

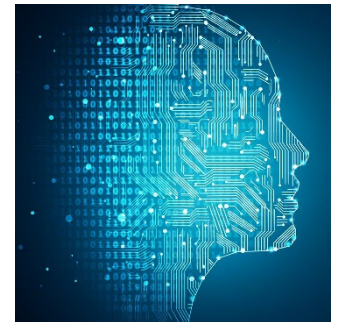
# Machine Learning

# Unconstrained Optimization



Kevin Moon (kevin.moon@usu.edu)

STAT 6655



# Outline



1. Minimums and necessary conditions
2. Convexity
3. Methods for solving optimization problems

# Optimization in machine learning



- Many (all?) machine learning problems can be posed as a minimization or maximization problem
  - Empirical risk minimization
- In some cases a closed form solution exists
  - E.g. linear regression
- Most of the time, a closed form solution doesn't exist
- We can solve these problems using optimization theory

# Unconstrained Optimization



- An *unconstrained optimization problem* has the form

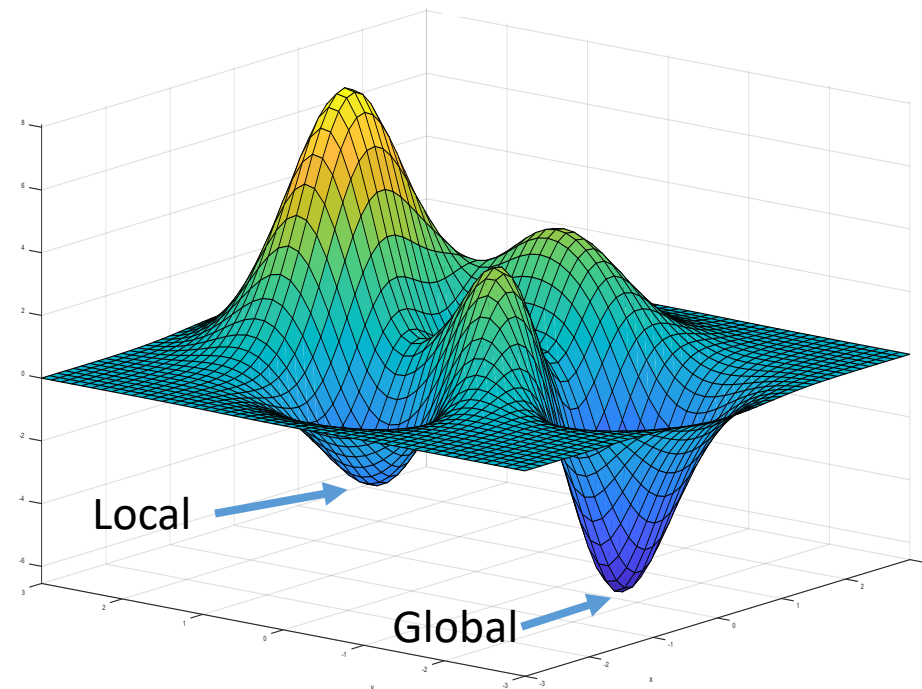
$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$$

$\exists$  = “there exists”

$\forall$  = “for all”

where  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is called the *objective function*.

- What about maximization?
- A point  $\mathbf{x}^* \in \mathbb{R}^d$  is called a *local minimizer* if  $\exists r > 0$  such that  $f(\mathbf{x}^*) \leq f(\mathbf{x}) \forall \mathbf{x}$  satisfying  $\|\mathbf{x} - \mathbf{x}^*\| < r$
- $\mathbf{x}^*$  is called a *global minimizer* if  $f(\mathbf{x}^*) \leq f(\mathbf{x}) \forall \mathbf{x} \in \mathbb{R}^d$
- Is a global minimizer also a local minimizer?



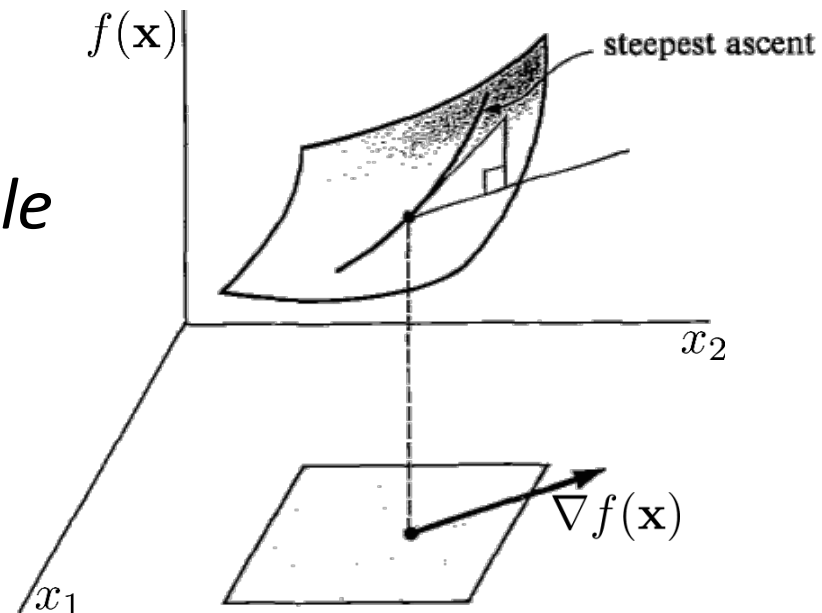
# Gradient



- Given a function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ , the *gradient* of  $f$  at  $\mathbf{x} = [x_1 \dots x_d]^T \in \mathbb{R}^d$  is defined by

$$\nabla f(\mathbf{x}) := \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_d} \end{bmatrix}$$

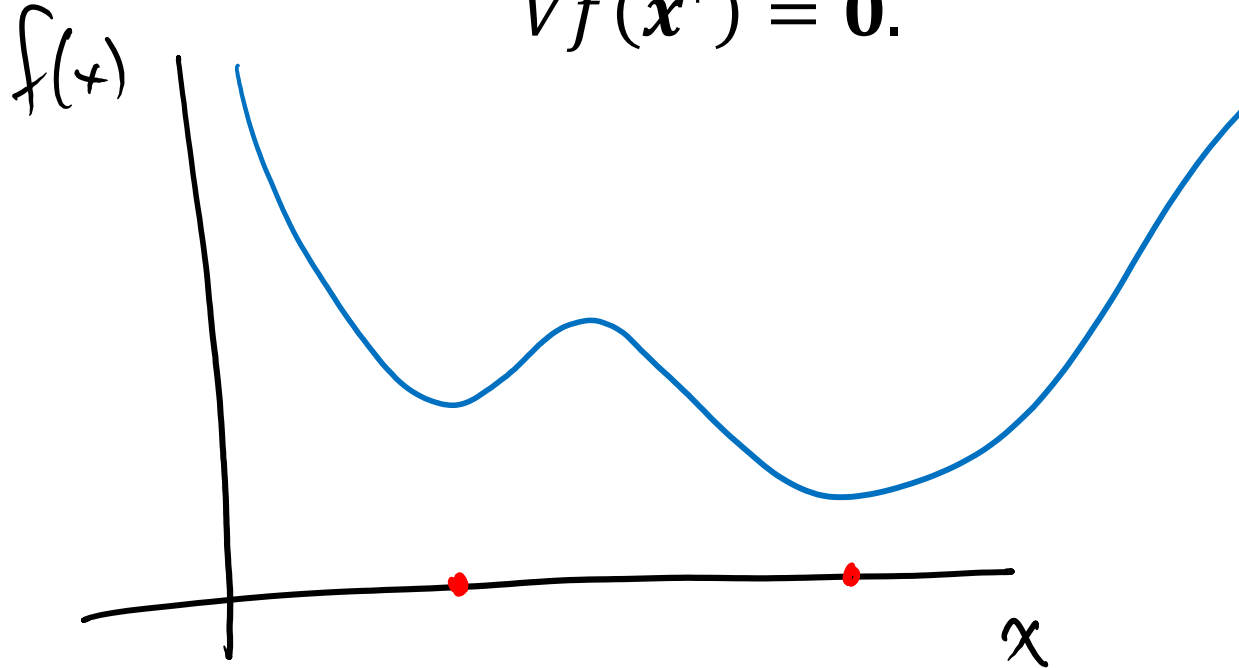
- If the gradient exists for all  $\mathbf{x} \in \mathbb{R}^d$ , we say  $f$  is *differentiable*
- The gradient gives the direction of *steepest ascent*



# First Order Necessary Condition



- If  $f$  is differentiable and  $\mathbf{x}^*$  is a local minimizer of  $f$ , then  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ .



- Note that this condition is necessary but not sufficient for  $\mathbf{x}^*$  to be a local minimizer
  - Why?
- If  $\nabla f(\mathbf{x}) = 0$  for some  $\mathbf{x}$ , then  $\mathbf{x}$  is said to be a critical point or a stationary point of  $f$

# First Order Necessary Condition



- *Proof:* Define the scalar valued function  $\phi(t) = f(\mathbf{x}^* + \mathbf{y}t)$ , where  $\mathbf{y} \in \mathbb{R}^d$  is arbitrary. Then

$$\begin{aligned}\phi'(0) &= \lim_{t \searrow 0} \frac{f(\mathbf{x}^* + \mathbf{y}t) - f(\mathbf{x}^*)}{t} \\ &= \langle \nabla f(\mathbf{x}^*), \mathbf{y} \rangle\end{aligned}$$

by the chain rule. Since  $\mathbf{x}^*$  is a local min, we know

$$f(\mathbf{x}^* + \mathbf{y}t) \geq f(\mathbf{x}^*)$$

for  $t$  sufficiently small. Therefore,  $\langle \nabla f(\mathbf{x}^*), \mathbf{y} \rangle \geq 0$ . Now choose  $\mathbf{y} = -\nabla f(\mathbf{x}^*)$ . Then

$$0 \leq \langle \nabla f(\mathbf{x}^*), -\nabla f(\mathbf{x}^*) \rangle = -\|\nabla f(\mathbf{x}^*)\|^2 \leq 0,$$

so we must have  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ .

# Hessian



- The *Hessian* of  $f$  at  $\mathbf{x}$  is the  $d \times d$  matrix

$$\nabla^2 f(\mathbf{x}) := \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_d \partial x_1} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_d^2} \end{bmatrix}$$

- We say that  $f$  is *twice differentiable* if  $\nabla^2 f(\mathbf{x})$  exists  $\forall \mathbf{x} \in \mathbb{R}^d$ .
- We say  $f$  is *twice continuously differentiable* if it is twice differentiable and all of the second derivatives are continuous functions of  $\mathbf{x}$ .
- If  $f$  is twice continuously differentiable, then  $\nabla^2 f(\mathbf{x})$  is a symmetric matrix  $\forall \mathbf{x}$ , i.e.,

$$\begin{aligned} \frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j} &= \frac{\partial^2 f(\mathbf{x})}{\partial x_j \partial x_i} \\ &\forall \mathbf{x} \in \mathbb{R}^d \\ &\forall i, j = 1, \dots, d \end{aligned}$$



# Positive (Semi-)Definite Matrices



- Let  $A$  be a  $d \times d$  matrix. We say that  $A$  is *positive definite* (PD) if  $\mathbf{x}^T A \mathbf{x} > 0$  for all  $\mathbf{x} \neq \mathbf{0}$ . We say that  $A$  is *positive semi-definite* (PSD) if  $\mathbf{x}^T A \mathbf{x} \geq 0$  for all  $\mathbf{x} \in \mathbb{R}^d$ .
- PD and PSD arise frequently in ML, for example
  - Gram matrices
  - Kernel matrices
  - Covariance matrices
  - Hessian matrices (sometimes)
- PD/PSD matrices are not necessarily symmetric, e.g.
$$\begin{bmatrix} x \\ y \end{bmatrix} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x & y \end{bmatrix} = x^2 + y^2$$
  - However, we will only consider PD/PSD matrices that are also symmetric

# Second Order Necessary Condition



- If  $f$  is twice continuously differentiable and  $\mathbf{x}^*$  is a local min, then  $\nabla^2 f(\mathbf{x}^*)$  is positive semi-definite, i.e.,

$$\mathbf{z}^T \nabla^2 f(\mathbf{x}^*) \mathbf{z} \geq 0, \quad \forall \mathbf{z} \in \mathbb{R}^d$$

- This generalizes the result from single-variable calculus that the second derivative is nonnegative at a local min
- *Proof:*  
From the definition of local optimality, there exists a neighborhood  $A$  of  $\mathbf{x}^*$  such that  $f(\mathbf{x}^*)$  is the minimum inside  $A$ .
  - I.e.,  $\exists r > 0$  such that  $f(\mathbf{x}^*) \leq f(\mathbf{x}) \forall \mathbf{x}$  satisfying  $\|\mathbf{x} - \mathbf{x}^*\| < r$  and  $A = \{\mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x} - \mathbf{x}^*\| < r\}$

# Proof continued



- By multidimensional Taylor series expansion, we can write for any  $\mathbf{y}$  and  $t$  such that  $\mathbf{x}^* + t\mathbf{y} \in A$ :

$$\begin{aligned} f(\mathbf{x}^* + t\mathbf{y}) \\ = f(\mathbf{x}^*) + \langle \nabla f(\mathbf{x}^*), t\mathbf{y} \rangle + \frac{t^2}{2} \langle \mathbf{y}, \nabla^2 f(\mathbf{x}^*) \mathbf{y} \rangle + o(t^2 \|\mathbf{y}\|^2) \end{aligned}$$

- $o(t)$  denotes a function satisfying  $\lim_{t \rightarrow 0} \frac{o(t)}{t} = 0$
- Noting that  $\nabla f(\mathbf{x}^*) = 0$  and rearranging gives for  $\mathbf{y} \neq 0$ :
$$0 \leq \frac{f(\mathbf{x}^* + t\mathbf{y}) - f(\mathbf{x}^*)}{t^2 \|\mathbf{y}\|^2} = \frac{o(t^2 \|\mathbf{y}\|^2)}{t^2 \|\mathbf{y}\|^2} + \frac{1}{2 \|\mathbf{y}\|^2} \langle \mathbf{y}, \nabla^2 f(\mathbf{x}^*) \mathbf{y} \rangle$$
  - The inequality follows from the local optimality of  $\mathbf{x}^*$
- Taking the limit of  $t \rightarrow 0$  of both sides gives
$$\frac{1}{2 \|\mathbf{y}\|^2} \langle \mathbf{y}, \nabla^2 f(\mathbf{x}^*) \mathbf{y} \rangle \geq 0$$
  - This statement proves that the Hessian is PSD

# Group Exercise



Notation:

$$\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$$

Consider the function

$$f(x, y) = x^2 + 4xy - y^2 - 8x - 6y + 10$$

1. Determine  $\nabla f(x, y)$
2. Determine  $\nabla^2 f(x, y)$
3. Determine a critical point  $\mathbf{x}^*$
4. Is  $\mathbf{x}^*$  a local min, a local max, or neither?

# Convexity



- We say that  $f$  is *convex* if

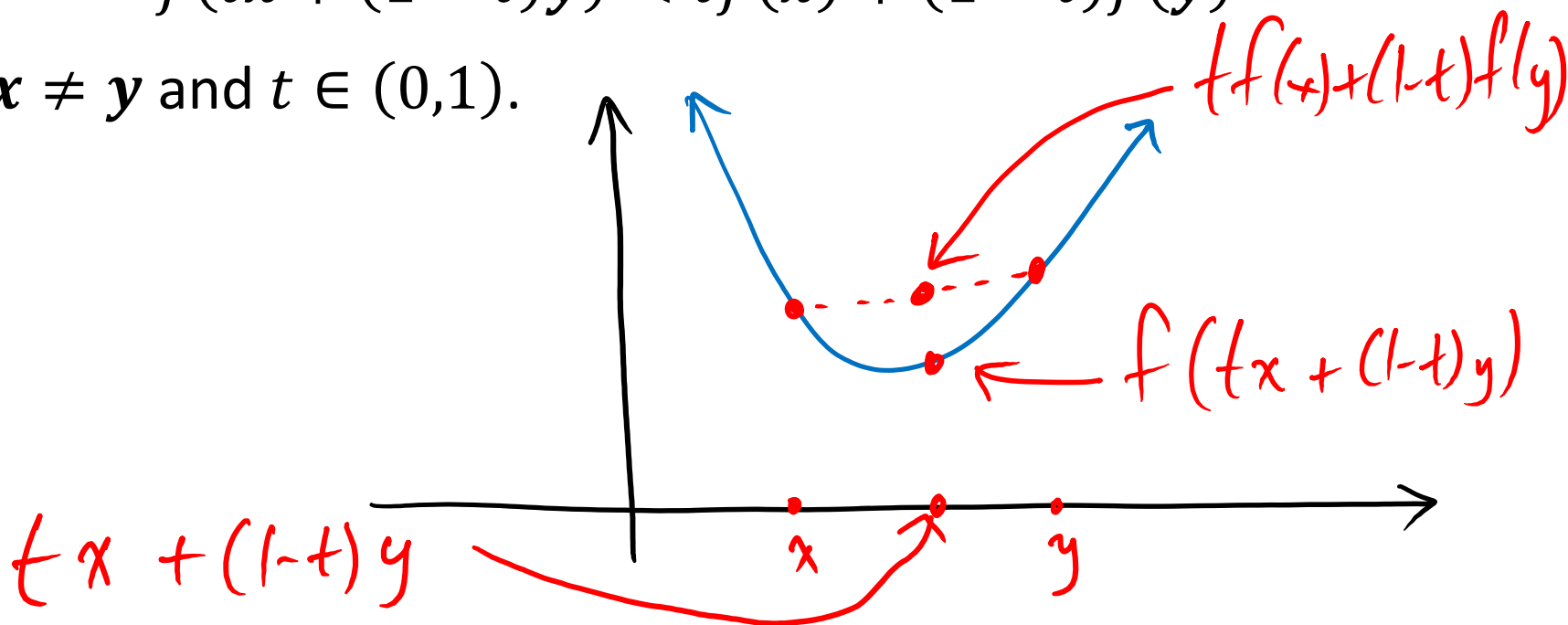
$$f(t\mathbf{x} + (1 - t)\mathbf{y}) \leq tf(\mathbf{x}) + (1 - t)f(\mathbf{y})$$

for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$  and  $t \in [0, 1]$ .

- We say  $f$  is *strictly convex* if

$$f(t\mathbf{x} + (1 - t)\mathbf{y}) < tf(\mathbf{x}) + (1 - t)f(\mathbf{y})$$

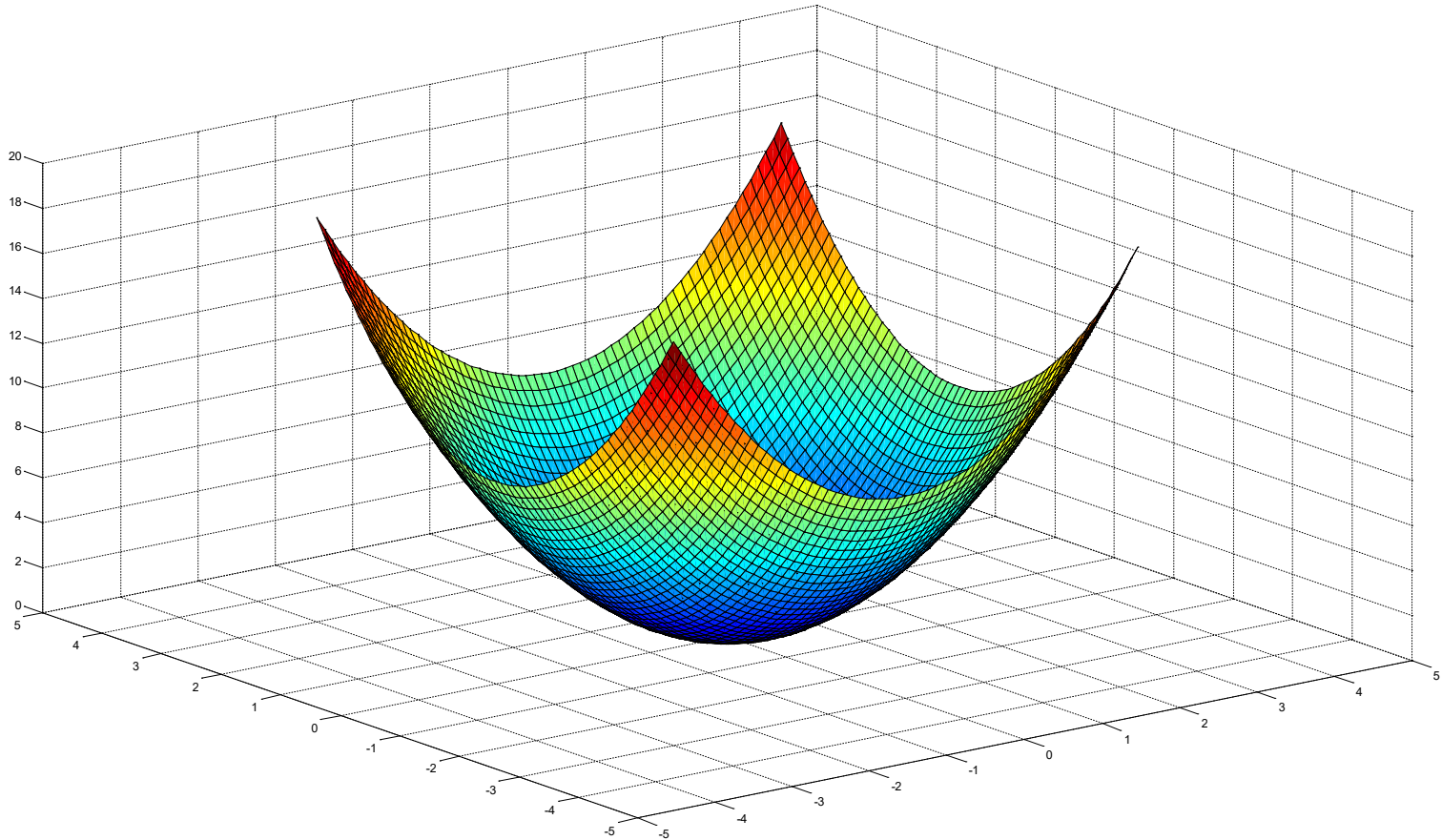
$\forall \mathbf{x} \neq \mathbf{y}$  and  $t \in (0, 1)$ .



# Convexity



$$f(x, y) = x^2 + y^2$$

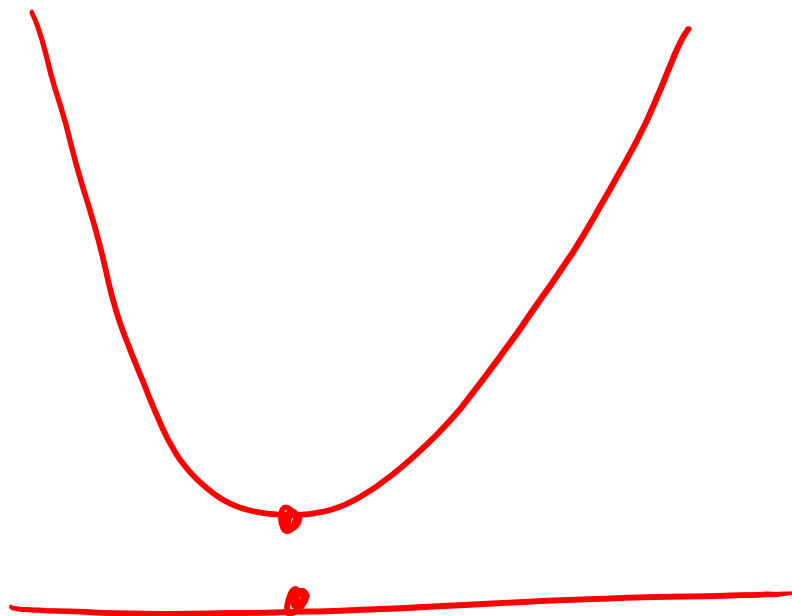


# Convex functions are nice



## Properties of convex functions

1. If  $f$  is convex, then every local min is a global min
2. If  $f$  is strictly convex, then  $f$  has at most one global min



# Convex functions are nice



*Proof of 1:* Suppose  $\mathbf{x}^*$  is a local min but not a global min. Then  $\exists \mathbf{y}^* \in \mathbb{R}$  such that  $f(\mathbf{y}^*) < f(\mathbf{x}^*)$ . By convexity,  $\forall t \in [0, 1)$  we have

$$\begin{aligned} f(t\mathbf{x}^* + (1-t)\mathbf{y}^*) &\leq tf(\mathbf{x}^*) + (1-t)f(\mathbf{y}^*) \\ &< tf(\mathbf{x}^*) + (1-t)f(\mathbf{x}^*) \\ &= f(\mathbf{x}^*) \end{aligned}$$

Taking  $t \nearrow 1$ , the above strict inequality contradicts local minimality of  $\mathbf{x}^*$ . Thus  $\mathbf{x}^*$  is a global min.



# Group Exercise



1. Give an example of a function  $f$  that is
  - convex but not strictly convex
  - convex and has more than one global minimum
  - strictly convex, but has no global minimum
2. Is the sum of convex functions necessarily convex?  
Prove or provide a counter example
3. Is the product of convex functions necessarily convex?  
Prove or provide a counter example

# 1<sup>st</sup> order characterization of convexity

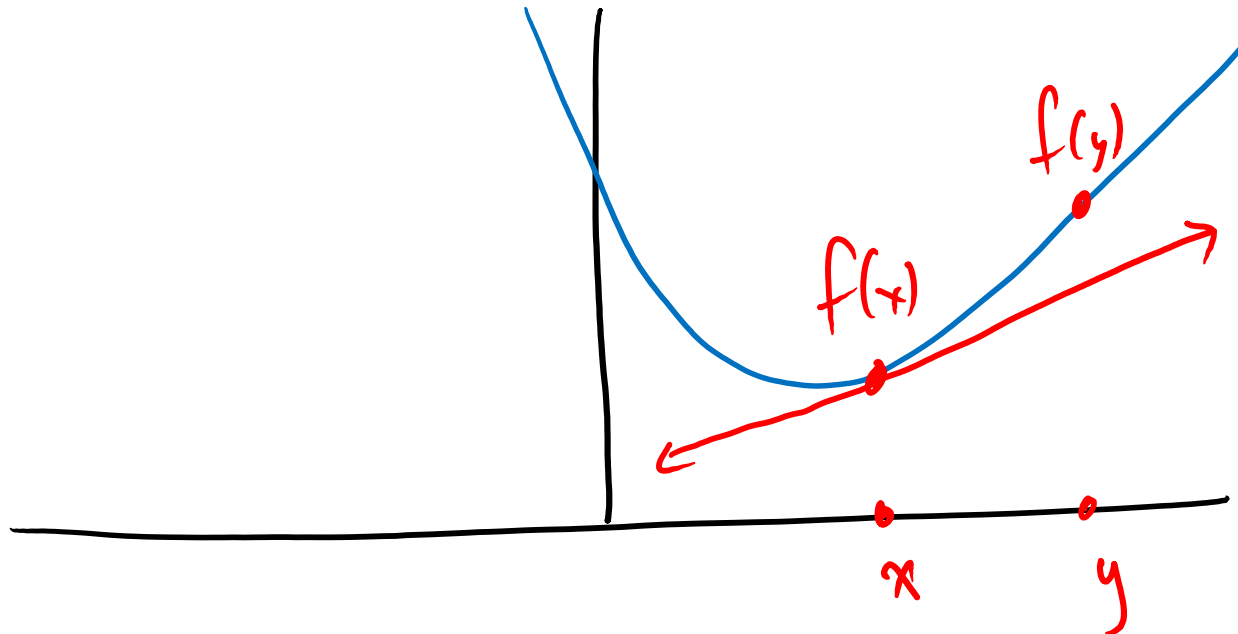


1. Suppose  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is differentiable. Then  $f$  is convex iff  $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ ,

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

2. Similarly,  $f$  is strictly convex iff  $\forall \mathbf{x} \neq \mathbf{y}$ ,

$$f(\mathbf{y}) > f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle.$$



# 1<sup>st</sup> order characterization of convexity



*Proof of 1 ( $\implies$ ):* First, assume  $\mathbf{x}$  is convex. For any  $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d, t \in [0, 1]$ ,

$$\begin{aligned} f(t\mathbf{y} + (1 - t)\mathbf{x}) &\leq tf(\mathbf{y}) + (1 - t)f(\mathbf{x}) \\ &= f(\mathbf{x}) + t(f(\mathbf{y}) - f(\mathbf{x})). \end{aligned} \tag{1}$$

Rearranging,

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

The limit of the LHS is a directional derivative and equal to  $\langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$  by the chain rule. Therefore  $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$ .

# 1<sup>st</sup> order characterization of convexity



*Proof of 1 ( $\Leftarrow$ ):* Now suppose conversely that  $\forall \mathbf{x}, \mathbf{y}$

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

Let  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$  and  $t \in [0, 1]$ . Denote  $\mathbf{z} = t\mathbf{x} + (1 - t)\mathbf{y}$ . Applying  $*$  twice we have

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

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

Now multiply  $*a$  by  $t$ ,  $*b$  by  $(1 - t)$  and add:

$$\begin{aligned} tf(\mathbf{x}) + (1 - t)f(\mathbf{y}) &\geq f(\mathbf{z}) + \langle \nabla f(\mathbf{z}), t\mathbf{x} + (1 - t)\mathbf{y} - \mathbf{z} \rangle \\ &= f(t\mathbf{x} + (1 - t)\mathbf{y}) \end{aligned}$$

as  $t\mathbf{x} + (1 - t)\mathbf{y} - \mathbf{z} = \mathbf{0}$ . This establishes convexity.

# 1<sup>st</sup> order condition for local min, revisited



- Let  $f$  be convex and continuously differentiable. Then  $\mathbf{x}^*$  is a global min iff  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ .

- Proof:

( $\Rightarrow$ ) Already discussed

$$(\Leftarrow) \quad \begin{aligned} \forall \mathbf{y}, f(\mathbf{y}) &\geq f(\mathbf{x}^*) + \langle \nabla f(\mathbf{x}^*), \mathbf{y} - \mathbf{x}^* \rangle \\ &= f(\mathbf{x}^*) \end{aligned}$$

- Thus for convex functions,  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  is both necessary and sufficient for  $\mathbf{x}^*$  to be a global min.

# Second order conditions of convexity



More important properties:

- $f$  is convex  $\Leftrightarrow \nabla^2 f(\mathbf{x})$  is positive semidefinite  $\forall \mathbf{x} \in \mathbb{R}^d$
- $f$  is strictly convex  $\Leftarrow \nabla^2 f(\mathbf{x})$  is positive definite  $\forall \mathbf{x} \in \mathbb{R}^d$
- Proofs follow from multidimensional Taylor series expansions and the property that if  $f$  is convex, then  $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$

# Group Exercises



1. Give an example of a function  $f$  and a point  $\mathbf{x}$  such that  $\nabla^2 f(\mathbf{x})$  is PSD but  $\mathbf{x}$  is not a local minimizer
2. Give an example of a function  $f$  that is strictly convex, but such that there exists  $\mathbf{x}$  for which  $\nabla^2 f(\mathbf{x})$  is not PD.
3. Numerically determine a critical point of

$$f(x, y) = x^2 + 2xy + 3y^2 + 4x + 5y + 6$$

and also determine if it is a local/global min or max. *Note:* If you don't have access to Matlab/Python/etc. in class, you can also use Wolfram Alpha for many calculations like eigenvalue decompositions

# Regularized Logistic Regression



- Unless  $n \gg d$ , it is preferable to minimize the modified objective function in logistic regression:

$$J(\boldsymbol{\theta}) = -\ell(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|^2$$

- $\lambda > 0$  is a fixed, user-specified constant called a *regularization parameter*
- Why introduce the regularization term?
  - So the Hessian is PD (see a future HW)



# Methods for solving optimization problems

1. Gradient Descent
2. Stochastic Gradient Descent
3. Newton's method
4. Subgradient methods

# Iterative methods



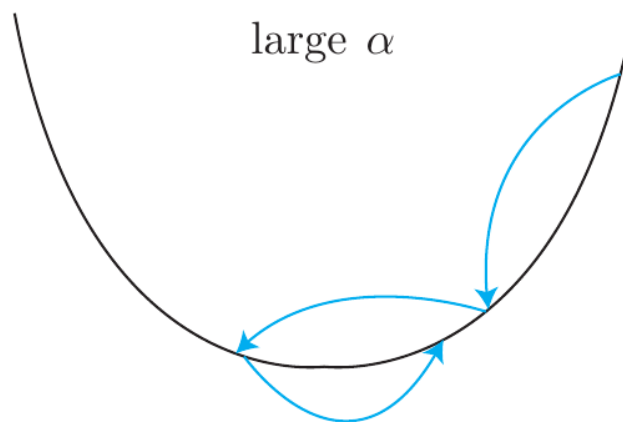
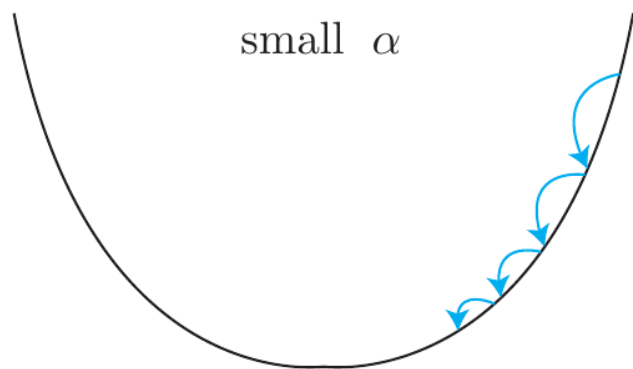
- In many problems, an analytical solution to the minimization problem does not exist
  - Including (regularized) logistic regression
- However, we can use an iterative method to solve it
  - Convexity guarantees that the iterative solution is the global minimum

# Gradient Descent (GD)



- Consider minimizing the generic objective function  $J(\boldsymbol{\theta})$
- What is a geometric interpretation of  $\nabla J(\boldsymbol{\theta})$ ?
  - Direction of steepest ascent
  - Thus iterative approaches take steps in opposite direction, i.e., the direction of steepest descent
- Initial guess  $\boldsymbol{\theta}_0$
- For  $t = 1, \dots, \text{max\_iter}$ 
  - $\boldsymbol{\theta}_t \leftarrow \boldsymbol{\theta}_{t-1} - \alpha \nabla J(\boldsymbol{\theta}_{t-1})$
  - If convergence condition satisfied, exit

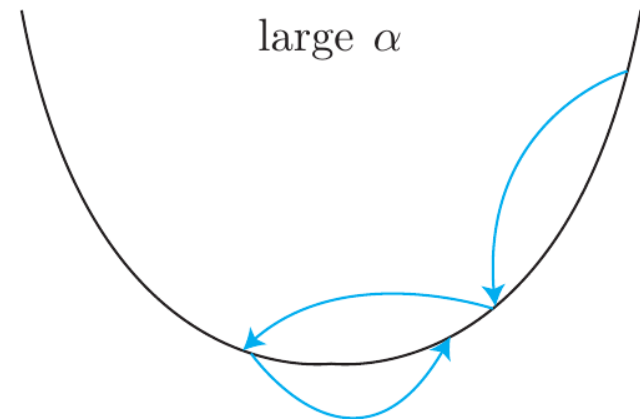
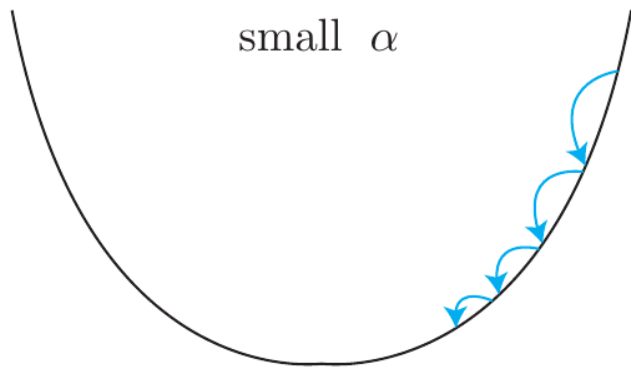
End



# Step Size



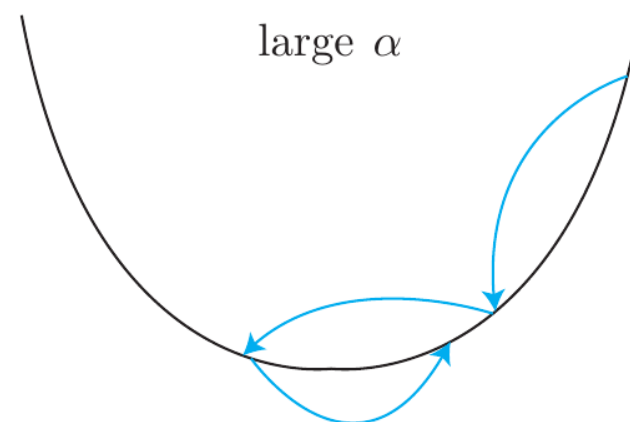
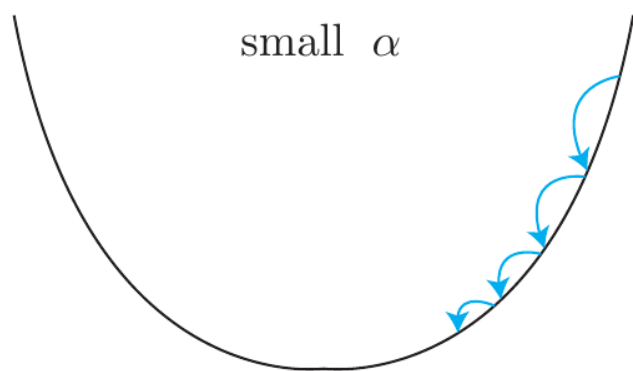
- In practice, typically need  $\alpha$  to decrease as a function of  $t$ 
  - Example:  $\alpha_t = \frac{1}{t}$



# Group Exercise



1. What would be a good convergence condition (also known as a stopping criterion) for terminating GD?
2. What are some potential problems with setting  $\alpha_t = \frac{1}{t}$ ?
3. How could you choose  $\alpha_t$  to counter these problems?
4. Does the initialization  $\theta_0$  affect the result?



# Stochastic Gradient Descent (SGD)



- In many ML problems, we can write

$$J(\boldsymbol{\theta}) = \sum_{i=1}^n J_i(\boldsymbol{\theta})$$

- For example, in regularized ERM, we can write

$$J_i(\boldsymbol{\theta}) = \frac{1}{n} \left( L(y_i, f_{\boldsymbol{\theta}}(\mathbf{x}_i)) + \lambda \Omega(f_{\boldsymbol{\theta}}) \right)$$

- Thus

$$\nabla J(\boldsymbol{\theta}) = \sum_{i=1}^n \nabla J_i(\boldsymbol{\theta})$$

# Stochastic Gradient Descent (SGD)



- In GD, we calculate the gradient using all of the data
  - For large  $n$ , each step in the parameter space  $\theta$  takes a lot of computations
  - Thus learning can occur slowly
- For SGD, we estimate the gradient using a small random sample of training inputs
  - Speeds up gradient descent and learning

# Stochastic Gradient Descent



- Pick a random set of  $m < n$  training inputs
  - Referred to as a *mini-batch*
- Estimate the gradient using the samples from the mini-batch and take a step in the direction of the negative gradient
- If  $m$  is large enough, then the estimated gradient will be roughly equal to the gradient using the full data



# Stochastic Gradient Descent Summary



A training ***epoch***:

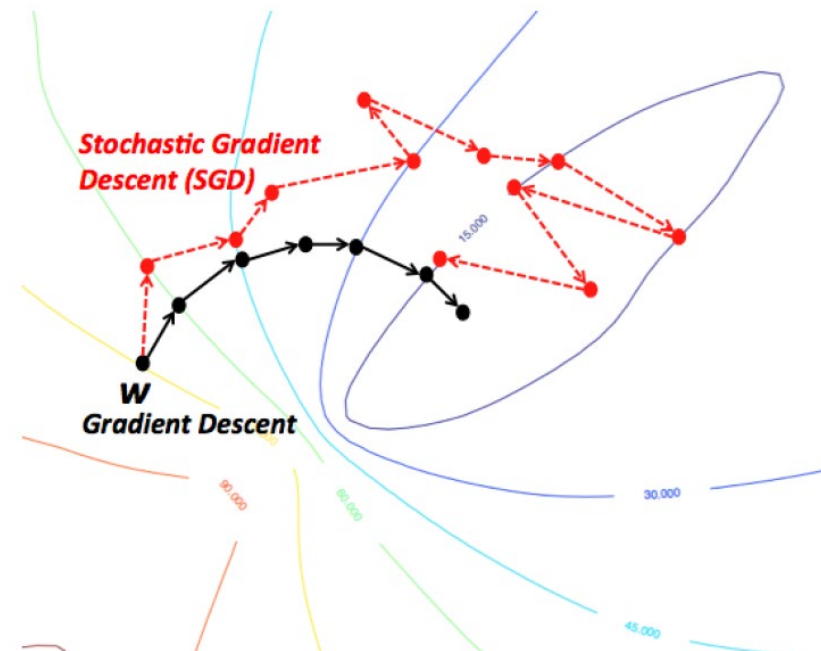
1. Pick a random subset of the training data
  - Referred to as a mini-batch
2. Update the parameters using the gradient estimates from the mini-batch
3. Pick another random mini-batch from the remaining training points and repeat step 2
  - Repeat until all training inputs have been used

Repeat multiple epochs until stopping conditions are met

# Analogy to political polling



- Much easier to carry out a poll than run a full election
- Similarly, it's much easier to estimate gradients from mini-batches than the entire training set
- Downside: gradient estimates will be noisier in SGD
- Generally ok as we only need to move in a general direction that decreases  $J$ 
  - This can be an advantage when the problem is nonconvex
- In practice, SGD (or its variants) are used extensively in learning neural networks



# Newton's method

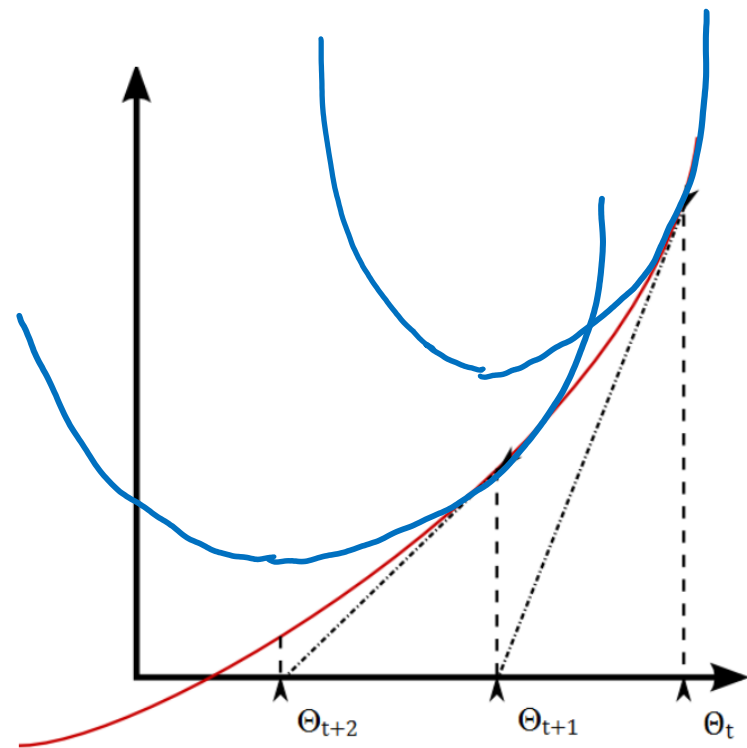


- Goal is to solve  $\nabla J(\boldsymbol{\theta}) = \mathbf{0}$  for an objective function  $J$
- Newton's method AKA Newton-Raphson:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - (\nabla^2 J(\boldsymbol{\theta}_t))^{-1} \nabla J(\boldsymbol{\theta}_t)$$

- Newton's method can be viewed as minimizing the second order approximation:

$$J(\boldsymbol{\theta}) \approx J(\boldsymbol{\theta}_t) + \nabla J(\boldsymbol{\theta}_t)^T (\boldsymbol{\theta} - \boldsymbol{\theta}_t) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_t)^T \nabla^2 J(\boldsymbol{\theta}_t) (\boldsymbol{\theta} - \boldsymbol{\theta}_t)$$

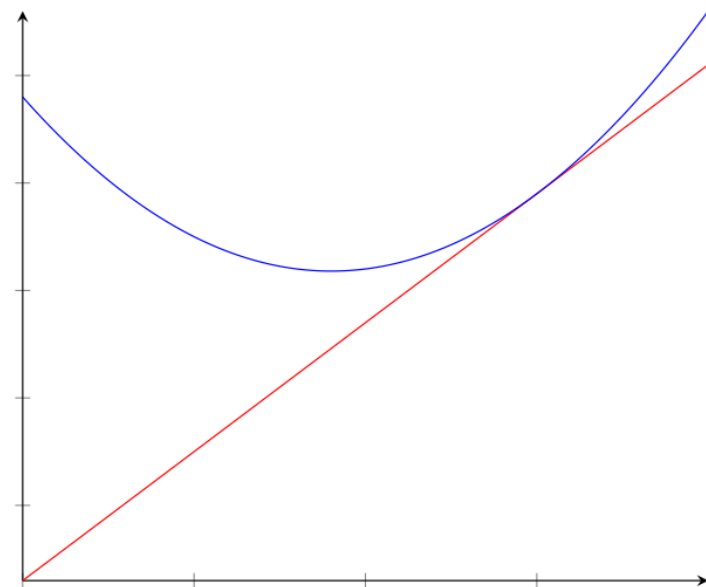


# Subgradient methods



- Subgradient methods are generalizations of gradient descent that can be applied to *nondifferentiable, convex* objective functions, like ERM with hinge loss
  - Hinge loss:  $L(y, t) = \max(0, 1 - yt)$
- Let  $g: \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and let  $\theta \in \mathbb{R}^d$ . If  $g$  is differentiable, then  $u = \nabla g(\theta)$  is the only vector such that

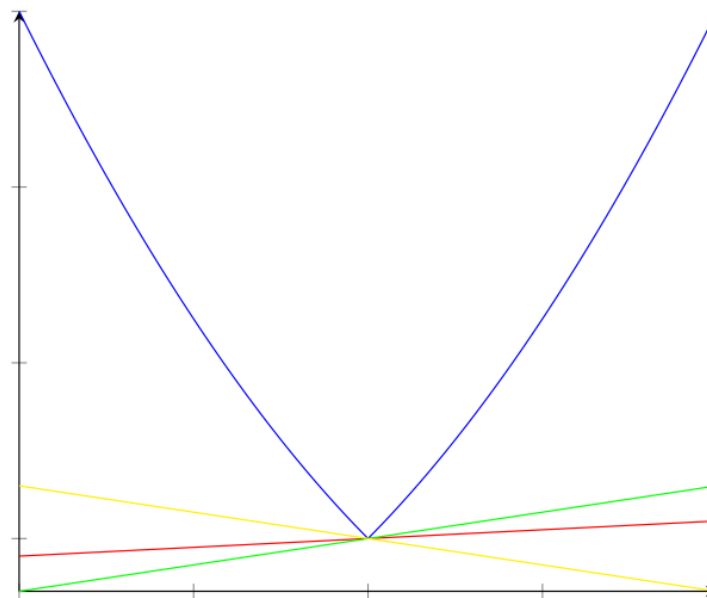
$$g(\theta') \leq g(\theta) + u^T(\theta' - \theta) \quad \forall \theta'$$



# Subgradients



- If  $g$  is convex but *not* differentiable, then for some  $\theta$ , there may be many  $u$  satisfying the previous inequality.
- We define the *subdifferential* of  $g$  at  $\theta$ , denoted  $\partial g(\theta)$ , to be the set of all  $u$  satisfying the inequality.
- A *subgradient* is any element of the subdifferential.
- In the figure, the subdifferential is the interval  $[g'_-(\theta), g'_+(\theta)]$  where  $g'_-(\theta), g'_+(\theta)$  denote the left and right derivatives.



# Subgradient Methods



- In a *subgradient method*, we update the parameter just as in gradient descent, but where the gradient is replaced by *any* subgradient
- Pseudo-code for minimizing  $g(\boldsymbol{\theta})$ 
  - Initialize  $\boldsymbol{\theta}_0$
  - $t \leftarrow 0$
  - Repeat
    - Select  $\mathbf{u}_t \in \partial g(\boldsymbol{\theta}_t)$
    - $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t - \alpha_t \mathbf{u}_t$
    - $t \leftarrow t + 1$
  - Until stopping criterion satisfied
- If we can write  $g(\boldsymbol{\theta}) = \sum_{i=1}^n g_i(\boldsymbol{\theta})$ , then we can also have a *stochastic subgradient method*, analogous to SGD.

# Big Takeaways



- The combination of convexity and differentiability enable us to determine if a global minimum exists
- Differentiability also provides us with methods for finding a minimum when a closed form solution may not exist
  - If convex but not differentiable, then subgradient methods may be useful
  - Other methods also exist such as the Majorize-Minimize Algorithm
- Thus when considering loss functions for an ML problem, convexity and differentiability are important properties to consider

# Further reading



- ESL Section 10.10 and 11.4