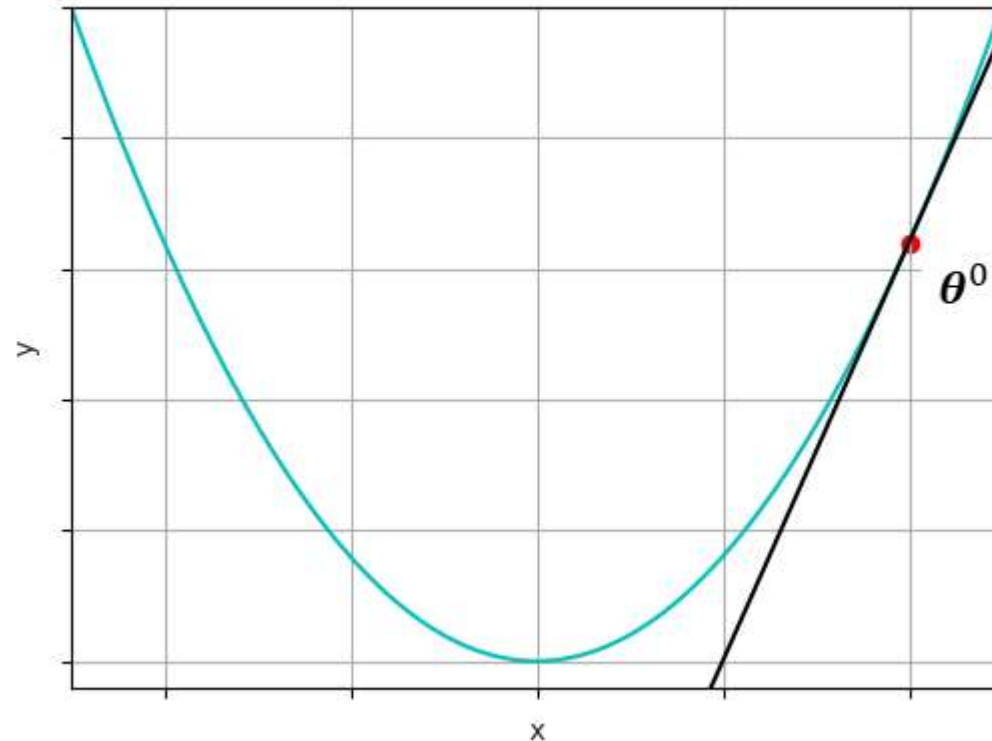
Then why not choose, $\Delta \theta = -J'(\theta^{k-1})$!!

Then, $(\Delta \theta)^T J'(\theta^{k-1}) = -||J'(\theta^{k-1})||^2$ surely is a negative quantity and satisfies the condition.

Gradient Descent

Start with an initial guess θ^0

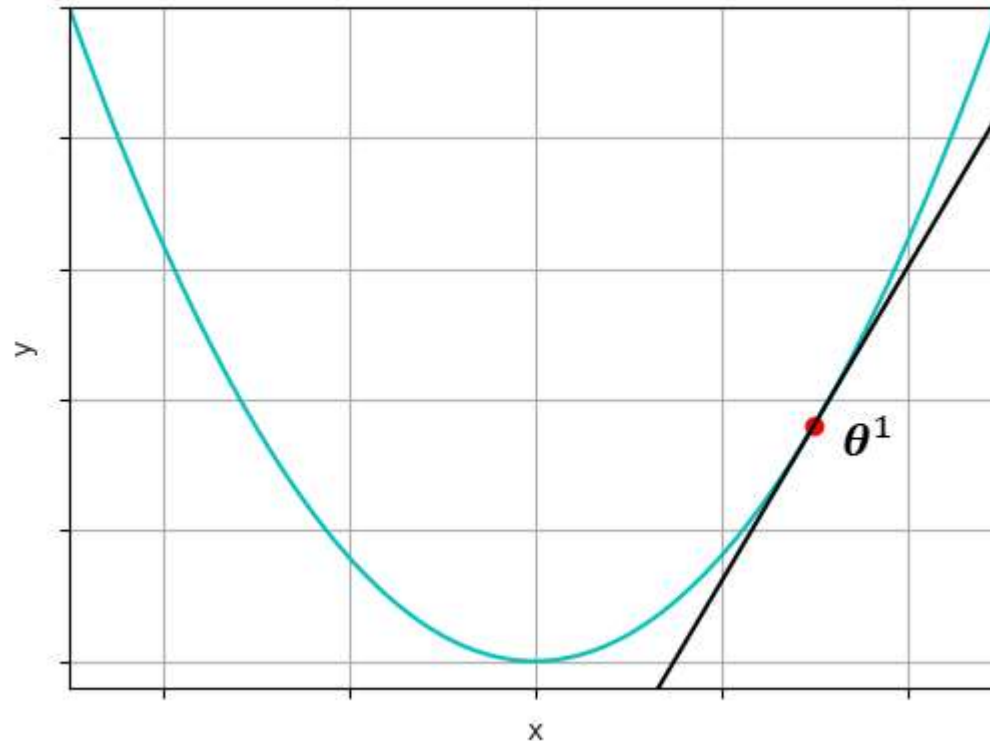
Repeatedly update θ by taking a small step: $\theta^k = \theta^{k-1} - \eta J'(\theta^{k-1})$ until convergence ($J'(\theta^{k-1})$ is very small)



Gradient Descent

Start with an initial guess θ^0

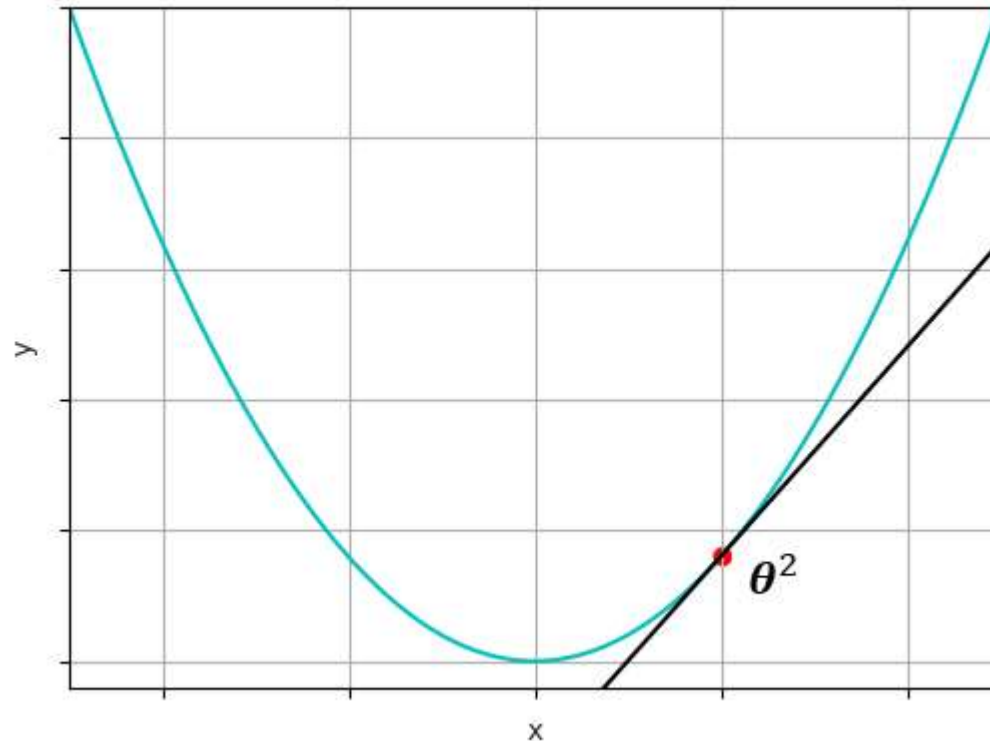
Repeatedly update θ by taking a small step: $\theta^k = \theta^{k-1} - \eta J'(\theta^{k-1})$ until convergence ($J'(\theta^{k-1})$ is very small)



Gradient Descent

Start with an initial guess θ^0

Repeatedly update θ by taking a small step: $\theta^k = \theta^{k-1} - \eta J'(\theta^{k-1})$ until convergence ($J'(\theta^{k-1})$ is very small)

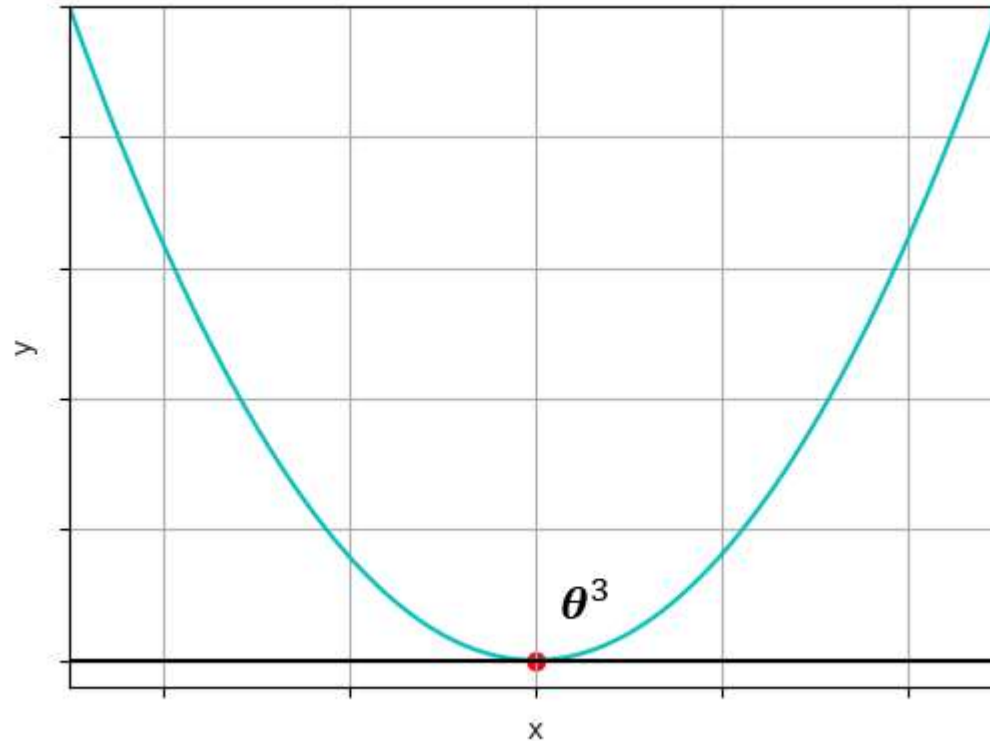


Gradient Descent

Start with an initial guess θ^0

Repeatedly update θ by taking a small step: $\theta^k = \theta^{k-1} - \eta J'(\theta^{k-1})$ until convergence ($J'(\theta^{k-1})$ is very small)

Remember, in Neural Networks,
the loss is computed by averaging
the losses for all examples



But Imagine This

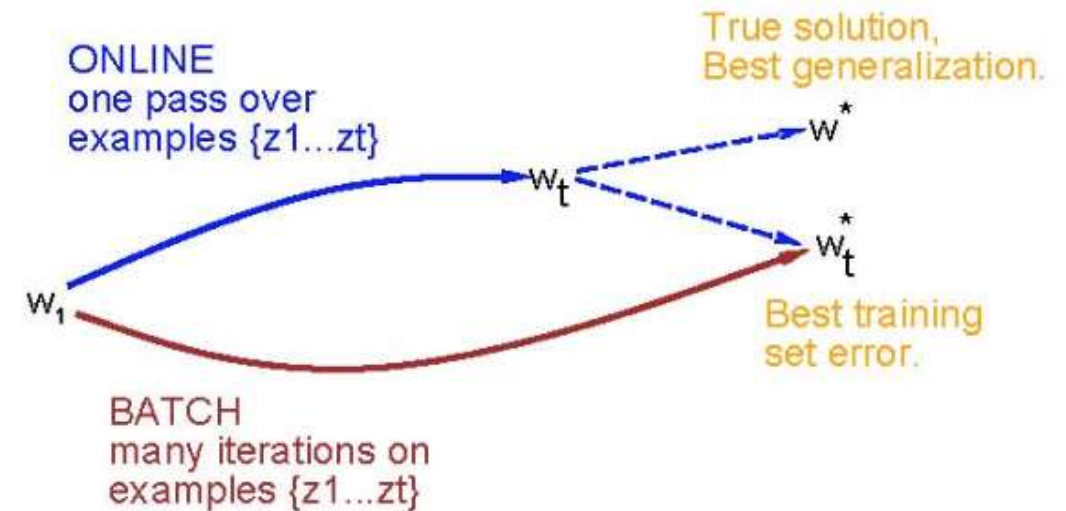
Speed Advantage.

Assume training set contains
10 copies of the 100 same examples.

- **Batch** Blindly computes redundant gradients.
1 epoch on large set \equiv 1 epochs on small set.
- **Online** Take advantage of redundancy.
1 epoch on large set \equiv 10 epochs on small set.

In practice, stochastic gradient
can be orders of magnitude faster.

Online vs. Batch



Adapted from Olivier Bousquet's invited talk at NeurIPS 2018 after winning Test of Time Award for NIPS 2007 paper: "The Trade-Offs of Large Scale Learning" by Leon Bottou and Olivier Bousquet. Link of the [talk](#).



Batch, Stochastic and Minibatch

- Optimization algorithms that use the entire training set to compute the gradient are called batch or deterministic gradient methods. Ones that use a single training example for that task are called stochastic or online gradient methods
- Most of the algorithms we use for deep learning fall somewhere in between!
- These are called minibatch or minibatch stochastic methods

Batch, Stochastic and Mini-batch Stochastic Gradient Descent

Algorithm 1 Batch Gradient Descent at Iteration k

Require: Learning rate ϵ_k

Require: Initial Parameter θ

- 1: **while** stopping criteria not met **do**
- 2: Compute gradient estimate over N examples:
- 3: $\hat{\mathbf{g}} \leftarrow +\frac{1}{N} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$
- 4: Apply Update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$
- 5: **end while**

Algorithm 2 Stochastic Gradient Descent at Iteration k

Require: Learning rate ϵ_k

Require: Initial Parameter θ

- 1: **while** stopping criteria not met **do**
- 2: Sample example $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$ from training set
- 3: Compute gradient estimate:
- 4: $\hat{\mathbf{g}} \leftarrow +\nabla_{\theta} L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$
- 5: Apply Update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$
- 6: **end while**

Mini-batch

Algorithm 8.1 Stochastic gradient descent (SGD) update at training iteration k

Require: Learning rate ϵ_k .

Require: Initial parameter θ

while stopping criterion not met **do**

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

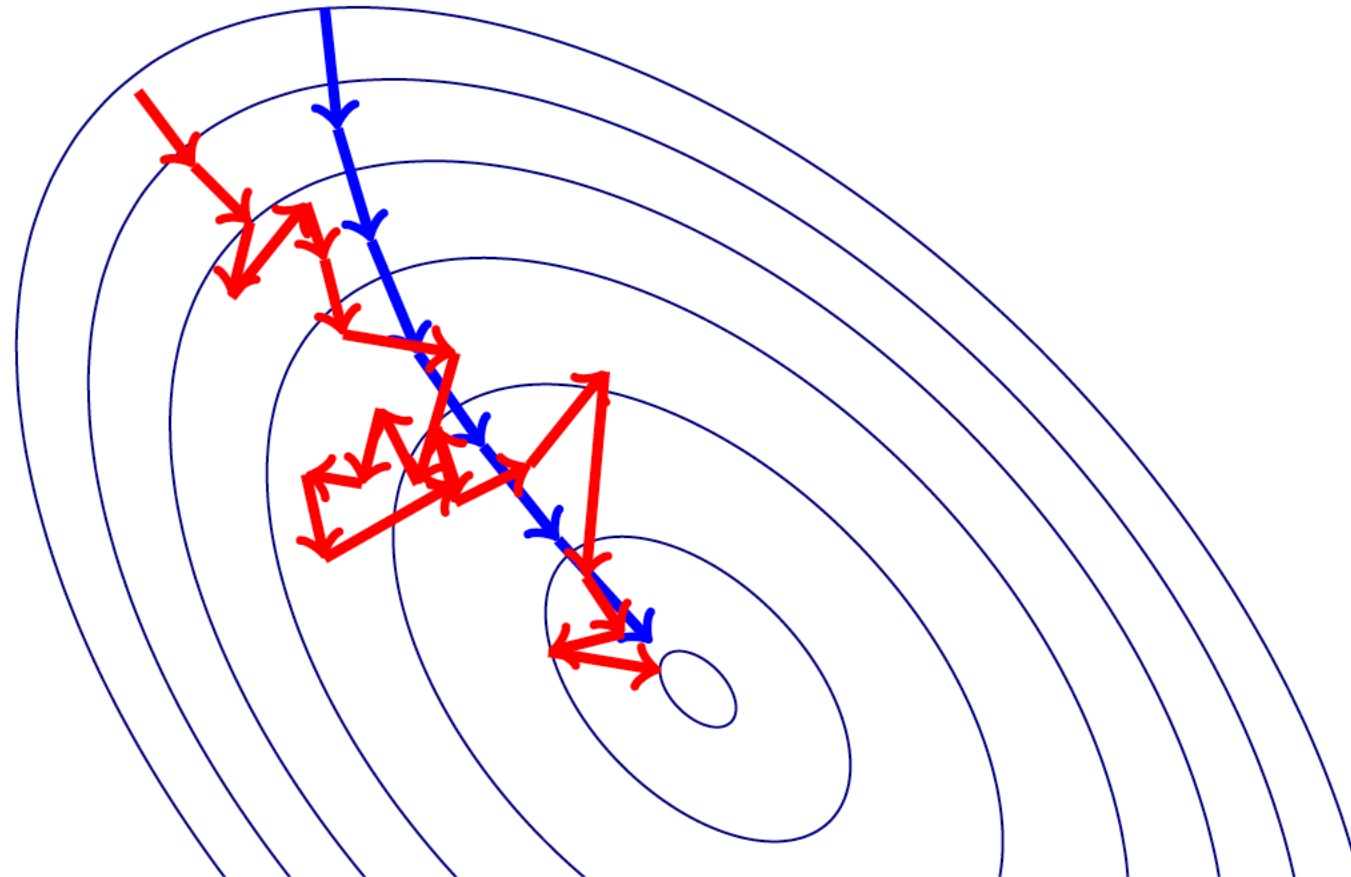
 Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Apply update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$

end while

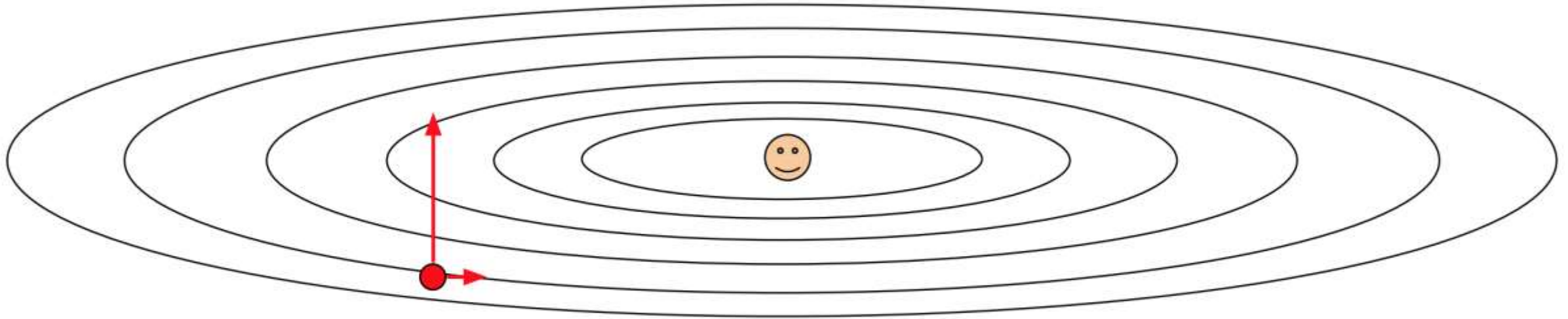
Images courtesy: Shubhendu Trivedi et. al, Goodfellow et. al..

Batch and Stochastic Gradient Descent



Images courtesy: Shubhendu Trivedi et. al.

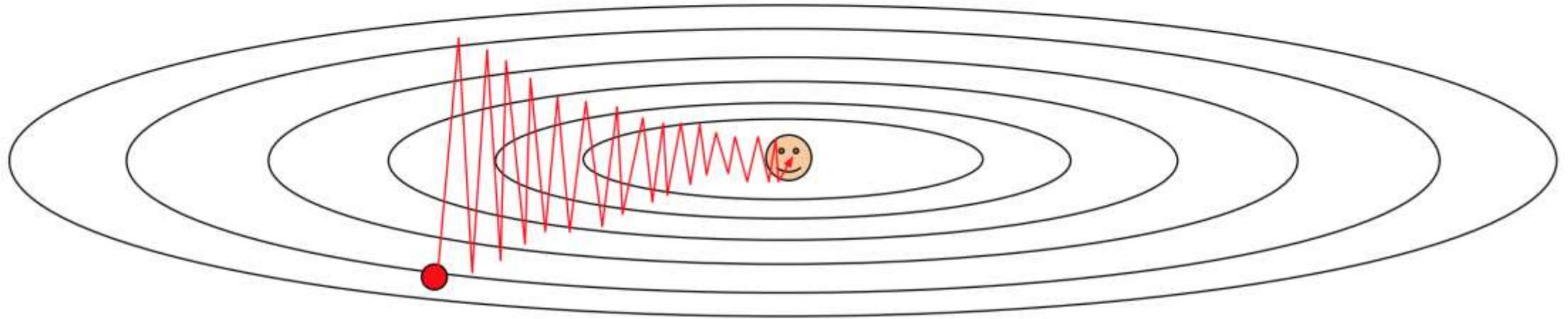
Suppose loss function is steep vertically but shallow horizontally:



Q: What is the trajectory along which we converge towards the minimum with SGD?

Images courtesy: Karpathy et. al.

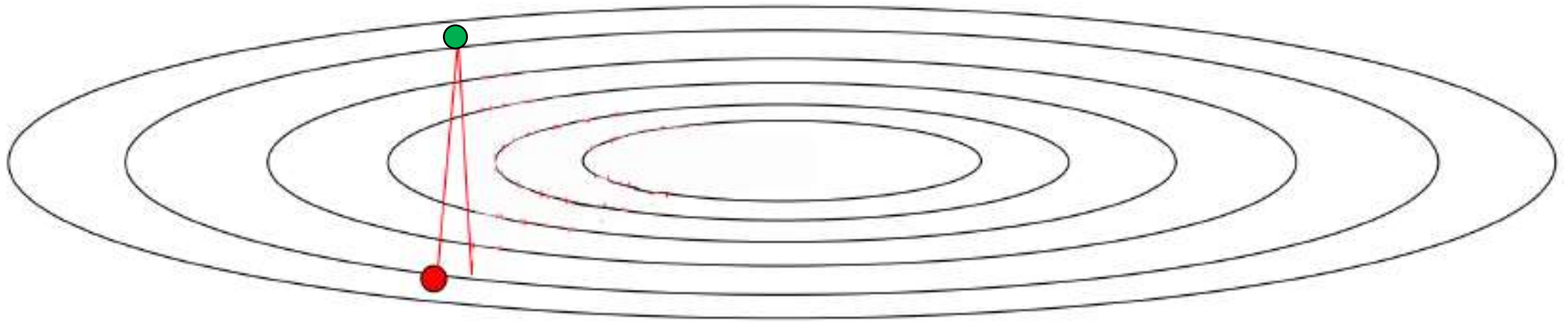
Suppose loss function is steep vertically but shallow horizontally:



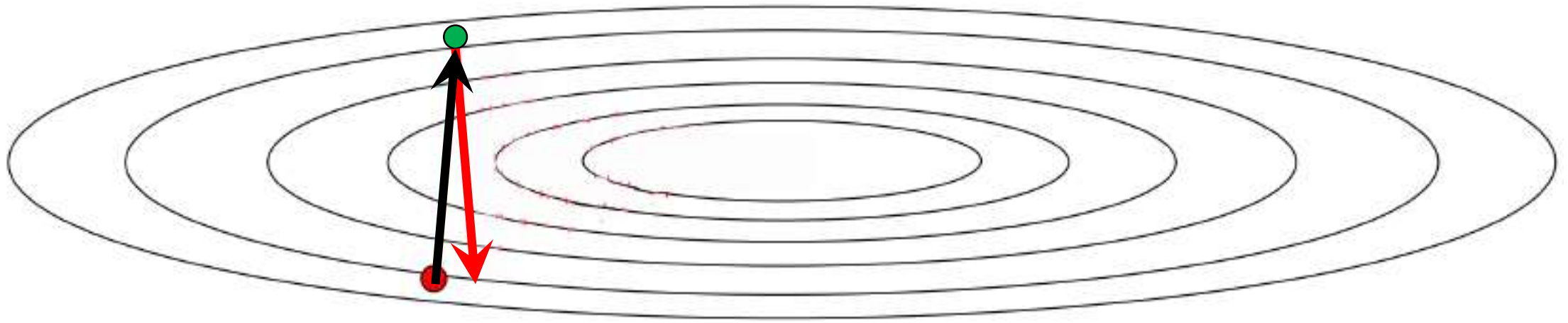
Q: What is the trajectory along which we converge towards the minimum with SGD? **very slow progress along flat direction, jitter along steep one**

Images courtesy: Karpathy et. al.

Suppose loss function is steep vertically but shallow horizontally:

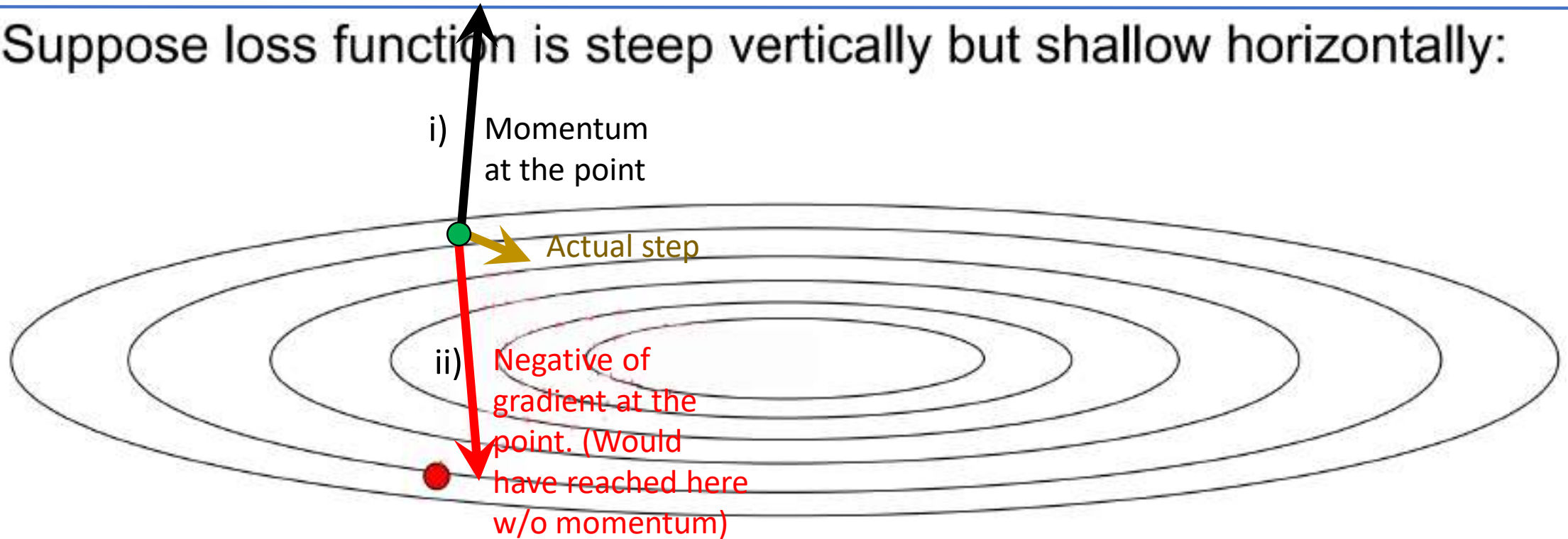


Suppose loss function is steep vertically but shallow horizontally:



Gradient Descent: only one force acting at any point.

Suppose loss function is steep vertically but shallow horizontally:



Momentum: Two forces acting at any point.

- i) Momentum built up due to gradients pushing the particle at that point
- ii) Gradient computed at that point

Stochastic Gradient Descent with Momentum

Algorithm 8.2 Stochastic gradient descent (SGD) with momentum

Require: Learning rate ϵ , momentum parameter α .

Require: Initial parameter θ , initial velocity v .

while stopping criterion not met **do**

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

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

 Compute velocity update: $v \leftarrow \alpha v - \epsilon g$

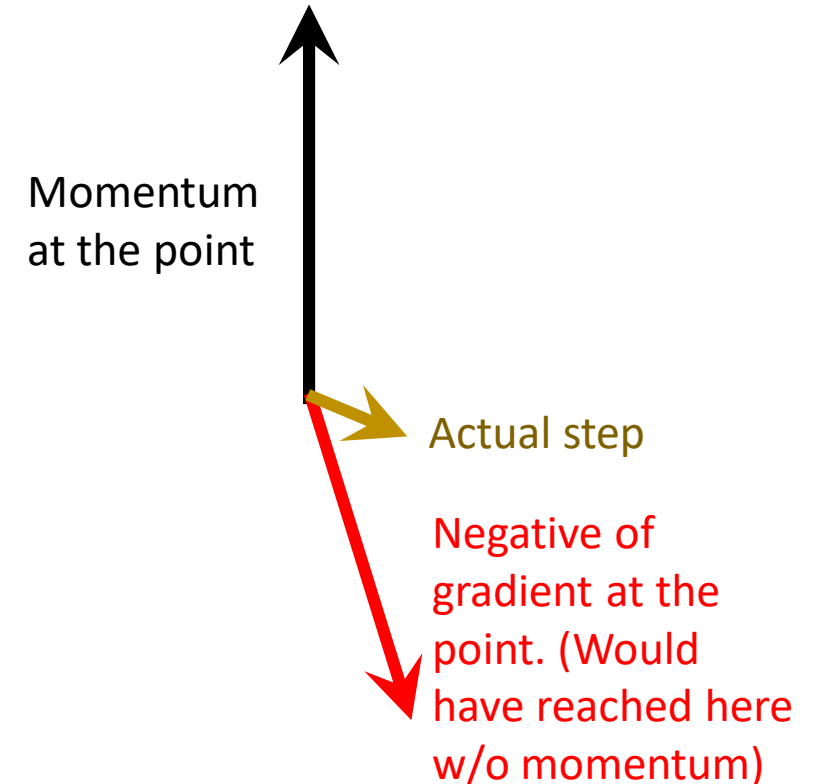
 Apply update: $\theta \leftarrow \theta + v$

end while

$\alpha \in [0, 1)$

What is the role of α ?

- If α is larger than ϵ the current update is more affected by the previous gradients
- Usually values for α are set high $\approx 0.8, 0.9$



Images courtesy: Goodfellow et. al.

Nesterov Momentum

Algorithm 8.3 Stochastic gradient descent (SGD) with Nesterov momentum

Require: Learning rate ϵ , momentum parameter α .

Require: Initial parameter θ , initial velocity v .

while stopping criterion not met **do**

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

 Apply interim update: $\tilde{\theta} \leftarrow \theta + \alpha v$

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

 Compute velocity update: $v \leftarrow \alpha v - \epsilon g$

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

end while

- The difference between Nesterov momentum and standard momentum is where the gradient is evaluated.
- With Nesterov momentum the gradient is evaluated after the current velocity is applied.
- In practice, Nesterov momentum speeds up the convergence only for well behaved loss functions (convex with consistent curvature)

Images courtesy: Goodfellow et. al.

Distill Pub Article on Why Momentum Really Works

Why Momentum Really Works



Step-size $\alpha = 0.02$



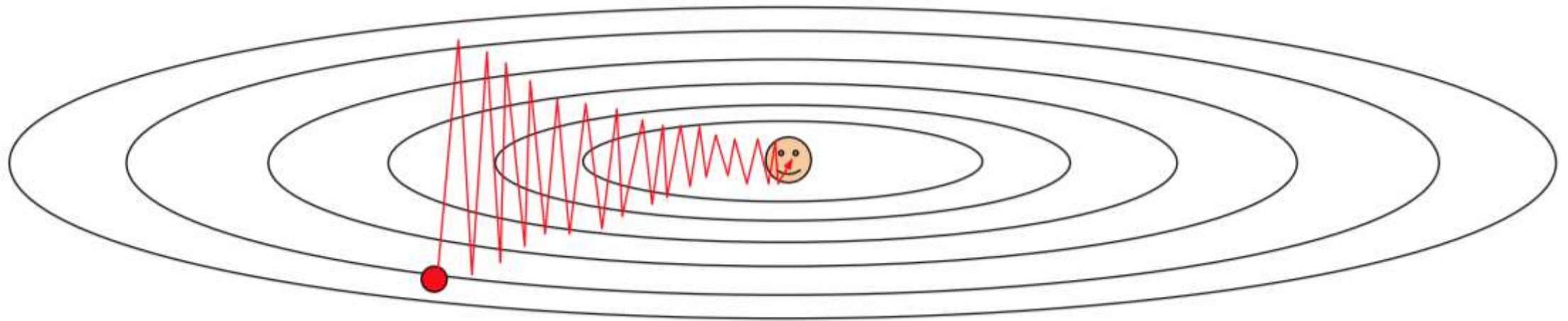
Momentum $\beta = 0.85$



We often think of Momentum as a means of dampening oscillations and speeding up the iterations, leading to faster convergence. But it has other interesting behavior. It allows a larger range of step-sizes to be used, and creates its own oscillations. What is going on?

Adaptive Learning Rate Methods

- Till now we assign the same learning rate in all directions.
- Would it not be a good idea if we can move slowly in a steeper direction whereas move fast in a shallower direction?



Images courtesy: Karpathy et. al.

AdaGrad

- Downscale learning rate by square-root of sum of squares of all the historical gradient values
- Parameters that have large partial derivative of the loss – learning rates for them are rapidly declined

Algorithm 4 AdaGrad

Require: Global Learning rate ϵ , Initial Parameter θ , δ

Initialize $\mathbf{r} = 0$

- 1: **while** stopping criteria not met **do**
 - 2: Sample example $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$ from training set
 - 3: Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow +\nabla_{\theta} L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$
 - 4: Accumulate: $\mathbf{r} \leftarrow \mathbf{r} + \hat{\mathbf{g}} \odot \hat{\mathbf{g}}$
 - 5: Compute update: $\Delta\theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \hat{\mathbf{g}}$
 - 6: Apply Update: $\theta \leftarrow \theta + \Delta\theta$
 - 7: **end while**
-

AdaGrad

$$\bullet \hat{\mathbf{g}}_{(1)} = \begin{bmatrix} g_{(1),1} \\ g_{(1),2} \\ \vdots \\ g_{(1),d} \end{bmatrix}, \hat{\mathbf{g}}_{(2)} = \begin{bmatrix} g_{(2),1} \\ g_{(2),2} \\ \vdots \\ g_{(2),d} \end{bmatrix}, \dots, \hat{\mathbf{g}}_{(t)} = \begin{bmatrix} g_{(t),1} \\ g_{(t),2} \\ \vdots \\ g_{(t),d} \end{bmatrix}$$

$$\bullet \mathbf{r}_{(1)} = \begin{bmatrix} g_{(1),1}^2 \\ g_{(1),2}^2 \\ \vdots \\ g_{(1),d}^2 \end{bmatrix}, \mathbf{r}_{(2)} = \begin{bmatrix} g_{(1),1}^2 + g_{(2),1}^2 \\ g_{(1),2}^2 + g_{(2),2}^2 \\ \vdots \\ g_{(1),d}^2 + g_{(2),d}^2 \end{bmatrix}, \dots, \mathbf{r}_{(t)} = \begin{bmatrix} g_{(1),1}^2 + g_{(2),1}^2 + \dots + g_{(t),1}^2 \\ g_{(1),2}^2 + g_{(2),2}^2 + \dots + g_{(t),2}^2 \\ \vdots \\ g_{(1),d}^2 + g_{(2),d}^2 + \dots + g_{(t),d}^2 \end{bmatrix}$$

$$\bullet \Delta \boldsymbol{\theta}_{(t)} = - \begin{bmatrix} \frac{\epsilon \cdot g_{(t),1}}{\delta + \sqrt{g_{(1),1}^2 + g_{(2),1}^2 + \dots + g_{(t),1}^2}} \\ \frac{\epsilon \cdot g_{(t),2}}{\delta + \sqrt{g_{(1),2}^2 + g_{(2),2}^2 + \dots + g_{(t),2}^2}} \\ \vdots \\ \frac{\epsilon \cdot g_{(t),d}}{\delta + \sqrt{g_{(1),d}^2 + g_{(2),d}^2 + \dots + g_{(t),d}^2}} \end{bmatrix}$$

RMSProp

- One problem of AdaGrad is that the 'r' vector continues to build up and grow its value.
- This shrinks the learning rate too aggressively.
- RMSProp strikes a balance by exponentially decaying contributions from past gradients.

Algorithm 5 RMSProp

Require: Global Learning rate ϵ , decay parameter ρ, δ

Initialize $\mathbf{r} = 0$

- 1: **while** stopping criteria not met **do**
 - 2: Sample example $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$ from training set
 - 3: Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow +\nabla_{\theta} L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$
 - 4: Accumulate: $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \hat{\mathbf{g}} \odot \hat{\mathbf{g}}$
 - 5: Compute update: $\Delta \theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \hat{\mathbf{g}}$
 - 6: Apply Update: $\theta \leftarrow \theta + \Delta \theta$
 - 7: **end while**
-

RMSProp

- $\hat{\mathbf{g}}_{(1)}, \hat{\mathbf{g}}_{(2)}, \dots, \hat{\mathbf{g}}_{(t)} \quad \mathbf{r}_{(t)} = \rho \mathbf{r}_{(t-1)} + (1 - \rho) \hat{\mathbf{g}}_{(t)} \odot \hat{\mathbf{g}}_{(t)}$
 - $\mathbf{r}_{(0)} = 0$
 - $\mathbf{r}_{(1)} = \rho \mathbf{r}_{(0)} + (1 - \rho) \hat{\mathbf{g}}_{(1)} \odot \hat{\mathbf{g}}_{(1)} = (1 - \rho) \hat{\mathbf{g}}_{(1)} \odot \hat{\mathbf{g}}_{(1)}$
 - $\mathbf{r}_{(2)} = \rho \mathbf{r}_{(1)} + (1 - \rho) \hat{\mathbf{g}}_{(2)} \odot \hat{\mathbf{g}}_{(2)} = \rho(1 - \rho) \hat{\mathbf{g}}_{(1)} \odot \hat{\mathbf{g}}_{(1)} + (1 - \rho) \hat{\mathbf{g}}_{(2)} \odot \hat{\mathbf{g}}_{(2)}$
 - $\mathbf{r}_{(3)} = \rho \mathbf{r}_{(2)} + (1 - \rho) \hat{\mathbf{g}}_{(3)} \odot \hat{\mathbf{g}}_{(3)} = \rho^2(1 - \rho) \hat{\mathbf{g}}_{(1)} \odot \hat{\mathbf{g}}_{(1)} + \rho(1 - \rho) \hat{\mathbf{g}}_{(2)} \odot \hat{\mathbf{g}}_{(2)} + (1 - \rho) \hat{\mathbf{g}}_{(3)} \odot \hat{\mathbf{g}}_{(3)}$
-
- RMSProp uses an exponentially decaying average to discard history from the extreme past so that the accumulation of gradients do not stall the learning.
 - AdaDelta is another variant where instead of exponentially decaying average, a moving window average of the past gradients is taken
 - Nesterov acceleration can also be applied to both these variants by computing the gradients at a 'look ahead' position (i.e., at a place where the momentum would have taken the parameters).



ADAM (Adaptive Moments)

- Variant of the combination of RMSProp and Momentum.
- Incorporates first order moment of the gradient which can be thought of as equivalent to taking advantage of the momentum strategy. Here momentum is also added with exponential averaging.
- It also incorporates the second order term which can be thought of as the RMSProp like exponential averaging of the past gradients.
- Both first and second moments are corrected for bias to account for their initialization to zero.



ADAM (Adaptive Moments)

- Biased first order moment $\mathbf{s}_{(t)} = \rho_1 \mathbf{s}_{(t-1)} + (1 - \rho_1) \mathbf{g}_t$
- Biased second order moment $\mathbf{r}_{(t)} = \rho_2 \mathbf{r}_{(t-1)} + (1 - \rho_2) \mathbf{g}_t \odot \mathbf{g}_t$
- Bias corrected first order moment $\hat{\mathbf{s}}_{(t)} = \frac{\mathbf{s}_{(t)}}{1 - \rho_1}$
- Bias corrected second order moment $\hat{\mathbf{r}}_{(t)} = \frac{\mathbf{r}_{(t)}}{1 - \rho_2}$
- Weight update $\Delta \boldsymbol{\theta}_{(t)} = -\epsilon \frac{\hat{\mathbf{s}}_{(t)}}{\sqrt{\hat{\mathbf{r}}_{(t)} + \delta}}$ (operations applied elementwise)

ADAM (Adaptive Moments)

Algorithm 8.7 The Adam algorithm

Require: Step size ϵ (Suggested default: 0.001)

Require: Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1)$.
(Suggested defaults: 0.9 and 0.999 respectively)

Require: Small constant δ used for numerical stabilization. (Suggested default: 10^{-8})

Require: Initial parameters θ

Initialize 1st and 2nd moment variables $\mathbf{s} = \mathbf{0}$, $\mathbf{r} = \mathbf{0}$

Initialize time step $t = 0$

while stopping criterion not met **do**

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

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

$t \leftarrow t + 1$

 Update biased first moment estimate: $\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$

 Update biased second moment estimate: $\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$

 Correct bias in first moment: $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$

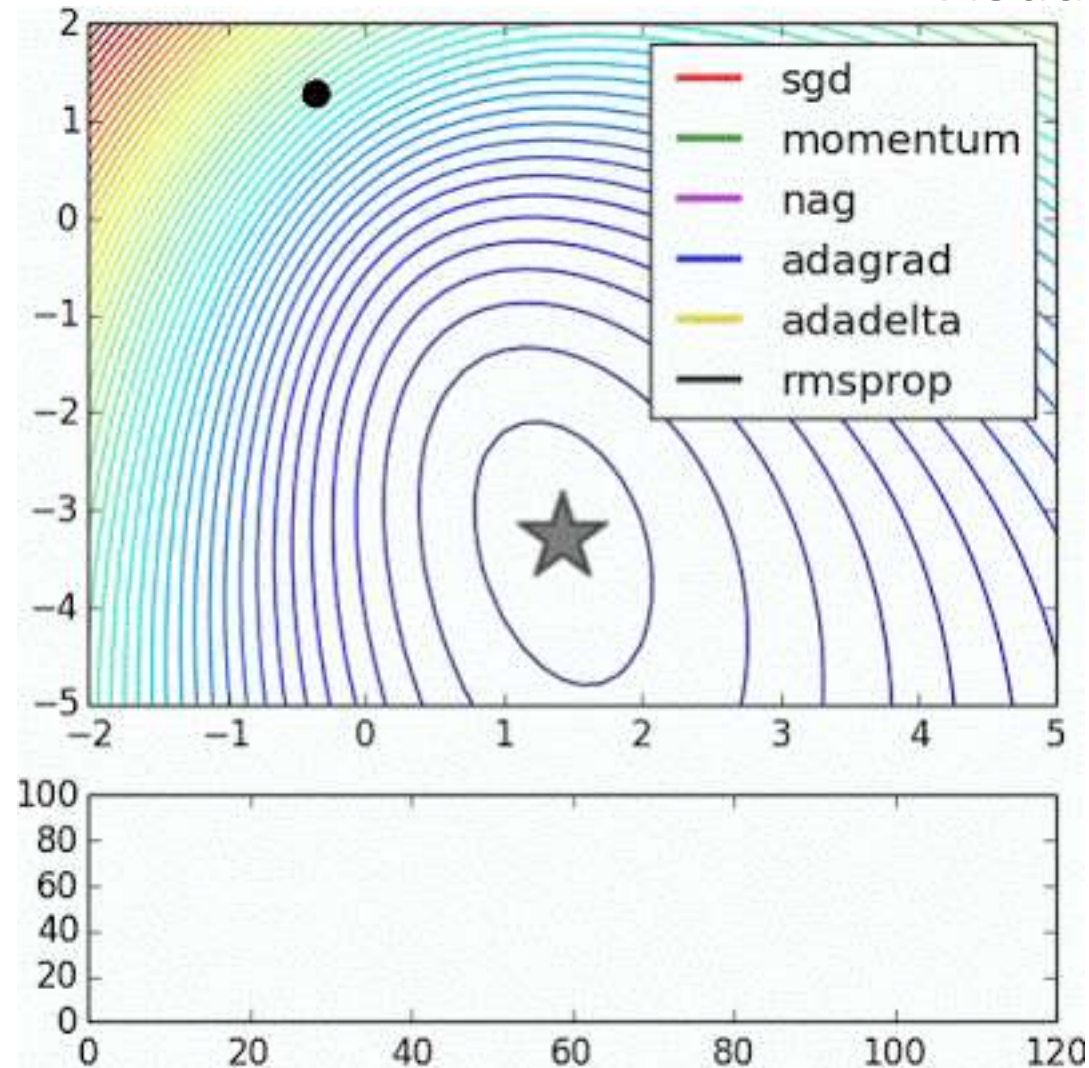
 Correct bias in second moment: $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$

 Compute update: $\Delta \theta = -\epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \delta}}$ (operations applied element-wise)

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

end while

Visualization



Find more animations at
<https://tinyurl.com/y6tkf4f8>

animation courtesy: Alec Radford



Thank You!!