

Mathematics of Data: From Theory to Computation

Prof. Volkan Cevher
volkan.cevher@epfl.ch

Lecture 8: From variance reduction to deep learning...

Laboratory for Information and Inference Systems (LIONS)
École Polytechnique Fédérale de Lausanne (EPFL)

EE-556 (Fall 2022)



License Information for Mathematics of Data Slides

- ▶ This work is released under a [Creative Commons License](#) with the following terms:
- ▶ **Attribution**
 - ▶ The licensor permits others to copy, distribute, display, and perform the work. In return, licensees must give the original authors credit.
- ▶ **Non-Commercial**
 - ▶ The licensor permits others to copy, distribute, display, and perform the work. In return, licensees may not use the work for commercial purposes – unless they get the licensor's permission.
- ▶ **Share Alike**
 - ▶ The licensor permits others to distribute derivative works only under a license identical to the one that governs the licensor's work.
- ▶ [Full Text of the License](#)

An observation of GD vs. SGD step

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k \nabla f(\mathbf{x}^k) \quad (\text{GD})$$

Lemma

Assume f is Lipschitz smooth with constant L . Then,

$$f(\mathbf{x}^{k+1}) - f(\mathbf{x}^k) \leq \left(\frac{\gamma_k^2 L}{2} - \gamma_k \right) \|\nabla f(\mathbf{x}^k)\|^2.$$

An observation of GD vs. SGD step

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k G(\mathbf{x}^k, \theta_k) \quad (\text{SGD})$$

Lemma

Assume f is Lipschitz smooth with constant L . Then,

$$\mathbb{E}[f(\mathbf{x}^{k+1}) - f(\mathbf{x}^k)] \leq \left(\frac{\gamma_k^2 L}{2} - \gamma_k \right) \mathbb{E}[\|\nabla f(\mathbf{x}^k)\|^2] + \frac{L\gamma_k^2}{2} \mathbb{E}[\|G(\mathbf{x}^k, \theta_k) - \nabla f(\mathbf{x}^k)\|^2]$$

An observation of GD vs. SGD step

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k G(\mathbf{x}^k, \theta_k) \quad (\text{SGD})$$

Lemma

Assume f is Lipschitz smooth with constant L . Then,

$$\mathbb{E}[f(\mathbf{x}^{k+1}) - f(\mathbf{x}^k)] \leq \left(\frac{\gamma_k^2 L}{2} - \gamma_k \right) \mathbb{E}[\|\nabla f(\mathbf{x}^k)\|^2] + \frac{L\gamma_k^2}{2} \mathbb{E}[\|G(\mathbf{x}^k, \theta_k) - \nabla f(\mathbf{x}^k)\|^2]$$

- Observations:**
- The variance of gradient estimate dominates as $\nabla f(\mathbf{x}^k) \rightarrow 0$.
 - To ensure convergence we need to control variance.

$$\gamma_k \rightarrow 0 \implies \text{Slow convergence!}$$

Can we decrease the variance while using a constant step-size?

An observation of GD vs. SGD step

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k G(\mathbf{x}^k, \theta_k) \quad (\text{SGD})$$

Lemma

Assume f is Lipschitz smooth with constant L . Then,

$$\mathbb{E}[f(\mathbf{x}^{k+1}) - f(\mathbf{x}^k)] \leq \left(\frac{\gamma_k^2 L}{2} - \gamma_k \right) \mathbb{E}[\|\nabla f(\mathbf{x}^k)\|^2] + \frac{L\gamma_k^2}{2} \mathbb{E}[\|G(\mathbf{x}^k, \theta_k) - \nabla f(\mathbf{x}^k)\|^2]$$

- Observations:**
- The variance of gradient estimate dominates as $\nabla f(\mathbf{x}^k) \rightarrow 0$.
 - To ensure convergence we need to control variance.

$\gamma_k \rightarrow 0 \implies$ Slow convergence!

Can we decrease the variance while using a constant step-size?

Choose a stochastic gradient, s.t. $\mathbb{E}[\|G(\mathbf{x}^k; \theta_k)\|^2] \rightarrow 0$.

A simple approach: Mini-batch SGD

- More samples imply a better estimate for full gradient.

SGD with mini batches

Let $G(\mathbf{x}, \theta)$ be an unbiased gradient estimate ($\mathbb{E}[G(\mathbf{x}, \theta)] = \nabla f(\mathbf{x})$) and B_k be the batch size. Then, we have

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k \frac{1}{B_k} \sum_{j=1}^{B_k} G(\mathbf{x}^k, \theta_{k,j}).$$

Theorem

Let $B_k > 0$ be the batch size and $G(\mathbf{x}, \theta)$ be an unbiased gradient estimate with bounded variance, i.e., $\mathbb{E}[\|G(\mathbf{x}, \theta) - \nabla f(\mathbf{x})\|^2 \mid \mathbf{x}] \leq \sigma^2$. Then, the mini-batch estimate has the following properties:

$$\mathbb{E} \left[\frac{1}{B_k} \sum_{j=1}^{B_k} G(\mathbf{x}, \theta_{k,j}) \right] = \nabla f(\mathbf{x}) \quad \text{and} \quad \mathbb{E} \left[\left\| \frac{1}{B_k} \sum_{j=1}^{B_k} G(\mathbf{x}, \theta_{k,j}) - \nabla f(\mathbf{x}) \right\|^2 \mid \mathbf{x} \right] \leq \frac{\sigma^2}{B_k}.$$

- Remarks:**
- We might need to increase the batch size over time to take variance to 0.
 - We can come up with a “smarter” estimate for $\nabla f(\mathbf{x})$.

How to construct a new estimate $G(\mathbf{x}^k; \theta_k)$? [6]

Finite sum structure:	SGD update rule:
$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) := \frac{1}{n} \sum_{j=1}^n f_j(\mathbf{x}) \right\}$	$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k \nabla f_j(\mathbf{x}^k)$

- Let $X = \nabla f_j(\mathbf{x}^k)$ be a random variable (due to $j \sim \text{Uniform}(\{1, \dots, n\})$).
- Let $Y = \nabla f_j(\tilde{\mathbf{x}})$ be another random variable, and $\tilde{\mathbf{x}}$ is a particularly selected point.

- Remarks:**
- We want X and Y to be correlated (we will see why!).
 - Given Y , we should be able to estimate $\mathbb{E}[X]$ with more confidence.

- Observations:**
- Choice of $\tilde{\mathbf{x}}$ affects how correlated X and Y are.
 - We can compute $\mathbb{E}[Y] = \frac{1}{n} \sum_{j=1}^n \nabla f_j(\tilde{\mathbf{x}}) = \nabla f(\tilde{\mathbf{x}})$.

- Goal:**
- Find a **good** estimate of $\mathbb{E}[X] = \frac{1}{n} \sum_{j=1}^n \nabla f_j(\mathbf{x}^k) = \nabla f(\mathbf{x}^k)$.

How to construct a new estimate $G(\mathbf{x}^k; \theta_k)$? [6]

Finite sum structure:	SGD update rule:
$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) := \frac{1}{n} \sum_{j=1}^n f_j(\mathbf{x}) \right\}$	$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k \nabla f_j(\mathbf{x}^k)$

- Let $X = \nabla f_j(\mathbf{x}^k)$ be a random variable (due to $j \sim \text{Uniform}(\{1, \dots, n\})$).
- Let $Y = \nabla f_j(\tilde{\mathbf{x}})$ be another random variable, and $\tilde{\mathbf{x}}$ is a particularly selected point.

A generalized estimator: $R_\alpha = \alpha(X - Y) + \mathbb{E}[Y]$

- $\mathbb{E}[R_\alpha] = \alpha \mathbb{E}[X] + (1 - \alpha) \mathbb{E}[Y]$
- $\text{Var}(R_\alpha) = \alpha^2 (\text{Var}(X) + \text{Var}(Y) - 2\text{Cov}(X, Y))$

- Observations:**
- When $\alpha = 1$, R_α becomes unbiased, i.e., $\mathbb{E}[R_\alpha] = \mathbb{E}[X]$.
 - If $\text{Cov}(X, Y)$ is large enough (X and Y are correlated enough), $\text{Var}(R_\alpha) \leq \text{Var}(X)$.

How could we use this information to construct our estimate?

Variance reduction techniques: SVRG

- Select the stochastic gradient ∇f_{i_k} , and compute a gradient estimate

$$\mathbf{r}_k = \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}) + \nabla f(\tilde{\mathbf{x}}).$$

- As $\tilde{\mathbf{x}} \rightarrow \mathbf{x}^*$ and $\mathbf{x}^k \rightarrow \mathbf{x}^*$, we have

$$\nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}) + \nabla f(\tilde{\mathbf{x}}) \rightarrow 0.$$

- As a result, we can ensure the following

$$\mathbb{E}[\|\nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}) + \nabla f(\tilde{\mathbf{x}})\|^2] \rightarrow 0.$$

Remarks:

- Remember the generalized estimator: $R_\alpha = \alpha(X - Y) + \mathbb{E}[Y]$.
- For SVRG, $\alpha = 1$, $X = \nabla f_{i_k}(\mathbf{x}^k)$ and $Y = \nabla f_{i_k}(\tilde{\mathbf{x}})$.
- We will see how $\tilde{\mathbf{x}}$ is computed!

Stochastic gradient algorithm with variance reduction

Stochastic gradient with variance reduction (SVRG) [11, 21]

1. Choose $\tilde{\mathbf{x}}^0 \in \mathbb{R}^p$ as a starting point and $\gamma > 0$ and $q \in \mathbb{N}_+$.

2. For $s = 0, 1, 2, \dots$, perform:

2a. $\tilde{\mathbf{x}} = \tilde{\mathbf{x}}^s$, $\tilde{\mathbf{v}} = \nabla f(\tilde{\mathbf{x}})$, $\mathbf{x}^0 = \tilde{\mathbf{x}}$.

2b. For $k = 0, 1, \dots, q-1$, perform:

$$\begin{cases} \text{Pick } i_k \in \{1, \dots, n\} \text{ uniformly at random} \\ \mathbf{r}_k = \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}) + \tilde{\mathbf{v}} \\ \mathbf{x}^{k+1} := \mathbf{x}^k - \gamma \mathbf{r}_k, \end{cases} \quad (1)$$

2c. Update $\tilde{\mathbf{x}}^{s+1} = \frac{1}{m} \sum_{j=0}^{q-1} \mathbf{x}^j$.

Features

- ▶ The SVRG method uses a multistage scheme to reduce the **variance** of the **stochastic gradient** \mathbf{r}_k .
- ▶ **Learning rate** γ does not necessarily tend to 0 while \mathbf{x}^k and $\tilde{\mathbf{x}}^s$ tend to \mathbf{x}_* .
- ▶ Each stage, SVRG uses $n + 2q$ component **gradient** evaluations.
- ▶ n for the **full gradient** at the beginning of each stage, and $2q$ for each of the q **stochastic gradient steps**.

Convergence analysis

Assumption A5.

- (i) f is μ -strongly convex
- (ii) The learning rate $0 < \gamma < 1/(4L_{\max})$, where $L_{\max} = \max_{1 \leq j \leq n} L_j$.
- (iii) q is large enough such that

$$\kappa = \frac{1}{\mu\gamma(1 - 4\gamma L_{\max})q} + \frac{4\gamma L_{\max}(q + 1)}{(1 - 4\gamma L_{\max})q} < 1.$$

Theorem

Assumptions:

- ▶ The sequence $\{\tilde{\mathbf{x}}^s\}_{k \geq 0}$ is generated by SVRG.
- ▶ Assumption A5 is satisfied.

Conclusion: Linear convergence is obtained:

$$\mathbb{E}f(\tilde{\mathbf{x}}^s) - f(\mathbf{x}^*) \leq \kappa^s (f(\tilde{\mathbf{x}}^0) - f(\mathbf{x}^*)).$$

Choice of γ and q , and complexity

Chose γ and q such that $\kappa \in (0, 1)$:

For example

$$\gamma = 0.1/L_{\max}, q = 100(L_{\max}/\mu) \implies \kappa \approx 5/6.$$

Complexity

$$\mathbb{E}f(\tilde{\mathbf{x}}^s) - f(\mathbf{x}^*) \leq \epsilon, \quad \text{when } s \geq \log((f(\tilde{\mathbf{x}}^0) - f(\mathbf{x}^*))/\epsilon) / \log(\kappa^{-1})$$

- ▶ Each stage needs $n + 2q$ **component gradient evaluations**
- ▶ With $q = \mathcal{O}(L_{\max}/\mu)$, we obtain an **overall complexity** of

$$\mathcal{O}\left((n + L_{\max}/\mu) \log(1/\epsilon)\right).$$

Comparison: GD vs. SGD vs. SVRG

- GD update:

$$\left\{ \begin{array}{l} \mathbf{x}^{k+1} := \mathbf{x}^k - \gamma \nabla f(\mathbf{x}^k), \end{array} \right.$$

- SGD update:

$$\left\{ \begin{array}{l} \mathbf{x}^{k+1} := \mathbf{x}^k - \gamma \nabla f_{i_k}(\mathbf{x}^k), \end{array} \right.$$

- SVRG update:

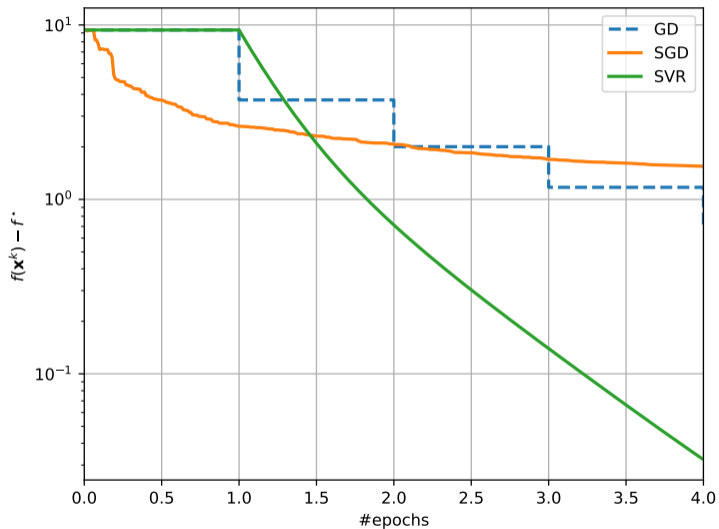
$$\left\{ \begin{array}{l} \mathbf{r}_k = \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}) + \nabla f(\tilde{\mathbf{x}}) \\ \mathbf{x}^{k+1} := \mathbf{x}^k - \gamma \mathbf{r}_k, \end{array} \right.$$

	SGD	SVRG	GD
Requires gradient storage?	no	no	no
Epoch-based	no	yes	no
Parameters	stepsize	stepsize & epoch length	stepsize
Gradient evaluations	1 per iteration	$n + 2q$ per epoch	n per iteration

Table: Comparisons of SGD, SVRG and GD [6]

- Recall that $q = \mathcal{O}(L_{\max}/\mu)$ is the epoch length for SVRG.

Example: ℓ_2 -regularized least squares with synthetic data



Taxonomy of algorithms

$$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) := \frac{1}{n} \sum_{j=1}^n f_j(\mathbf{x}) \right\}.$$

- $f(\mathbf{x}) = \frac{1}{n} \sum_{j=1}^n f_j(\mathbf{x})$: μ -strongly convex with L -Lipschitz continuous gradient.

SVRG	GD	SGD
Linear	Linear	Sublinear

Table: Rate of convergence.

- $\kappa = L/\mu$.

SVRG	AGD	SGD
$\mathcal{O}((n + \kappa) \log(1/\varepsilon))$	$\mathcal{O}(n\kappa \log(1/\varepsilon))$	$1/\varepsilon$

Table: Complexity to obtain ε -solution.

The variance reduction zoo: convex

Setting	Algorithm	Lower bound	Complexity bound
L -smooth f_i 's with bounded variance	Gradient descent SVRG ($B_k = 1$) [16] SVRG ($B_k = \Omega(n^{2/3})$) [16] SAGA ($B_k = 1$) [16] SAGA ($B_k = \Omega(n^{2/3})$) [16] SpiderBoost [19] SpiderBoost-M [19] Spider [10] PAGE [15]	$L\Delta_0 \min\{\sigma/\epsilon^3, \sqrt{n}/\epsilon^2\}$ [10]	$nL\Delta_0/\epsilon^2$ $nL\Delta_0/\epsilon^2$ $n^{2/3}L\Delta_0/\epsilon^2$ $nL\Delta_0/\epsilon^2$ $n^{2/3}L\Delta_0/\epsilon^2$ $\sqrt{n}L\Delta_0/\epsilon^2$ $\sqrt{n}L\Delta_0/\epsilon^2$ $L\Delta_0 \min\{\sigma/\epsilon^3, \sqrt{n}/\epsilon^2\}$ $L\Delta_0 \min\{\sigma/\epsilon^3, \sqrt{n}/\epsilon^2\}$
f is μ -SCVX and L -smooth f_i 's are average L -smooth	KatyushaX [3]	$(n + n^{3/4} \sqrt{\frac{L}{\mu}}) \log \frac{\Delta_0}{\epsilon}$ [22]	$(n + n^{3/4} \sqrt{\frac{L}{\mu}}) \log \frac{\Delta_0}{\epsilon}$
f is CVX and L -smooth f_i 's are average L -smooth	KatyushaX [3]	$n + n^{3/4} \sqrt{\frac{LD_0^2}{\epsilon}}$ [23]	$n + n^{3/4} \sqrt{\frac{LD_0^2}{\epsilon}}$

- Remarks:**
- Complexity ((S)CVX f): total number of stochastic first-order oracle calls to find \mathbf{x}_ϵ^* with $\mathbb{E}[f(\mathbf{x}_\epsilon^*) - f(\mathbf{x}^*)] \leq \epsilon$.
 - $\Delta_0 = f(\mathbf{x}^0) - f^*$, $D_0 = \|\mathbf{x}^0 - \mathbf{x}^*\|$.
 - Bounded variance: $\mathbb{E}_i[\|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|^2] \leq \sigma^2 \quad \forall \mathbf{x}$.
 - Average L -smooth: $\mathbb{E}_i[\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\|^2] \leq L^2 \|\mathbf{x} - \mathbf{y}\|^2 \quad \forall \mathbf{x}, \mathbf{y}$.

Variance-reduction for **non-convex** problems

SVRG estimator vs. a **recursive** estimator

o SVRG update:

$$\begin{cases} \mathbf{r}_1 = \nabla f(\tilde{\mathbf{x}}) \\ \mathbf{r}_k := \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}) + \nabla f(\tilde{\mathbf{x}}) \\ \mathbf{x}^{k+1} := \mathbf{x}^k - \gamma \mathbf{r}_k, \end{cases}$$

o Spider [10] update:

$$\begin{cases} \mathbf{r}_1 = \nabla f(\tilde{\mathbf{x}}) \\ \mathbf{r}_k := \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}) + \mathbf{r}_{k-1} \\ \mathbf{x}^{k+1} := \mathbf{x}^k - \gamma \mathbf{r}_k, \end{cases}$$

Spider [10]

1. Choose $\mathbf{x}^0 \in \mathbb{R}^p$ as a starting point and $\gamma = \epsilon/L$.

2. For $k = 0, 1, 2, \dots$, perform:

2a. If $t \bmod n = 0$, do:

$$\mathbf{r}_k = \nabla f(\mathbf{x}^k)$$

else:

Pick $i_k \in \{1, \dots, n\}$ uniformly at random

$$\mathbf{r}_k = \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\mathbf{x}^{k-1}) + \mathbf{r}_{k-1}$$

2b. Update $\mathbf{x}^{k+1} := \mathbf{x}^k - \frac{\gamma}{\|\mathbf{r}_k\|} \mathbf{r}_k$

3. Return \mathbf{x}^k

Remarks:

o Sample complexity: $O\left(n + \sqrt{n} \frac{\Delta L}{\epsilon^2}\right)$.

o Sets the final accuracy a priori.

o Step-size depends on ϵ and L .

Adaptive variance-reduction for non-convex problems

AdaSpider [13]

1. Choose $\mathbf{x}^0 \in \mathbb{R}^p$ as a starting point.
2. For $k = 0, 1, 2, \dots$, perform:
 - 2a. If $t \bmod n = 0$, do:
 $\mathbf{r}_k = \nabla f(\mathbf{x}^k)$
else:
Pick $i_k \in \{1, \dots, n\}$ uniformly at random
 $\mathbf{r}_k = \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\mathbf{x}^{k-1}) + \mathbf{r}_{k-1}$
 - 2b. Compute $\gamma_k := 1 / \left(n^{1/4} \sqrt{n^{1/2} + \sum_{i=0}^k \|\mathbf{r}_i\|^2} \right)$
 - 2c. Update $\mathbf{x}^{k+1} := \mathbf{x}^k - \gamma_k \mathbf{r}_k$
3. Return \mathbf{x}^k

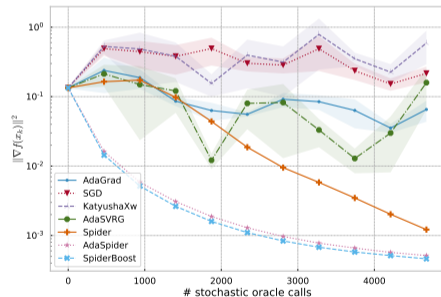
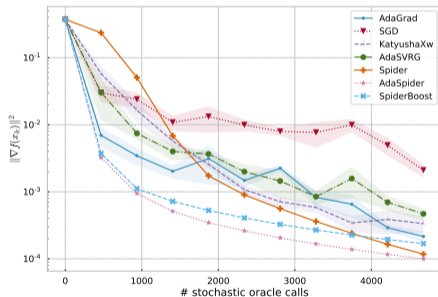
Theorem

Let $\Delta_0 = f(\mathbf{x}^0) - \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$. The sequence $\mathbf{x}^0, \dots, \mathbf{x}^k$ generated by AdaSpider satisfies:

$$\frac{1}{k} \sum_{i=0}^{k-1} \mathbb{E}[\|\nabla f(\mathbf{x}^i)\|] \leq O \left(n^{1/4} \frac{\Delta_0 + L^2}{\sqrt{k}} \log(k) \right), \quad \text{with sample complexity } \tilde{O} \left(n + \sqrt{n} \frac{\Delta_0^2 + L^4}{\varepsilon^2} \right).$$

Performance of AdaSpider

- Image classification with neural networks (spoiler alert!) trained with cross entropy loss.
- AdaGrad [8], KatyushaXw [2], AdaSVRG[7], Spider [10], SpiderBoost [20].



The variance reduction zoo: non-convex

Setting	Algorithm	Lower bound	Complexity bound
f is α -weakly CVX and L -smooth f_i 's are average L -smooth	Spider [10]	$\frac{\Delta_0}{\epsilon^2} \min\{n^{3/4} \sqrt{\alpha L}, \sqrt{n}L\}$ [23]	$\frac{\Delta_0}{\epsilon^2} \min\{n^{3/4} \sqrt{\alpha L}, \sqrt{n}L\}$
f_i 's are α -weakly CVX and L -smooth	Natasha [1]	$\frac{\Delta_0}{\epsilon^2} \min\{\sqrt{n\alpha L}, L\}$ [23]	$\frac{\Delta_0}{\epsilon^2} \min\{\sqrt{n\alpha L}, \sqrt{n}L\}$
f is non-CVX f_i 's are non-CVX and L -smooth	AdaSpider [13]	$\frac{\Delta_0 L}{\epsilon^2} \sqrt{n}$ [23, 10]	$\tilde{O}\left(n + \frac{\Delta_0^2 + L^4}{\epsilon^2} \sqrt{n}\right)$

- Remarks:**
- Complexity (nonCVX f): total number of stochastic first-order oracle calls to find \mathbf{x}_ϵ^* with $\mathbb{E}[\|\nabla f(\mathbf{x}_\epsilon^*)\|^2] \leq \epsilon^2$.
 - $\Delta_0 = f(\mathbf{x}^0) - f^*$, $D_0 = \|\mathbf{x}^0 - \mathbf{x}^*\|$.
 - Bounded variance: $\mathbb{E}_i[\|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|^2] \leq \sigma^2 \quad \forall \mathbf{x}$.
 - L -smooth: $\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y}$.
 - $f(\mathbf{x})$ is α -weakly convex if $f(\mathbf{x}) + \frac{\alpha}{2}\|\mathbf{x}\|^2$ is convex $\forall \mathbf{x}$.

Deep learning outline

- In the sequel,
 - ▶ Introduction to deep learning
 - ▶ The deep learning paradigm
 - ▶ Challenges in deep learning theory and applications
- Next class
 - ▶ Generalization in deep learning

Remark about notation

The Deep Learning literature might use a different notation:

	Our lectures	DL literature
data/sample	\mathbf{a}	\mathbf{x}
label	b	y
bias	μ	b
weight	\mathbf{x}, \mathbf{X}	\mathbf{w}, \mathbf{W}

Power of linear classifiers–I

Problem (Recall: Logistic regression)

Given a sample vector $\mathbf{a}_i \in \mathbb{R}^d$ and a binary class label $b_i \in \{-1, +1\}$ ($i = 1, \dots, n$), we define the conditional probability of b_i given \mathbf{a}_i as follows:

$$\mathbb{P}(b_i | \mathbf{a}_i, \mathbf{x}) \propto 1 / (1 + e^{-b_i \langle \mathbf{x}, \mathbf{a}_i \rangle}),$$

where $\mathbf{x} \in \mathbb{R}^d$ is some weight vector.

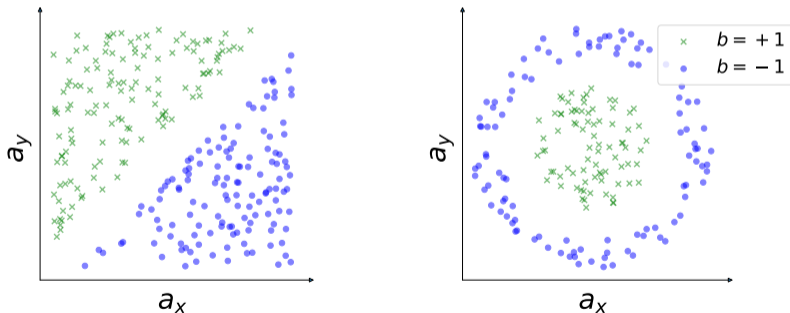


Figure: Linearly separable versus non-linearly separable dataset

Power of linear classifiers–II

- Lifting dimensions to the rescue
 - ▶ Convex optimization objective
 - ▶ Side effect: **The curse-of-dimensionality**
 - ▶ Possible to avoid via kernel methods, such as SVMs

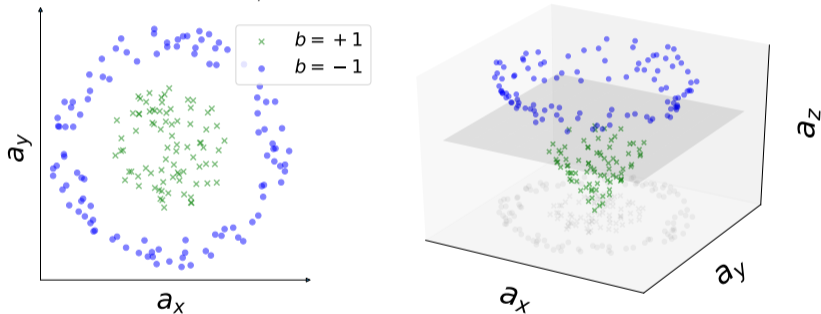


Figure: Non-linearly separable data (left). Linearly separable in \mathbb{R}^3 via $\mathbf{a}_z = \sqrt{\mathbf{a}_x^2 + \mathbf{a}_y^2}$ (right).

An important alternative for non-linearly separable data

1-hidden-layer neural network with m neurons (fully-connected architecture):

- Parameters: $\mathbf{X}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{X}_2 \in \mathbb{R}^{c \times m}$ (weights), $\mu_1 \in \mathbb{R}^m$, $\mu_2 \in \mathbb{R}^c$ (biases)
- Activation function: $\sigma : \mathbb{R} \rightarrow \mathbb{R}$

$$h_{\mathbf{x}}(\mathbf{a}) :=$$



An important alternative for non-linearly separable data

1-hidden-layer neural network with m neurons (fully-connected architecture):

- Parameters: $\mathbf{X}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{X}_2 \in \mathbb{R}^{c \times m}$ (weights), $\mu_1 \in \mathbb{R}^m$, $\mu_2 \in \mathbb{R}^c$ (biases)
- Activation function: $\sigma : \mathbb{R} \rightarrow \mathbb{R}$

$$h_{\mathbf{x}}(\mathbf{a}) := \left[\begin{array}{c} \text{weight} \\ \downarrow \\ \mathbf{X}_1 \end{array} \right] \left[\begin{array}{c} \text{input} \\ \downarrow \\ \mathbf{a} \end{array} \right]$$

An important alternative for non-linearly separable data

1-hidden-layer neural network with m neurons (fully-connected architecture):

- Parameters: $\mathbf{X}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{X}_2 \in \mathbb{R}^{c \times m}$ (weights), $\mu_1 \in \mathbb{R}^m$, $\mu_2 \in \mathbb{R}^c$ (biases)
- Activation function: $\sigma : \mathbb{R} \rightarrow \mathbb{R}$

$$h_{\mathbf{x}}(\mathbf{a}) := \left[\begin{array}{c} \text{weight} \\ \downarrow \\ \mathbf{X}_1 \end{array} \right] \left[\begin{array}{c} \text{input} \\ \downarrow \\ \mathbf{a} \end{array} \right] + \left[\begin{array}{c} \text{bias} \\ \downarrow \\ \mu_1 \end{array} \right]$$

An important alternative for non-linearly separable data

1-hidden-layer neural network with m neurons (fully-connected architecture):

- Parameters: $\mathbf{X}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{X}_2 \in \mathbb{R}^{c \times m}$ (weights), $\mu_1 \in \mathbb{R}^m$, $\mu_2 \in \mathbb{R}^c$ (biases)
- Activation function: $\sigma : \mathbb{R} \rightarrow \mathbb{R}$

$$h_{\mathbf{x}}(\mathbf{a}) := \underbrace{\sigma \left(\mathbf{X}_1 \begin{bmatrix} \mathbf{a} \end{bmatrix} + \begin{bmatrix} \mu_1 \end{bmatrix} \right)}_{\text{hidden layer = learned features}}$$

The diagram illustrates the computation of the hidden layer output. It shows the activation function σ applied to the linear combination of the input \mathbf{a} (green vector) and bias μ_1 (blue vector) through the weight matrix \mathbf{X}_1 (black matrix). The weight matrix is labeled "weight" with a downward arrow. The input vector is labeled "input" with a downward arrow. The bias vector is labeled "bias" with a downward arrow. The activation function is labeled "activation" with a downward arrow. The entire expression is enclosed in large red parentheses, and a red bracket underneath labels it as "hidden layer = learned features".

An important alternative for non-linearly separable data

1-hidden-layer neural network with m neurons (fully-connected architecture):

- Parameters: $\mathbf{X}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{X}_2 \in \mathbb{R}^{c \times m}$ (weights), $\mu_1 \in \mathbb{R}^m$, $\mu_2 \in \mathbb{R}^c$ (biases)
- Activation function: $\sigma : \mathbb{R} \rightarrow \mathbb{R}$

$$h_{\mathbf{x}}(\mathbf{a}) := \left[\mathbf{X}_2 \right] \sigma \left(\underbrace{\left[\mathbf{X}_1 \right] \left[\mathbf{a} \right] + \left[\mu_1 \right]}_{\text{hidden layer = learned features}} \right)$$

The diagram illustrates the computation of the hidden layer output. The input vector \mathbf{a} (green) is multiplied by the weight matrix \mathbf{X}_1 (black) and the bias vector μ_1 (blue) is added. This sum is then passed through the activation function σ (red). The resulting vector is then multiplied by the weight matrix \mathbf{X}_2 (black) to produce the final output. The entire process is labeled as "hidden layer = learned features".

An important alternative for non-linearly separable data

1-hidden-layer neural network with m neurons (fully-connected architecture):

- Parameters: $\mathbf{X}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{X}_2 \in \mathbb{R}^{c \times m}$ (weights), $\mu_1 \in \mathbb{R}^m$, $\mu_2 \in \mathbb{R}^c$ (biases)
- Activation function: $\sigma : \mathbb{R} \rightarrow \mathbb{R}$

$$h_{\mathbf{x}}(\mathbf{a}) := \left[\mathbf{X}_2 \right] \underbrace{\left(\sigma \left(\left[\mathbf{X}_1 \right] \left[\mathbf{a} \right] + \left[\mu_1 \right] \right) \right)}_{\text{hidden layer = learned features}} + \left[\mu_2 \right]$$

The diagram illustrates the computation of the hidden layer output. The input vector \mathbf{a} (green) is multiplied by the weight matrix \mathbf{X}_1 (black) and the bias vector μ_1 (blue) is added. This sum is then passed through the activation function σ (red). The result is then multiplied by the weight matrix \mathbf{X}_2 (black) and the bias vector μ_2 (blue) is added to produce the final output. The entire process from \mathbf{X}_1 to μ_2 is labeled as the "hidden layer = learned features".

An important alternative for non-linearly separable data

1-hidden-layer neural network with m neurons (fully-connected architecture):

- Parameters: $\mathbf{X}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{X}_2 \in \mathbb{R}^{c \times m}$ (weights), $\mu_1 \in \mathbb{R}^m$, $\mu_2 \in \mathbb{R}^c$ (biases)
- Activation function: $\sigma : \mathbb{R} \rightarrow \mathbb{R}$

$$h_{\mathbf{x}}(\mathbf{a}) := \left[\begin{array}{c} \mathbf{X}_2 \end{array} \right] \underbrace{\left(\begin{array}{c} \text{activation} \\ \downarrow \\ \sigma \left(\begin{array}{c} \text{weight} \\ \downarrow \\ \left[\begin{array}{c} \mathbf{X}_1 \end{array} \right] \left[\begin{array}{c} \text{input} \\ \downarrow \\ \mathbf{a} \end{array} \right] + \left[\begin{array}{c} \text{bias} \\ \downarrow \\ \mu_1 \end{array} \right] \end{array} \right) \\ \text{hidden layer = learned features} \end{array} \right) + \left[\begin{array}{c} \text{bias} \\ \downarrow \\ \mu_2 \end{array} \right], \quad \mathbf{x} := [\mathbf{X}_1, \mathbf{X}_2, \mu_1, \mu_2]$$

An important alternative for non-linearly separable data

1-hidden-layer neural network with m neurons (fully-connected architecture):

- Parameters: $\mathbf{X}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{X}_2 \in \mathbb{R}^{c \times m}$ (weights), $\mu_1 \in \mathbb{R}^m$, $\mu_2 \in \mathbb{R}^c$ (biases)
- Activation function: $\sigma : \mathbb{R} \rightarrow \mathbb{R}$

$$h_{\mathbf{x}}(\mathbf{a}) := \left[\begin{array}{c} \mathbf{X}_2 \end{array} \right] \underbrace{\left(\begin{array}{c} \text{activation} \\ \downarrow \\ \sigma \left(\begin{array}{c} \text{weight} \\ \downarrow \\ \left[\begin{array}{c} \mathbf{X}_1 \end{array} \right] \left[\begin{array}{c} \text{input} \\ \downarrow \\ \mathbf{a} \end{array} \right] + \left[\begin{array}{c} \text{bias} \\ \downarrow \\ \mu_1 \end{array} \right] \end{array} \right) \\ \text{hidden layer = learned features} \end{array} \right) + \left[\begin{array}{c} \text{bias} \\ \downarrow \\ \mu_2 \end{array} \right], \quad \mathbf{x} := [\mathbf{X}_1, \mathbf{X}_2, \mu_1, \mu_2]$$

recursively repeat activation + affine transformation to obtain “deeper” networks.

Why neural networks?: An approximation theoretic motivation

Theorem (Universal approximation [5])

Let $\sigma(\cdot)$ be a nonconstant, bounded, and increasing continuous function. Let $I_d = [0, 1]^d$. The space of continuous functions on I_d is denoted by $\mathcal{C}(I_d)$.

Given $\epsilon > 0$ and $g \in \mathcal{C}(I_d)$ there exists a 1-hidden-layer network h with m neurons such that h is an ϵ -approximation of g , i.e.,

$$\sup_{\mathbf{a} \in I_d} |g(\mathbf{a}) - h(\mathbf{a})| \leq \epsilon$$

Caveat

The number of neurons m needed to approximate some function g can be arbitrarily large!

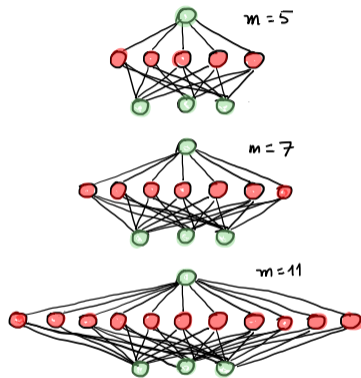


Figure: networks of increasing width

Why were NNs not popular before 2010?

- too big to optimize!
- did not have enough data
- could not find the optimum via algorithms

Why were NNs not popular before 2010?

- too big to optimize!
- did not have enough data
- could not find the optimum via algorithms

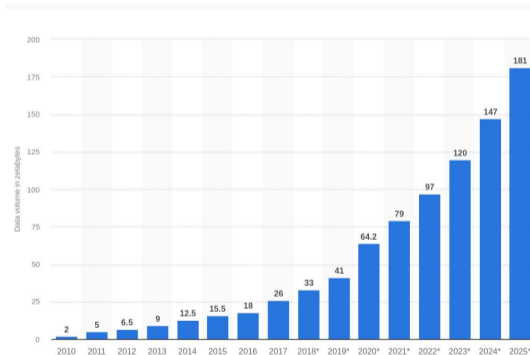
Market Summary > NVIDIA Corporation

121.01 USD

+71.78 (145.81%) ↑ past 5 years

19 Oct, 11:35 GMT-4 • Disclaimer

1D | 5D | 1M | 6M | YTD | 1Y | **5Y** | Max



Supervised learning: Multi-class classification

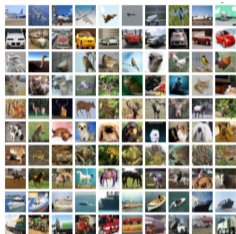


Figure: CIFAR10 dataset: 60000 32x32 color images (3 channels) from 10 classes

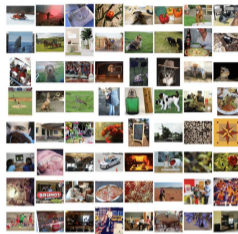


Figure: Imagenet dataset: 14 million color images (varying resolution, 3 channels) from 21K classes

Goal

Image-label pairs $(\mathbf{a}, b) \subseteq \mathbb{R}^d \times \{1, \dots, c\}$ follow an unknown distribution \mathbb{P} . Find $h : \mathbb{R}^d \rightarrow \{1, \dots, c\}$ with minimum *misclassification probability*

$$\min_{h \in \mathcal{H}} \mathbb{P}(h(\mathbf{a}) \neq b)$$

2010-today: Deep Learning becomes popular again

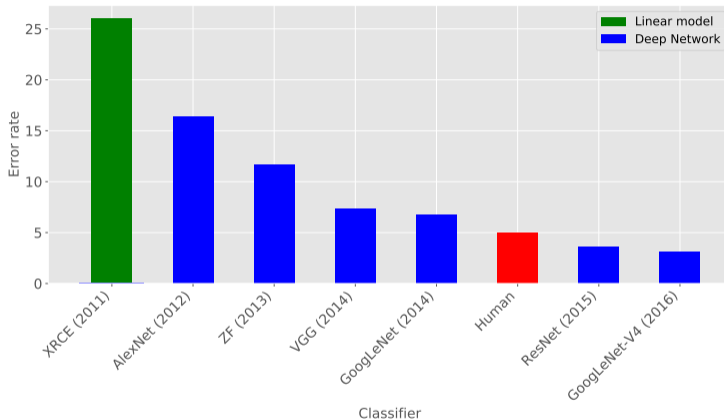


Figure: Error rate on the ImageNet challenge, for different network architectures.

2010-today: Deep Learning becomes popular again

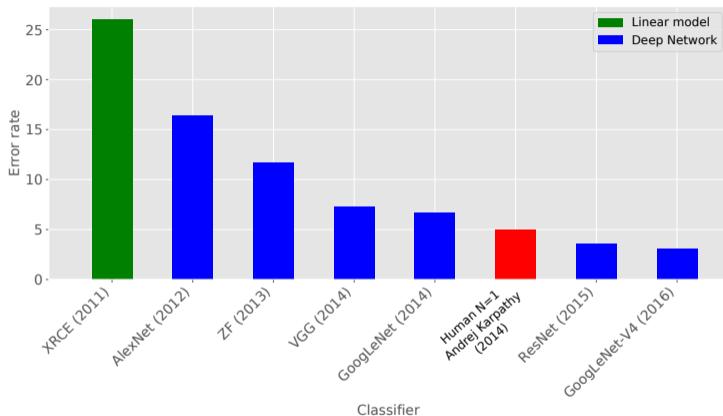


Figure: Error rate on the ImageNet challenge, for different network architectures [17, 12].

Convolutional architectures in Computer Vision tasks

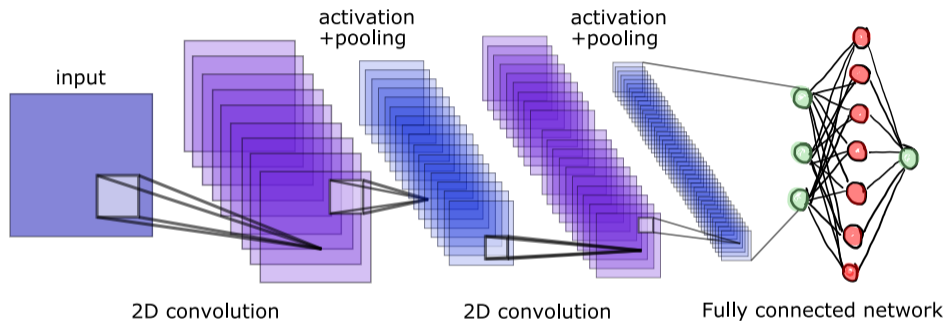
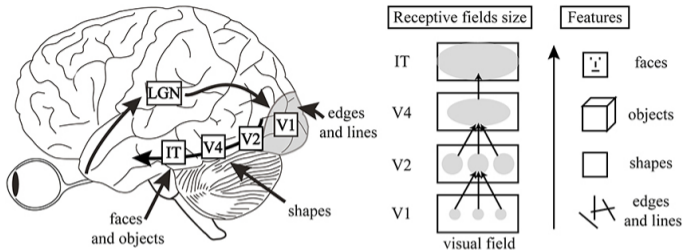
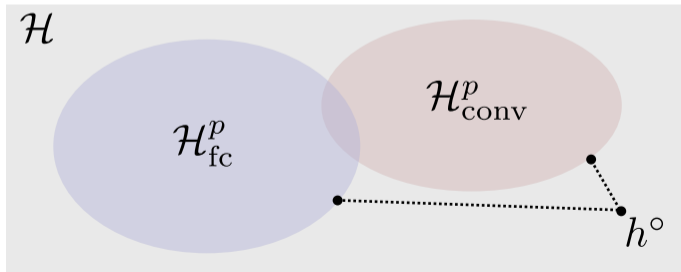


Figure: "Locality" structure of a 2D deep convolutional neural network.

Inductive Bias: Why convolution works so well in Computer Vision tasks?



h°	true unknown function
\mathcal{H}	space of all functions
\mathcal{H}_{fc}^p	fully-connected networks with p parameters
\mathcal{H}_{conv}^p	convolutional networks with p parameters



2010-today: Size of neural networks grows exponentially!

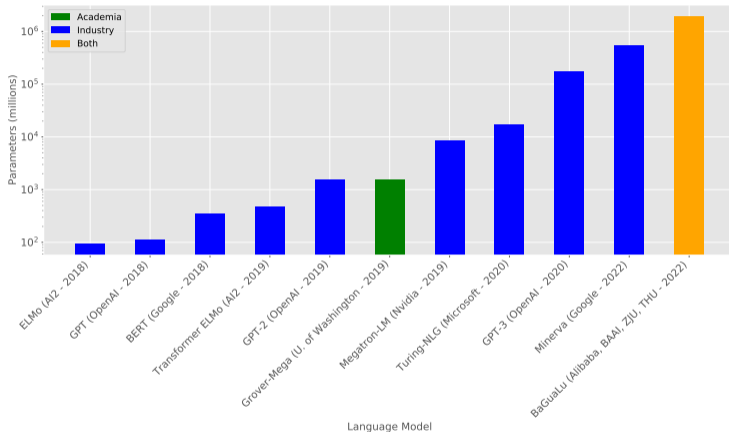


Figure: Number of parameters in language models based on Deep Learning.

The landscape of ERM with multilayer networks

Recall: Empirical risk minimization (ERM)

Let $h_{\mathbf{x}} : \mathbb{R}^n \rightarrow \mathbb{R}$ be network and let $\{(\mathbf{a}_i, b_i)\}_{i=1}^n$ be a sample with $b_i \in \{-1, 1\}$ and $\mathbf{a}_i \in \mathbb{R}^n$. The *empirical risk minimization* (ERM) is defined as follows

$$\min_{\mathbf{x}} \left\{ R_n(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n L(h_{\mathbf{x}}(\mathbf{a}_i), b_i) \right\} \quad (2)$$

where $L(h_{\mathbf{x}}(\mathbf{a}_i), b_i)$ is the loss on the sample (\mathbf{a}_i, b_i) and \mathbf{x} are the parameters of the network.

Some frequently used loss functions

- ▶ $L(h_{\mathbf{x}}(\mathbf{a}), b) = \log(1 + \exp(-b \cdot h_{\mathbf{x}}(\mathbf{a})))$ (logistic loss)
- ▶ $L(h_{\mathbf{x}}(\mathbf{a}), b) = (b - h_{\mathbf{x}}(\mathbf{a}))^2$ (squared error)
- ▶ $L(h_{\mathbf{x}}(\mathbf{a}), b) = \max(0, 1 - b \cdot h_{\mathbf{x}}(\mathbf{a}))$ (hinge loss)

The landscape of ERM with multilayer networks

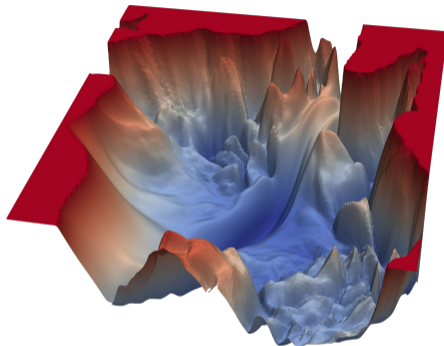
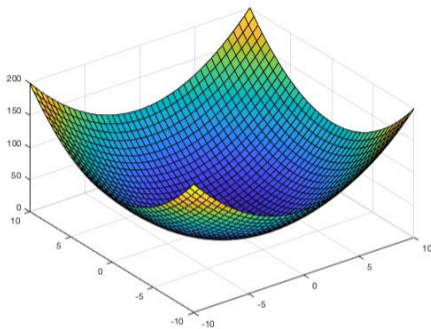


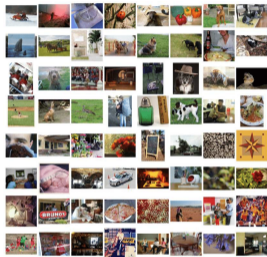
Figure: convex (left) vs non-convex (right) optimization landscape [14]

Conventional wisdom in ML until 2010:
Simple models + simple errors

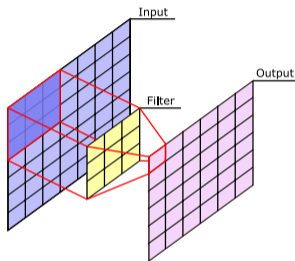
The landscape of ERM with multilayer networks

1. Simsek, Berfin, et al. *Geometry of the loss landscape in overparameterized neural networks: Symmetries and invariances*. International Conference on Machine Learning. PMLR, 2021.
2. Foret, Pierre, et al. *Sharpness-aware minimization for efficiently improving generalization*. arXiv preprint arXiv:2010.01412 (2020).
3. Jiang, Yiding, et al. *Fantastic generalization measures and where to find them*. arXiv preprint arXiv:1912.02178 (2019).
4. Dziugaite, Gintare Karolina, and Daniel M. Roy. *Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data*. arXiv preprint arXiv:1703.11008 (2017).
5. Keskar, Nitish Shirish, et al. *On large-batch training for deep learning: Generalization gap and sharp minima*. arXiv preprint arXiv:1609.04836 (2016).

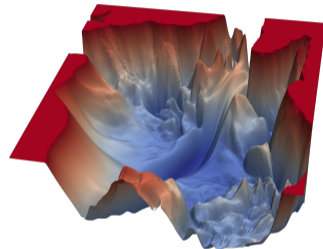
The deep learning paradigm



(a) Massive datasets



(b) Inductive bias from large and complex architectures



(c) ERM using stochastic non-convex first-order optimization algorithms (SGD)

Figure: Most common components in a Deep Learning Pipeline

Challenges in DL/ML applications: Robustness (I)



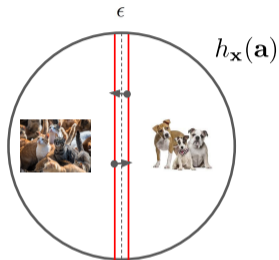
(a) Turtle classified as rifle [4].



(b) Stop sign classified as 45 mph sign [9].

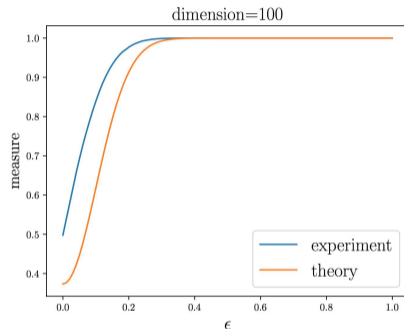
Figure: Natural or human-crafted modifications that trick neural networks used in computer vision tasks

Challenges in DL/ML applications: Robustness (II)



$$h_{\mathbf{x}}(\mathbf{a} + \epsilon) \approx h_{\mathbf{x}}(\mathbf{a}) + \langle \epsilon, \nabla h_{\mathbf{x}}(\mathbf{a}) \rangle$$
$$|h_{\mathbf{x}}(\mathbf{a} + \epsilon) - h_{\mathbf{x}}(\mathbf{a})| \leq \|\epsilon\| \|\nabla h_{\mathbf{x}}(\mathbf{a})\|$$

(a) Linear classifier on data distributed on a sphere



(b) Concentration of measure phenomenon on high dimensions

Figure: Understanding the robustness of a classifier in high-dimensional spaces [18]

Challenges in DL/ML applications: Robustness (References I)

1. Madry, Aleksander and Makelov, Aleksandar and Schmidt, Ludwig and Tsipras, Dimitris and Vladu, Adrian. *Towards Deep Learning Models Resistant to Adversarial Attacks*. ICLR 2018.
2. Raghunathan, A., Steinhardt, J., and Liang, P. S. *Semidefinite relaxations for certifying robustness to adversarial examples*. Neurips 2018.
3. Wong, E. and Kolter, Z. (2018). *Provable defenses against adversarial examples via the convex outer adversarial polytope*. ICML 2018.
4. Huang, X., Kwiatkowska, M., Wang, S., and Wu, M. *Safety verification of deep neural networks*. Computer Aided Verification 2017.
5. Athalye, A., et al. *Synthesizing robust adversarial examples*. International conference on machine learning. PMLR, 2018.
6. Eykholt, K., Evtimov, I., Fernandes, E., Li, B., Rahmati, A., Xiao, C., and Song, D. *Robust physical-world attacks on deep learning visual classification*. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 1625-1634). 2018.
7. Shafahi A., Ronny Huang, W., Studer, C., Feizi, S. and Goldstein, T. *Are adversarial examples inevitable?*. International Conference on Learning Representations. 2019.

Challenges in DL/ML applications: Robustness (References II)

1. Z. Charles, H. Rosenberg, and D. Papailiopoulos, *A Geometric Perspective on the Transferability of Adversarial Directions*, in Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics, Apr. 2019, pp. 1960–1968.
2. J. Gilmer, N. Ford, N. Carlini, and E. Cubuk, *Adversarial Examples Are a Natural Consequence of Test Error in Noise*, in Proceedings of the 36th International Conference on Machine Learning, May 2019, pp. 2280–2289.
3. A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry, *Adversarial Examples Are Not Bugs, They Are Features*, in Advances in Neural Information Processing Systems, 2019, vol. 32.
4. A. Fawzi, H. Fawzi, and O. Fawzi, *Adversarial vulnerability for any classifier*, in Advances in Neural Information Processing Systems, 2018, vol. 31.
5. L. Schmidt, S. Santurkar, D. Tsipras, K. Talwar, and A. Ma, *Adversarially Robust Generalization Requires More Data*, p. 13.
6. J.-H. Jacobsen, J. Behrmann, R. Zemel, and M. Bethge, *Excessive Invariance Causes Adversarial Vulnerability*, presented at the International Conference on Learning Representations, Sep. 2018.
7. F. Tramer, J. Behrmann, N. Carlini, N. Papernot, and J.-H. Jacobsen, *Fundamental Tradeoffs between Invariance and Sensitivity to Adversarial Perturbations*, in Proceedings of the 37th International Conference on Machine Learning, Nov. 2020, pp. 9561–9571.
8. C. Xie and A. Yuille, *Intriguing Properties of Adversarial Training at Scale*, presented at the International Conference on Learning Representations, Sep. 2019.

Challenges in DL/ML applications: Robustness (References III)

1. L. Chen, Y. Min, M. Zhang, and A. Karbasi, *More Data Can Expand the Generalization Gap Between Adversarially Robust and Standard Models*, p. 11.
2. F. Tramèr, N. Carlini, W. Brendel, and A. Ma, *On Adaptive Attacks to Adversarial Example Defenses*, p. 37.
3. D. Krueger et al., *Out-of-Distribution Generalization via Risk Extrapolation (REx)*, arXiv:2003.00688 [cs, stat], Feb. 2021,
4. T. Pang, K. Xu, Y. Dong, C. Du, N. Chen, and J. Zhu, *Rethinking Softmax Cross-Entropy Loss For Adversarial Robustness*, p. 19, 2020.
5. R. Bhattacharjee, S. Jha, and K. Chaudhuri, *Sample Complexity of Robust Linear Classification on Separated Data*, in Proceedings of the 38th International Conference on Machine Learning, Jul. 2021, pp. 884–893.
6. R. Geirhos et al., *Shortcut Learning in Deep Neural Networks*, Nat Mach Intell, vol. 2, no. 11, pp. 665–673, Nov. 2020, doi: 10.1038/s42256-020-00257-z.
7. Y. Min, L. Chen, and A. Karbasi, *The Curious Case of Adversarially Robust Models: More Data Can Help, Double Descend, or Hurt Generalization*, arXiv:2002.11080 [cs, stat], Jun. 2020
8. E. Rosenfeld, P. K. Ravikumar, and A. Risteski, *The Risks of Invariant Risk Minimization*, presented at the International Conference on Learning Representations, Sep. 2020.
9. T. Li, A. Beirami, M. Sanjabi, and V. Smith, *Tilted Empirical Risk Minimization*, presented at the International Conference on Learning Representations, Sep. 2020.
10. A. D'Amour et al., *Underspecification Presents Challenges for Credibility in Modern Machine Learning*, arXiv:2011.03395 [cs, stat], Nov. 2020.

Challenges in DL/ML applications: Surveillance/Privacy/Manipulation



Psychographics: the behavioural analysis that helped Cambridge Analytica know voters' minds

Prof. Michael Wade

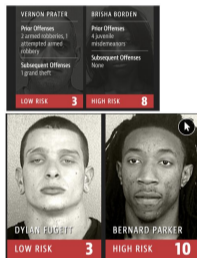


Figure: Political and societal concerns about some DL/ML applications

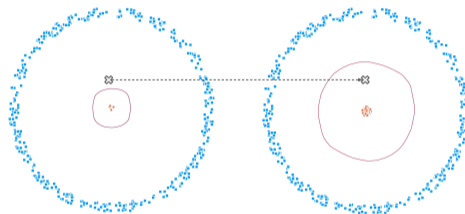
Challenges in DL/ML applications: Surveillance/Privacy/Manipulation (References)

1. Dwork, C., and Roth, A. *The Algorithmic Foundations of Differential Privacy*. Foundations and Trends in Theoretical Computer Science, 9, 2013.
2. Abadi, M., Chu, A., Goodfellow, I., McMahan, H. B., Mironov, I., Talwar, K., and Zhang, L. *Deep learning with differential privacy*. In Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security (pp. 308-318). 2016.
3. Sreenu, G., Saleem Durai, M.A. *Intelligent video surveillance: a review through deep learning techniques for crowd analysis*. J Big Data 6, 48. 2019.
4. O'Neil, C., *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy* Broadway Books, (2016);
5. Ali, R. E., So, J., Avestimehr, A. S. *On polynomial approximations for privacy-preserving and verifiable relu networks*. arXiv preprint arXiv:2011.05530.
6. Bagdasaryan, E., Poursaeed, O., Shmatikov, V. *Differential privacy has disparate impact on model accuracy*. In Neural Information Processing Systems (NeurIPS), 2019.
7. Bu, Z., Dong, J., Long, Q., Su, W. J. *Deep learning with Gaussian differential privacy*. Harvard data science review, 2020.
8. Pujol, D., McKenna, R., Kuppam, S., Hay, M., Machanavajjhala, A., Miklau, G. (textitFair decision making using privacy-protected data . In Proceedings of Fairness, Accountability, and Transparency (FAT), 2020.
9. Tran, C., Fioretto, F., Van Hentenryck, P. *Differentially private and fair deep learning: A lagrangian dual approach*. arXiv preprint arXiv:2009.12562.
10. Xu, D., Du, W., Wu, X. *Removing disparate impact on model accuracy in differentially private stochastic gradient descent*. In Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, 2021.

Challenges in DL/ML applications: Fairness



(a) Racist classifier



(b) Effect of unbalanced data

Figure: Unfair classifiers due to biased or unbalanced datasets/algorithms

Challenges in DL/ML applications: Fairness (References)

1. Barocas, S. Hardt, M. Narayanan, Arvind. *Fairness in Machine Learning Limitations and Opportunities*. <https://fairmlbook.org/pdf/fairmlbook.pdf> 2020.
2. Hardt, M. *How Big Data Is Unfair*. <https://medium.com/@mrtz/how-big-data-is-unfair-9aa544d739de> 2014.
3. Munoz, C., Smith, M., and Patil, D. *Big Data: A Report on Algorithmic Systems, Opportunity, and Civil Rights*. Executive Office of the President. The White House, 2016.
4. Campolo, A., Sanfilippo, M., Whittaker, M., Crawford, K. *AI Now 2017 Report*. AI Now Institute at New York University, 2017.
5. Friedman, B. and Nissenbaum, H. *Bias in Computer Systems*. ACM Transactions on Information Systems (TOIS) 14, no. 3. 1996: 330–47.
6. Pedreshi, D., Ruggieri, S. and Turini, F. *Discrimination-Aware Data Mining*. Proc. 14th SIGKDD. ACM 2008.
7. Noble, S.U. *Algorithms of Oppression: How Search Engines Reinforce Racism*. NYU Press. 2018.
8. Rolf, E., Simchowitz, M., Dean, S., Liu, L. T., Bjorkegren, D., Hardt, M., Blumenstock, J. *Balancing competing objectives with noisy data: Score-based classifiers for welfare-aware machine learning*. In International Conference on Machine Learning (ICML), 2021.
9. Hanna, A., Denton, E., Smart, A., Smith-Loud, J. *Towards a critical race methodology in algorithmic fairness*. In Proceedings of Fairness, Accountability, and Transparency (FAT), 2020.
10. Wang, A., Russakovsky, O. *Directional bias amplification*. In International Conference on Machine Learning (ICML), 2021.

Challenges in DL/ML applications: Fairness (References)

1. Yang, K., Qinami, K., Fei-Fei, L., Deng, J., Russakovsky, O. *Towards fairer datasets: Filtering and balancing the distribution of the people subtree in the imagenet hierarchy*. In Proceedings of Fairness, Accountability, and Transparency (FAT), 2020.
2. Mitchell, M., Baker, D., Moorosi, N., Denton, E., Hutchinson, B., Hanna, A., Morgenstern, J. *Diversity and inclusion metrics in subset selection*. In Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society, 2020.
3. Obermeyer, Z., Powers, B., Vogeli, C., Mullainathan, S. *Dissecting racial bias in an algorithm used to manage the health of populations*. Science, 366(6464), 447-453, 2019.
4. Jia, S., Meng, T., Zhao, J., Chang, K. W. *Mitigating gender bias amplification in distribution by posterior regularization*. Annual Meeting of the Association for Computational Linguistics (ACL), 2020.
5. Jalal, A., Karmalkar, S., Hoffmann, J., Dimakis, A., Price, E. *Fairness for Image Generation with Uncertain Sensitive Attributes*. In International Conference on Machine Learning (ICML), 2021.
6. Hoyle, A., Wallach, H., Augenstein, I., Cotterell, R. *Unsupervised discovery of gendered language through latent-variable modeling*. arXiv preprint arXiv:1906.04760.
7. Ramaswamy, V. V., Kim, S. S., Russakovsky, O. *Fair attribute classification through latent space de-biasing*. In Proceedings of the Conference on Computer Vision and Pattern Recognition, 2021.
8. Jacobs, A. Z., Wallach, H. *Measurement and fairness*. In Proceedings of Fairness, Accountability, and Transparency (FAT), 2021.

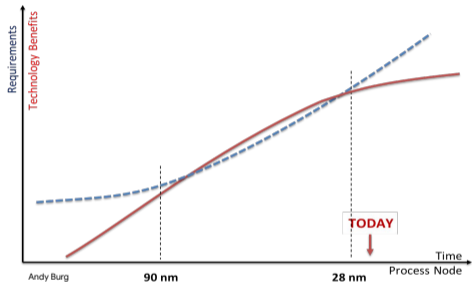
Challenges in DL/ML applications: Interpretability (References I)

1. Baehrens, David and Schroeter, Timon and Harmeling, Stefan and Kawanabe, Motoaki and Hansen, Katja and Mueller, Klaus-Robert. Simonyan, Karen and Vedaldi, Andrea and Zisserman, Andrew. *How to Explain Individual Classification Decisions*. JMLR 2010.
2. *Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps*. arXiv e-prints. arXiv:1312.6034. 2013.
3. Ribeiro, Marco and Singh, Sameer and Guestrin, Carlos. *Why Should I Trust You?": Explaining the Predictions of Any Classifier*. KDD 2016.
4. Sundararajan, Mukund and Taly, Ankur and Yan, Qiqi. *Axiomatic Attribution for Deep Networks*. ICML 2017.
5. Shrikumar, Avanti and Greenside, Peyton and Kundaje, Anshul. *Learning Important Features Through Propagating Activation Differences*. ICML 2017.
6. C. Rudin, C. Chen, Z. Chen, H. Huang, L. Semenova, and C. Zhong, *Interpretable Machine Learning: Fundamental Principles and 10 Grand Challenges* arXiv:2103.11251 [cs, stat], Jul. 2021
7. L. Semenova, C. Rudin, and R. Parr, *A study in Rashomon curves and volumes: A new perspective on generalization and model simplicity in machine learning* arXiv:1908.01755 [cs, stat], Apr. 2021

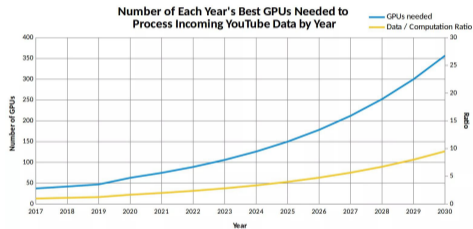
Challenges in DL/ML applications: Interpretability (References II)

1. S. Hooker, D. Erhan, P.-J. Kindermans, and B. Kim, *A Benchmark for Interpretability Methods in Deep Neural Networks* p. 12.
2. F. Poursabzi-Sangdeh, D. G. Goldstein, J. M. Hofman, J. W. Wortman Vaughan, and H. Wallach, *Manipulating and Measuring Model Interpretability* in Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems, New York, NY, USA: Association for Computing Machinery, 2021, pp. 1–52.
3. H. Kaur, H. Nori, S. Jenkins, R. Caruana, H. Wallach, and J. Wortman Vaughan, *Interpreting Interpretability: Understanding Data Scientists' Use of Interpretability Tools for Machine Learning* in Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems, New York, NY, USA, Apr. 2020, pp. 1–14. doi: 10.1145/3313831.3376219.
4. P. Hase and M. Bansal, *Evaluating Explainable AI: Which Algorithmic Explanations Help Users Predict Model Behavior?* in Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Online, 2020, pp. 5540–5552. doi: 10.18653/v1/2020.acl-main.491.
5. C. Rudin and J. Radin, *Why Are We Using Black Box Models in AI When We Don't Need To? A Lesson From An Explainable AI Competition* Harvard Data Science Review, vol. 1, no. 2, Nov. 2019, doi: 10.1162/99608f92.5a8a3a3d.
6. C. Rudin, *Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead* Nat Mach Intell, vol. 1, no. 5, pp. 206–215, May 2019, doi: 10.1038/s42256-019-0048-x.
7. S. Barocas, M. Hardt, and A. Narayanan, *Fairness in machine learning* Nips tutorial, vol. 1, p. 2017, 2017.
8. J. Lee, S. Park, and J. Shin, *Learning Bounds for Risk-sensitive Learning* arXiv:2006.08138 [cs, stat], Jan. 2021.

Challenges in DL/ML applications: Energy efficiency and cost



(a)



(b)

Figure: Efficiency and Scalability concerns in DL/ML

Challenges in DL/ML applications: Energy efficiency and cost (References)

1. García-Marín, E., Rodrigues, C. F., Riley, G., and Grahn, H. *Estimation of energy consumption in machine learning*. Journal of Parallel and Distributed Computing, 134, 75-88. 2019.
2. Strubell, E., Ganesh, A., and McCallum, A. *Energy and policy considerations for deep learning in NLP*. arXiv preprint arXiv:1906.02243. 2019.
3. Goel, A., Tung, C., Lu, Y. H., and Thiruvathukal, G. K. *A Survey of Methods for Low-Power Deep Learning and Computer Vision*. arXiv preprint arXiv:2003.11066. 2020.
4. Conti, F., Rusci, M., and Benini, L. *The Memory Challenge in Ultra-Low Power Deep Learning*. In NANO-CHIPS 2030 (pp. 323-349). Springer, Cham. 2020.
5. J. Launay, I. Poli, F. Boniface, and F. Krzakala, *Direct Feedback Alignment Scales to Modern Deep Learning Tasks and Architectures* NeurIPS 2020
6. J. Launay et al., *Hardware Beyond Backpropagation: a Photonic Co-Processor for Direct Feedback Alignment* arXiv:2012.06373
7. E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell, *On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?* in Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, New York, NY, USA, Mar. 2021, pp. 610–623. doi: 10.1145/3442188.3445922.

Wrap up!

- Learning deep continues!

References I

- [1] Zeyuan Allen-Zhu.
Natasha: Faster non-convex stochastic optimization via strongly non-convex parameter.
In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 89–97. PMLR, 06–11 Aug 2017.
(Cited on page 21.)
- [2] Zeyuan Allen-Zhu.
Katyusha x: Practical momentum method for stochastic sum-of-nonconvex optimization.
In *ICML*, 2018.
(Cited on page 20.)
- [3] Zeyuan Allen-Zhu.
Katyusha x: Simple momentum method for stochastic sum-of-nonconvex optimization.
In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 179–185. PMLR, 10–15 Jul 2018.
(Cited on page 17.)
- [4] Anish Athalye, Nicholas Carlini, and David Wagner.
Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples.
In *International Conference on Machine Learning*, pages 274–283. PMLR, 2018.
(Cited on page 48.)

References II

- [5] George Cybenko.
Approximation by superpositions of a sigmoidal function.
Mathematics of control, signals and systems, 2(4):303–314, 1989.
(Cited on page 34.)
- [6] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien.
Saga: A fast incremental gradient method with support for non-strongly convex composite objectives.
In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 1646–1654. Curran Associates, Inc., 2014.
(Cited on pages 8, 9, and 14.)
- [7] Benjamin Dubois-Taine, Sharan Vaswani, Reza Babanezhad, Mark Schmidt, and Simon Lacoste-Julien.
Svrg meets adagrad: Painless variance reduction, 2021.
(Cited on page 20.)
- [8] John Duchi, Elad Hazan, and Yoram Singer.
Adaptive subgradient methods for online learning and stochastic optimization.
J. Mach. Learn. Res., 12:2121–2159, July 2011.
(Cited on page 20.)

References III

- [9] Kevin Eykholt, Ivan Evtimov, Earlene Fernandes, Bo Li, Amir Rahmati, Chaowei Xiao, Atul Prakash, Tadayoshi Kohno, and Dawn Song.
Robust physical-world attacks on deep learning visual classification.
In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1625–1634, 2018.
(Cited on page 48.)
- [10] Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang.
SPIDER: near-optimal non-convex optimization via stochastic path-integrated differential estimator.
In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 687–697, 2018.
(Cited on pages 17, 18, 20, and 21.)
- [11] Rie Johnson and Tong Zhang.
Accelerating stochastic gradient descent using predictive variance reduction.
In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 315–323. Curran Associates, Inc., 2013.
(Cited on page 11.)

References IV

- [12] Andrej Karpathy.
What I learned from competing against a ConvNet on ImageNet.
(Cited on page 39.)
- [13] Ali Kavis, Stratis Skoulakis, Kimon Antonakopoulos, Leello Tadesse Dadi, and Volkan Cevher.
Adaptive stochastic variance reduction for non-convex finite-sum minimization.
Advances In Neural Information Processing Systems 36 (NeurIPS 2022), (CONF), 2022.
(Cited on pages 19 and 21.)
- [14] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein.
Visualizing the Loss Landscape of Neural Nets.
In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
(Cited on page 45.)
- [15] Zhize Li, Hongyan Bao, Xiangliang Zhang, and Peter Richtarik.
Page: A simple and optimal probabilistic gradient estimator for nonconvex optimization.
In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 6286–6295. PMLR, 18–24 Jul 2021.
(Cited on page 17.)

References V

- [16] Sashank J Reddi, Suvrit Sra, Barnabás Póczos, and Alex Smola.
Stochastic frank-wolfe methods for nonconvex optimization.
arXiv preprint arXiv:1607.08254, 2016.
(Cited on page 17.)
- [17] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al.
Imagenet large scale visual recognition challenge.
115(3):211–252, 2015.
(Cited on page 39.)
- [18] Ali Shafahi, W Ronny Huang, Christoph Studer, Soheil Feizi, and Tom Goldstein.
Are adversarial examples inevitable?
In *7th International Conference on Learning Representations (ICLR 2019)*, 2019.
(Cited on page 49.)
- [19] Zhe Wang, Kaiyi Ji, Yi Zhou, Yingbin Liang, and Vahid Tarokh.
Spiderboost and momentum: Faster stochastic variance reduction algorithms.
In *Advances in Neural Information Processing Systems*, 2019.
(Cited on page 17.)

References VI

- [20] Zhe Wang, Kaiyi Ji, Yi Zhou, Yingbin Liang, and Vahid Tarokh.
Spiderboost and momentum: Faster variance reduction algorithms.
In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alche Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
(Cited on page 20.)
- [21] Lin Xiao and Tong Zhang.
A proximal stochastic gradient method with progressive variance reduction.
SIAM Journal on Optimization, 24, 03 2014.
(Cited on page 11.)
- [22] Guangzeng Xie, Luo Luo, and Zhihua Zhang.
A general analysis framework of lower complexity bounds for finite-sum optimization, 2019.
(Cited on page 17.)
- [23] Dongruo Zhou and Quanquan Gu.
Lower bounds for smooth nonconvex finite-sum optimization.
In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7574–7583. PMLR, 09–15 Jun 2019.
(Cited on pages 17 and 21.)