

# CMPT 409/981: Optimization for Machine Learning

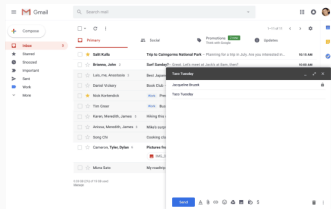
## Lecture 1

---

Sharan Vaswani

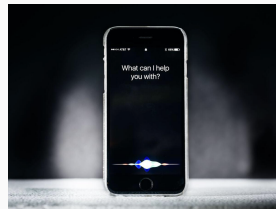
September 5, 2024

# Successes of Machine Learning



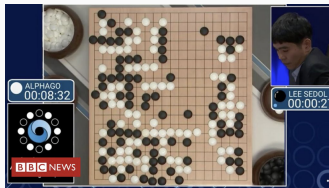
<https://www.blog.google/products/gmail/subject-write-emails-faster-smart-compose-gmail/>

(a) Natural language processing



<https://www.cnet.com/news/what-is-siri/>

(b) Speech recognition



<https://www.bbc.com/news/technology-35785875>

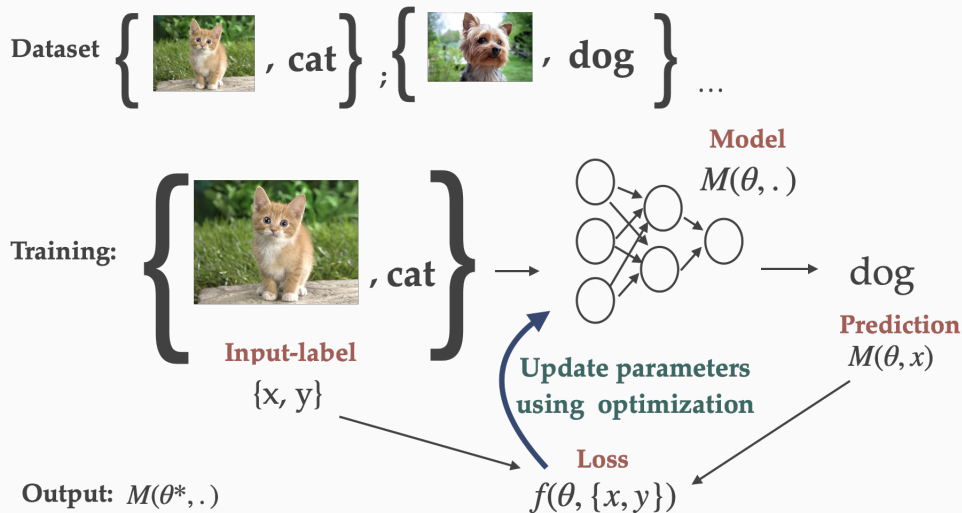
(c) Reinforcement learning



<https://www.pbs.org/newshour/science/in-a-crash-should-self-driving-cars-save-passengers-or-pedestrians-2-million-people-weigh-in>

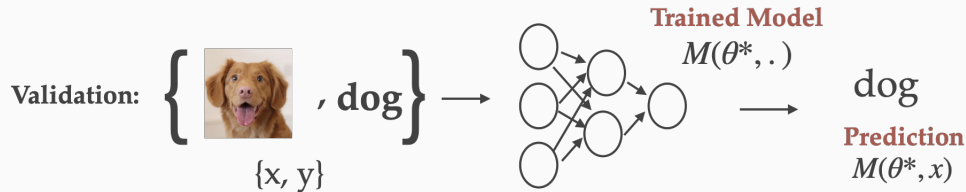
(d) Self-driving cars

# Machine Learning 101



# Machine Learning 101

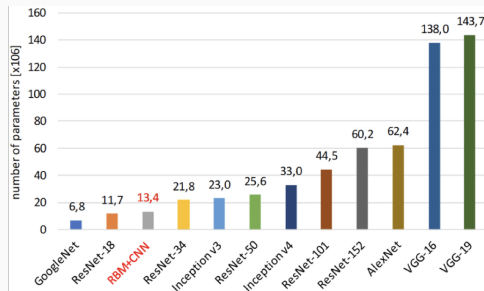
**Validation Dataset:**  $\left\{ \text{cat image}, \text{cat} \right\}; \left\{ \text{dog image}, \text{dog} \right\} \dots$



**Output:** Validation Accuracy

Measures how good the trained model is

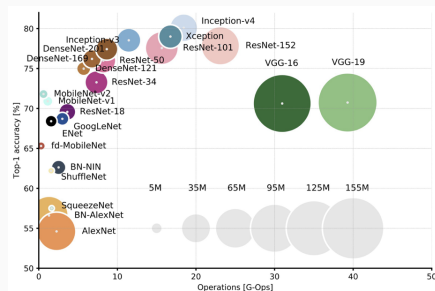
# Modern Machine Learning



Sobczak, Szymon, et al. "Restricted Boltzmann machine as an aggregation technique for binary descriptors.", 2019.

Model size

(a)



Canziani et al, "An Analysis of Deep Neural Network Models for Practical Applications", 2016.

Number of operations for computing the loss

(b)

Figure 1: Models for multi-class classification on Image-Net. Number of examples = 1.2 M

Faster optimization methods can have a big practical impact!

- **(Non)-Convex minimization:** Supervised learning (classification/regression), Matrix factorization for recommender systems, Image denoising.
- **Online optimization:** Learning how to play Go/Atari games, Imitating an expert and learning from demonstrations, Regulating control systems like industrial plants.
- **Min-Max optimization:** Generative Adversarial Networks, Adversarial Learning, Multi-agent RL.

**Objective:** Introduce foundational optimization concepts with applications to machine learning.

**Syllabus:**

- **(Non)-Convex minimization:** Gradient Descent, Momentum/Acceleration, Mirror Descent, Newton/Quasi-Newton methods, Stochastic gradient descent (SGD), Variance reduction
- **Online optimization:** Follow the (regularized) leader, Adaptive methods (AdaGrad, Adam)
- **Min-Max optimization:** (Stochastic) Gradient Descent-Ascent, (Stochastic) Extragradient

**What we won't get time to cover in detail:** Non-smooth optimization, Convex analysis, Global optimization.

**What we won't get time to cover:** Constrained optimization, Distributed optimization, Multi-objective optimization.

- **Instructor:** Sharan Vaswani (TASC-1 8221) Email: [sharan\\_vaswani@sfu.ca](mailto:sharan_vaswani@sfu.ca)
- **Instructor Office Hours:** Thursday, 2.30 pm - 3.30 pm (TASC-1 8221)
- **Teaching Assistant:** Qiushi Lin Email: [qla96@sfu.ca](mailto:qla96@sfu.ca)
- **TA Office Hours:** Monday, 9.30 am - 10.30 am (ASB 9814)
- **Course Webpage:** [https://vaswanis.github.io/409\\_981-F24.html](https://vaswanis.github.io/409_981-F24.html)
- **Piazza:** <https://piazza.com/sfu.ca/fall2024/cmpt409981/home>
- **Prerequisites:** Linear Algebra, Multivariable calculus, (Undergraduate) Machine Learning



## Assignments [48%]

- Individual assignments to be submitted online, typed up in Latex with accompanying code submitted as a zip file.
- **Assignment 0** [5%]: Out today. Assignment to recall prerequisite knowledge and get used to notation. Due next week.
- **Assignments 1 & 2** [21%]:
  - Due in 10 days (at 11.59 pm PST).
  - For some flexibility, each student is allowed 1 late-submission and can submit in the next class (no late submissions beyond that).
  - If you use up your late-submission and submit late again, you will lose 50% of the mark.
- **Assignments 3 & 4** [21%]: Released during the semester, but due only at the end of the term (December 10).

**Participation** [2%]: In class (during lectures, project presentations), on Piazza

## Final Project [50%]

- Aim is to give you a taste of research in Optimization.
- Projects to be done in groups of 3-4 (more details will be on Piazza)
- Will maintain a list on Piazza on possible project topics. You are free to choose from the list or propose a topic that combines Optimization with your own research area.
- Project Proposal [10%] – Discussion (before October 20) + Report (due October 22)
- Project Milestone [5%] – Update (before November 20)
- Project Presentation [10%] (December 3)
- Project Report [25%] (December 17)

Questions?

# Minimizing functions

Consider minimizing a function over the domain  $\mathcal{D}$

$$\min_{w \in \mathcal{D}} f(w).$$

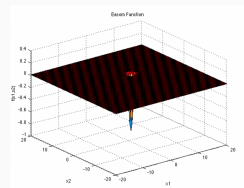
**Setting:** Have access to a *zero-order oracle* – querying the oracle at  $w \in \mathcal{D}$  returns  $f(w)$ .

**Objective:** For a target accuracy of  $\epsilon > 0$ , if  $f^*$  is the minimum value of  $f$  in  $\mathcal{D}$ , return a point  $\hat{w} \in \mathcal{D}$  s.t.  $f(\hat{w}) - f^* \leq \epsilon$ . Characterize the required number of oracle calls in terms of  $\epsilon$ .

*Example 1:* Minimize a one-dimensional function s.t.  $f(w) = 0$  for all  $x \neq w^*$ , and  $f(w^*) = -\epsilon$ .

*Example 2:* Easom function:

$$f(x_1, x_2) = -\cos(x_1) - \cos(x_2) \exp(-(x_1 - \pi)^2 - (x_2 - \pi)^2).$$



Minimizing generic functions is hard! We need to make assumptions on the structure.

# Lipschitz continuous functions

Consider minimizing a function over the domain  $\mathcal{D}$ :

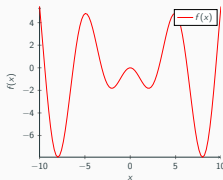
$$\min_{w \in \mathcal{D}} f(w).$$

**Assumption:**  $f$  is *Lipschitz continuous* (in  $\mathcal{D}$ ) meaning that  $f$  can not change arbitrarily fast as  $w$  changes. Formally, for any  $x, y \in \mathcal{D}$ ,

$$|f(x) - f(y)| \leq G \|x - y\|$$

where  $G$  is the Lipschitz constant.

*Example:*  $f(x) := -x \sin(x)$  in the  $[-10, 10]$  interval.



Lipschitz continuity of the function immediately implies that the gradients are *bounded* i.e. for all  $x \in \mathcal{D}$ ,  $\|\nabla f(x)\| \leq G$ .

# Global Minimization

Consider minimizing a  $G$ -Lipschitz continuous function over a unit hyper-cube:

$$\min_{w \in [0,1]^d} f(w).$$

**Objective:** For a target accuracy of  $\epsilon > 0$ , if  $w^* \in [0,1]^d$  is the minimizer of  $f$ , return a point  $\hat{w} \in [0,1]^d$  s.t.  $f(\hat{w}) - f(w^*) \leq \epsilon$ . Characterize the required number of zero-order oracle calls.

**Naive algorithm:** Divide the hyper-cube into cubes with length of each side equal to  $\epsilon' > 0$  (to be determined). Call the zero-order oracle on the centers of these  $\frac{1}{(\epsilon')^d}$  cubes and return the point  $\hat{w}$  with the minimum function value.

**Analysis:** The minimizer lies in/at the boundary of one of these cubes, and hence by returning the minimum  $\hat{w}$ , we guarantee that  $\hat{w}$  is at most  $\sqrt{d}\epsilon'$  away from  $w^*$  i.e.  $\|\hat{w} - w^*\| \leq \sqrt{d}\epsilon'$ . By  $G$ -Lipschitz continuity,  $f(\hat{w}) - f(w^*) \leq G \|\hat{w} - w^*\| \leq G\sqrt{d}\epsilon'$ . For a target accuracy of  $\epsilon$ , we can set  $\epsilon' = \frac{\epsilon}{G\sqrt{d}}$ . Hence, for this naive algorithm, total number of oracle calls =  $\left(\frac{G\sqrt{d}}{\epsilon}\right)^d$ .

Consider minimizing a differentiable,  $G$ -Lipschitz continuous function over a unit hyper-cube:

$$\min_{w \in [0,1]^d} f(w).$$

**Q:** Suppose we do a random search over the cubes – choosing a cube at random and then querying its centre? What is the expected number of function evaluations?

Is our naive algorithm good? Can we do better?

**Lower-Bound:** For minimizing a  $G$ -Lipschitz continuous function over a unit hyper-cube, any algorithm requires  $\Omega\left(\left(\frac{G}{\epsilon}\right)^d\right)$  calls to the zero-order oracle.

Our naive-algorithm is *sub-optimal* by a factor of  $O\left((\sqrt{d})^d\right)$ .

Questions?



# Smooth functions

Recall that Lipschitz continuous functions have bounded gradients i.e.  $\|\nabla f(w)\| \leq G$  and can still include *non-smooth* (not differentiable everywhere) functions.

For example,  $f(x) = |x|$  is 1-Lipschitz continuous but not differentiable at  $x = 0$  and the gradient changes from  $-1$  at  $0^-$  to  $+1$  at  $0^+$ .

An alternative assumption that we can make is that  $f$  is *smooth* – it is differentiable everywhere and its gradient is Lipschitz-continuous i.e. it can not change arbitrarily fast.

Formally, the gradient  $\nabla f$  is  $L$ -Lipschitz continuous if for all  $x, y \in \mathcal{D}$ ,

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$$

where  $L$  is the Lipschitz constant of the gradient (also called the smoothness constant of  $f$ ).

**Q:** Does Lipschitz-continuity of the gradient imply Lipschitz-continuity of the function?

If  $f$  is twice-differentiable and smooth, then for all  $x \in \mathcal{D}$ ,  $\nabla^2 f(x) \preceq L I_d$  i.e.  $\sigma_{\max}[\nabla^2 f(x)] \leq L$  where  $\sigma_{\max}$  is the maximum singular value.

Q: Does  $f(x) = x^3$  have a Lipschitz-continuous gradient over  $\mathbb{R}$ ?

Q: Does  $f(x) = x^3$  have a Lipschitz-continuous gradient over  $[0, 1]$ ?

Q: The *negative entropy function* is given by  $f(x) = x \log(x)$ . Does it have a Lipschitz-continuous gradient over  $[0, 1]$ ?

## Smooth functions – Examples

**Linear Regression** on  $n$  points with  $d$  features. Feature matrix:  $X \in \mathbb{R}^{n \times d}$ , vector of measurements:  $y \in \mathbb{R}^n$  and parameters  $w \in \mathbb{R}^d$ .

$$\min_{w \in \mathbb{R}^d} f(w) := \frac{1}{2} \|Xw - y\|^2$$

$$f(w) = \frac{1}{2} [w^\top (X^\top X) w - 2w^\top X^\top y + y^\top y] ; \nabla f(w) = X^\top X w - X^\top y ; \nabla^2 f(w) = X^\top X$$

(Prove in Assignment 0)

If  $f$  is  $L$ -smooth, then,  $\sigma_{\max}[\nabla^2 f(w)] \leq L$  for all  $w$ . Hence, for linear regression  $L = \lambda_{\max}[X^\top X]$ .

**Q:** Is the linear regression loss-function Lipschitz continuous?

**Q:** Compute  $L$  for *ridge regression* –  $\ell_2$ -regularized linear regression where  $f(w) := \frac{1}{2} \|Xw - y\|^2 + \frac{\lambda}{2} \|w\|^2$ .

# Smooth functions

**Claim:** For an  $L$ -smooth function,  $f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2$  for all  $x, y \in \mathcal{D}$ .

**Proof:**

$$f(y) = f(x) + \int_{t=0}^1 [\nabla f(x + t(y - x))] (y - x)^\top dt \quad (\text{Fundamental theorem of calculus})$$

$$= f(x) + \langle \nabla f(x), y - x \rangle + \int_{t=0}^1 [\nabla f(x + t(y - x))] (y - x)^\top dt - [\nabla f(x)] (y - x)^\top$$

$$= f(x) + \langle \nabla f(x), y - x \rangle + \int_{t=0}^1 [\nabla f(x + t(y - x)) - \nabla f(x)] (y - x)^\top dt$$

$$\leq f(x) + \langle \nabla f(x), y - x \rangle + \int_{t=0}^1 \|\nabla f(x + t(y - x)) - \nabla f(x)\| \|y - x\| dt$$

(Cauchy-Schwarz)

$$\leq f(x) + \langle \nabla f(x), y - x \rangle + L \int_{t=0}^1 \|x + t(y - x) - x\| \|y - x\| dt \quad (\text{Lipschitz continuity})$$

$$= f(x) + \langle \nabla f(x), y - x \rangle + L \|y - x\|^2 \int_{t=0}^1 t dt = f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2$$

The inequality  $f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2$  can be interpreted as a *global* quadratic upper-bound on  $f$  at point  $x$  i.e. the bound holds for all  $y \in \mathcal{D}$ .

There are other related (not always equivalent) ways to state the  $L$ -smoothness of  $f$  (you will need to prove these in Assignment 1).

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{1}{2L} \|\nabla f(y) - \nabla f(x)\|^2$$
$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \leq L \|x - y\|^2$$

Questions?

# Local Minimization

Smooth functions can include functions with multiple local/global minimum and stationary points. Eg:  $f(x) = -x \sin(x)$ .

Consider minimizing a smooth function over  $\mathbb{R}^d$  (unconstrained minimization)

$$\min_{w \in \mathbb{R}^d} f(w).$$

Since we have seen that global minimization can be impossible (without Lipschitz assumption on  $f$ ) or the number of oracle calls can be exponential in  $d$ , let us aim to solve an easier problem.

Access to a **first-order oracle** – query the oracle at point  $w$  and it returns  $f(w)$  and  $\nabla f(w)$ .

**Objective:** For a target accuracy of  $\epsilon > 0$ , return a point  $\hat{w}$  s.t.  $\|\nabla f(\hat{w})\|^2 \leq \epsilon$ ? Characterize the required number of oracle calls.

We only care about making the gradient small and finding an approximate stationary point.

## Local Minimization

Recall that  $L$ -smoothness of  $f$  implies that  $f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2$ .

**Idea:** Since the RHS is a global upper-bound on the true function, instead of minimizing the function directly, let us minimize the upper-bound at  $x$  w.r.t  $y$ .

Setting the gradient of the RHS w.r.t  $y$  to zero, we obtain  $\hat{y}$  as:

$$\nabla f(x) + L [\hat{y} - x] = 0 \implies \hat{y} = x - \frac{1}{L} \nabla f(x)$$

This is exactly the gradient descent update at  $x$ !

We can do this iteratively i.e. starting at  $w_0$ , form the upper-bound at  $w_0$ , minimize it by setting  $w_1 = w_0 - \frac{1}{L} \nabla f(w_0)$ , then form the quadratic upper-bound at  $w_1$  and repeat. Continue to do this until we find a point  $\hat{w}$  s.t.  $\|\nabla f(\hat{w})\|^2 \leq \epsilon$  and terminate.

This is exactly the gradient descent procedure – move in the direction of the negative gradient (“downhill”) with *step-size*  $\eta$  equal to  $1/L$ . Formally, at iteration  $k$ , the GD update is:

$$w_{k+1} = w_k - \eta \nabla f(w_k).$$



Is GD guaranteed to terminate? If so, can we characterize the number of iterations?

**Claim:** For  $L$ -smooth functions lower-bounded by  $f^*$ , gradient descent with  $\eta = \frac{1}{L}$  returns  $\hat{w}$  such that  $\|\nabla f(\hat{w})\|^2 \leq \epsilon$  and requires  $T = \frac{2L[f(w_0) - f^*]}{\epsilon}$  iterations (oracle calls).

**Proof:**

Using the  $L$ -smoothness of  $f$  with  $x = w_k$  and  $y = w_{k+1} = w_k - \frac{1}{L}\nabla f(w_k)$  in the quadratic bound (also referred to as the *descent lemma*),

$$\begin{aligned} f(w_{k+1}) &\leq f(w_k) + \langle \nabla f(w_k), -\frac{1}{L}\nabla f(w_k) \rangle + \frac{L}{2} \left\| \frac{1}{L}\nabla f(w_k) \right\|^2 \\ \implies f(w_{k+1}) &\leq f(w_k) - \frac{1}{2L} \|\nabla f(w_k)\|^2 \end{aligned}$$

By moving from  $w_k$  to  $w_{k+1}$ , we have decreased the value of  $f$  since  $f(w_{k+1}) \leq f(w_k)$ .

# Gradient Descent

Rearranging the inequality from the previous slide, for every iteration  $k$ ,

$$\frac{1}{2L} \|\nabla f(w_k)\|^2 \leq f(w_k) - f(w_{k+1})$$

By running GD for  $T$  iterations, adding up  $k = 0$  to  $T - 1$ ,

$$\frac{1}{2L} \sum_{k=0}^{T-1} \|\nabla f(w_k)\|^2 \leq \sum_{k=0}^{T-1} [f(w_k) - f(w_{k+1})] = f(w_0) - f(w_T) \leq [f(w_0) - f^*]$$

(Since  $f$  is lower-bound by  $f^*$ )

$$\implies \frac{\sum_{k=0}^{T-1} \|\nabla f(w_k)\|^2}{T} \leq \frac{2L[f(w_0) - f^*]}{T}$$

The LHS is the average of the gradient norms over the  $T$  iterates. Let

$\hat{w} := \arg \min_{k \in \{0, 1, \dots, T-1\}} \|\nabla f(w_k)\|^2$ . Since the minimum is smaller than the average,

$$\|\nabla f(\hat{w})\|^2 \leq \frac{2L[f(w_0) - f^*]}{T}$$

Since  $\|\nabla f(\hat{w})\|^2 \leq \frac{2L[f(w_0)-f^*]}{T}$ , the *rate of convergence* is  $O(1/T)$ .

If the RHS equal to  $\frac{2L[f(w_0)-f^*]}{T} \leq \epsilon$ , this would guarantee that  $\|\nabla f(\hat{w})\|^2 \leq \epsilon$  and we would achieve our objective.

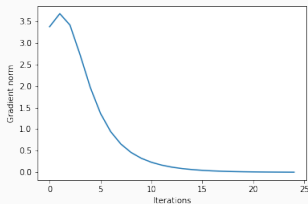
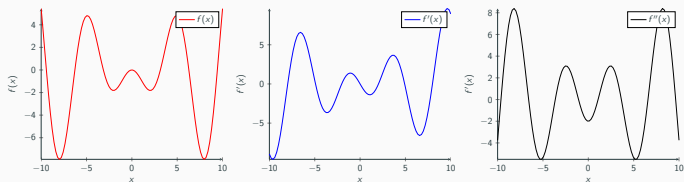
Hence, we need to run the algorithm for  $T \geq \frac{2L[f(w_0)-f^*]}{\epsilon}$  iterations. This is also referred to as an  $O\left(\frac{1}{\epsilon}\right)$  convergence rate.

**Lower-Bound:** When minimizing a smooth function (without additional assumptions), any *first-order* algorithm requires  $\Omega\left(\frac{1}{\epsilon}\right)$  oracle calls to return a point  $\hat{w}$  such that  $\|\nabla f(\hat{w})\|^2 \leq \epsilon$ .

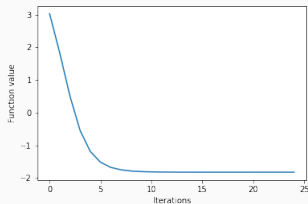
Hence, gradient descent is optimal for minimizing smooth functions!

# Gradient Descent – Example

$\min_{x \in [-10, 10]} f(x) := -x \sin(x)$ . Run GD with  $\eta = 1/L \approx 0.1$  and  $x_0 = 4$ .



(a) Gradient norm



(b) Function value

Questions?