# CPSC 340:
# Machine Learning and Data Mining

Gradient Descent

# Last Time: Change of Basis
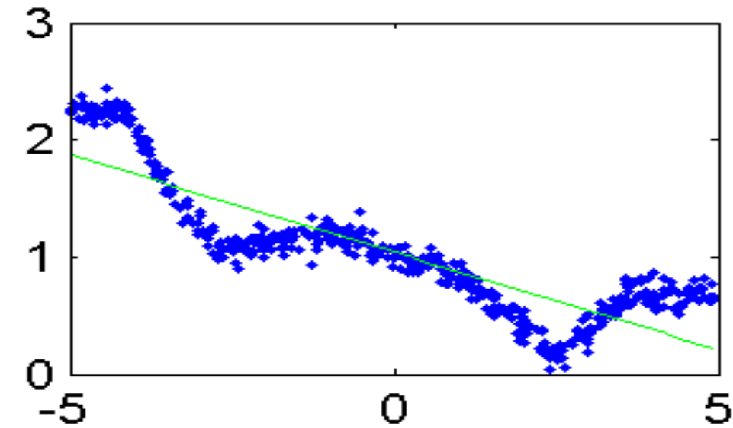
- Last time we discussed change of basis:
  - E.g., polynomial basis:

$$\text{Replace } X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \text{ with } Z = \begin{bmatrix} 1 & x_1 & (x_1)^2 & ----- & (x_1)^p \\ 1 & x_2 & (x_2)^2 & ----- & (x_2)^p \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & (x_n)^2 & ----- & (x_n)^p \end{bmatrix}$$

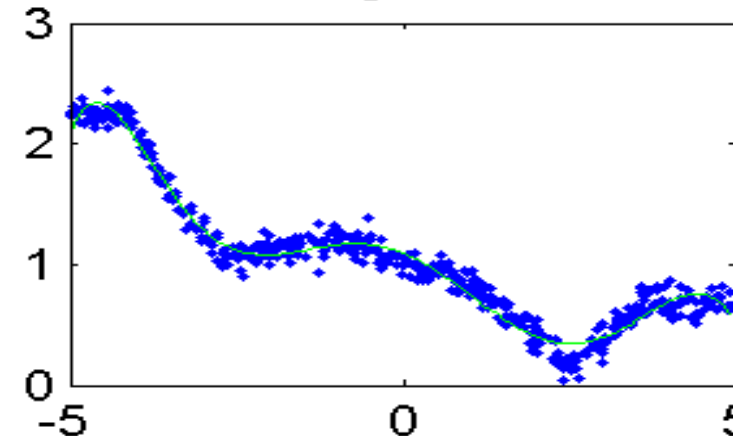  - You can fit non-linear models with linear regression.

$$\hat{y}_i = v^T z_i = w_0 + w_1 x_i + w_2 x_i^2 + w_3 x_i^3 + \cdots + w_p x_i^p$$

  - Just treat 'Z' as your data, then fit linear model.



Degree 7

# Optimization Terminology

- When we minimize or maximize a function we call it "optimization".
  - In least squares, we want to solve the "optimization problem":

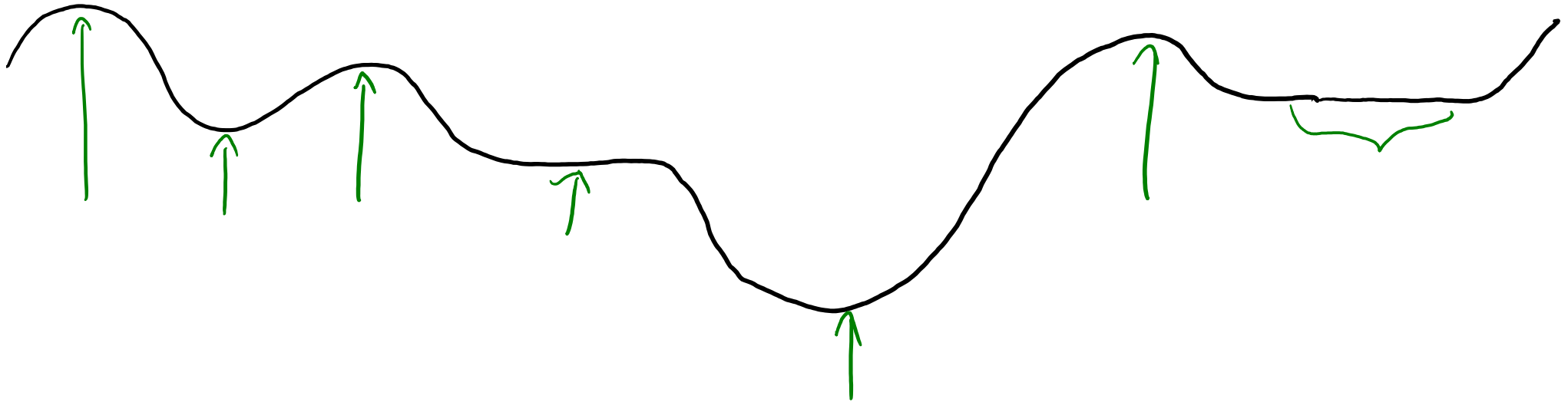$$\min_{w \in \mathbb{R}^d} \frac{1}{2} \sum_{j=1}^{n} \left( w^\top x_i - y_i \right)^2$$

  - The function being optimized is called the "objective".
    - Also sometimes called "loss" or "cost", but these have different meanings in ML.

  - The set over which we search for an optimum is called the domain.

  - Often, instead of the minimum objective value, you want a minimizer.
    - A set of parameters 'w' that achieves the minimum value.

# Discrete vs. Continuous Optimization

- We have seen examples of continuous optimization:
  - Least squares:
    - Domain is the real-valued set of parameters 'w'.
    - Objective is the sum of the squared training errors.

- We have seen examples of discrete optimization:
  - Fitting decision stumps:
    - Domain is the finite set of unique rules.
    - Objective is the number of classification errors (or infogain).

- We have also seen a mixture of discrete and continuous:
  - K-means: clusters are discrete and means are continuous.
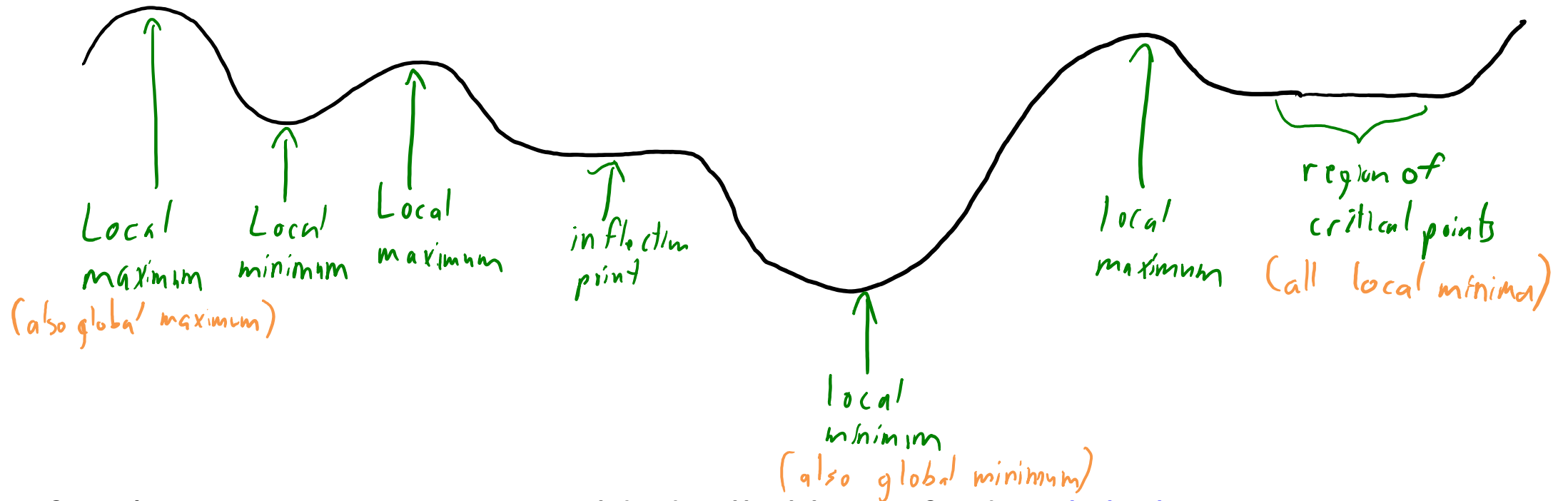
# Stationary/Critical Points

- A 'w' with $\nabla f(w) = 0$ is called a stationary point or critical point.
  - The slope is zero so the tangent plane is "flat".



Critical points

# Stationary/Critical Points

- A 'w' with $\nabla f(w) = 0$ is called a stationary point or critical point.
  - The slope is zero so the tangent plane is "flat".



Local maximum (also global maximum)    Local minimum    Local maximum    inflection point    local minimum (also global minimum)    local maximum    region of critical points (all local minima)

  - If we're minimizing, we would ideally like to find a global minimum.
    - But for some problems the best we can do is find a stationary point where $\nabla f(w)=0$.

# Motivation: Large-Scale Least Squares

- Normal equations find 'w' with $\nabla f(w) = 0$ in $O(nd^2 + d^3)$ time.

$$(X^T X)_w = X^T y$$

$\underbrace{\phantom{(X^T X)}}$    $\underbrace{\phantom{X^T y}}$

$O(nd^2)$      $O(nd)$

(matrix multiply)   (matrix-vector)

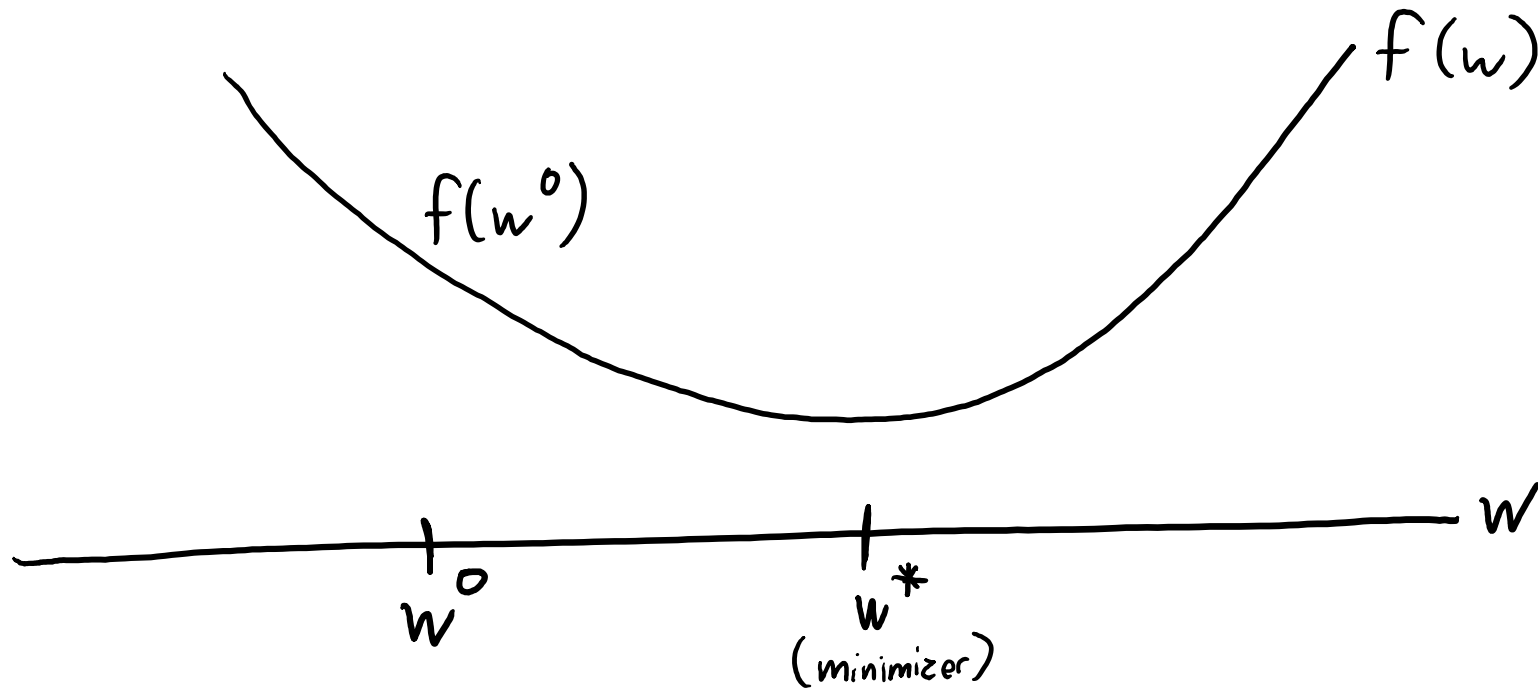$\rightarrow$ Solving a $d \times d$ system is $O(d^3)$

  – Very slow if 'd' is large.

- Alternative when 'd' is large is gradient descent methods.
  – Probably the most important class of algorithms in machine learning.

# Gradient Descent for Finding a Local Minimum

- Gradient descent is an iterative optimization algorithm:
  - It starts with a "guess" $w^0$.
  - It uses the gradient $\nabla f(w^0)$ to generate a better guess $w^1$.
  - It uses the gradient $\nabla f(w^1)$ to generate a better guess $w^2$.
  - It uses the gradient $\nabla f(w^2)$ to generate a better guess $w^3$.

    …
  - The limit of $w^t$ as 't' goes to $\infty$ has $\nabla f(w^t) = 0$.
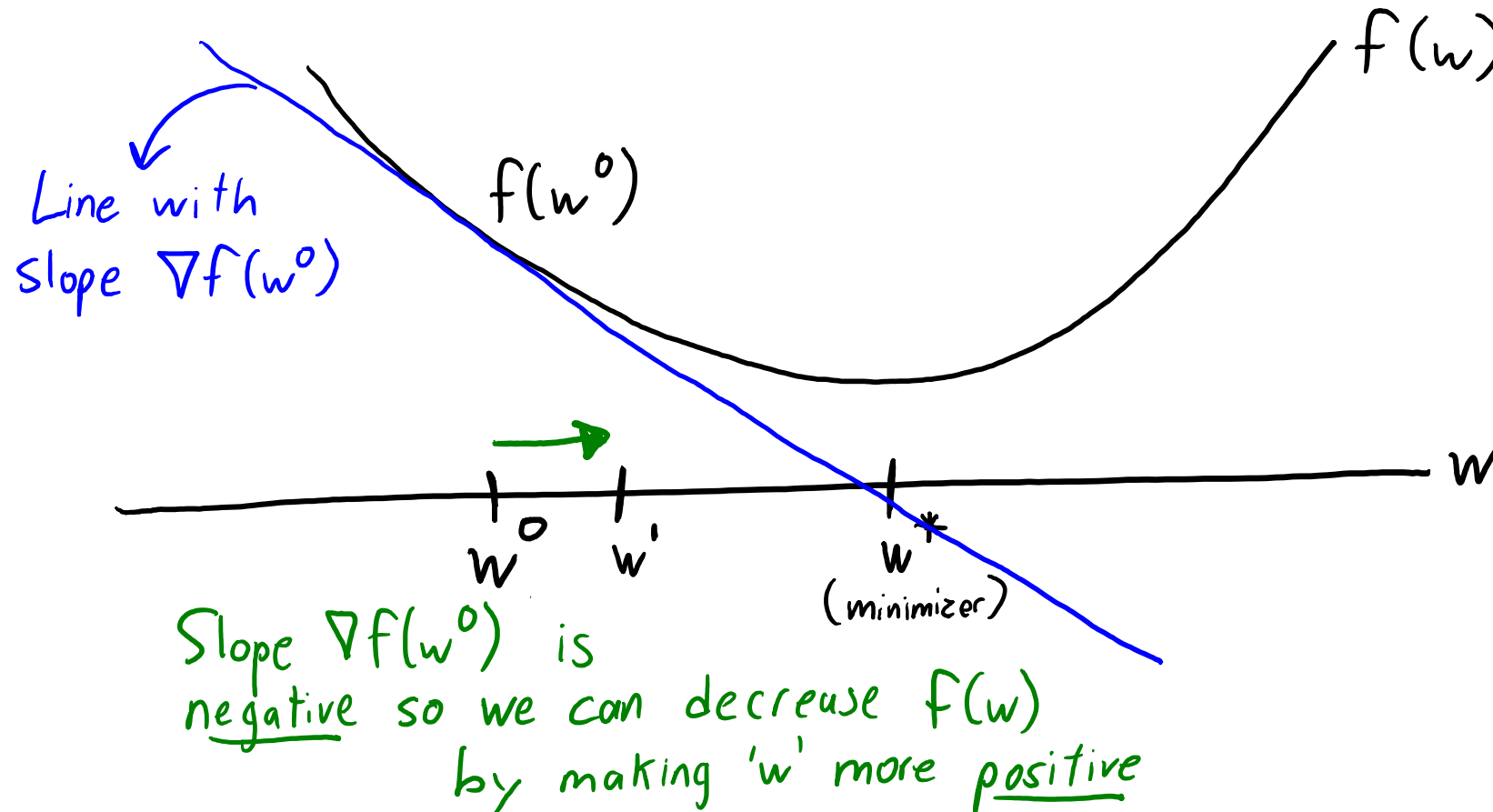

- It converges to a global optimum if 'f' is "convex".

# Gradient Descent for Finding a Local Minimum

- Gradient descent is based on a simple observation:
  - Given parameters 'w', the direction of largest decrease is $-\nabla f(w)$.

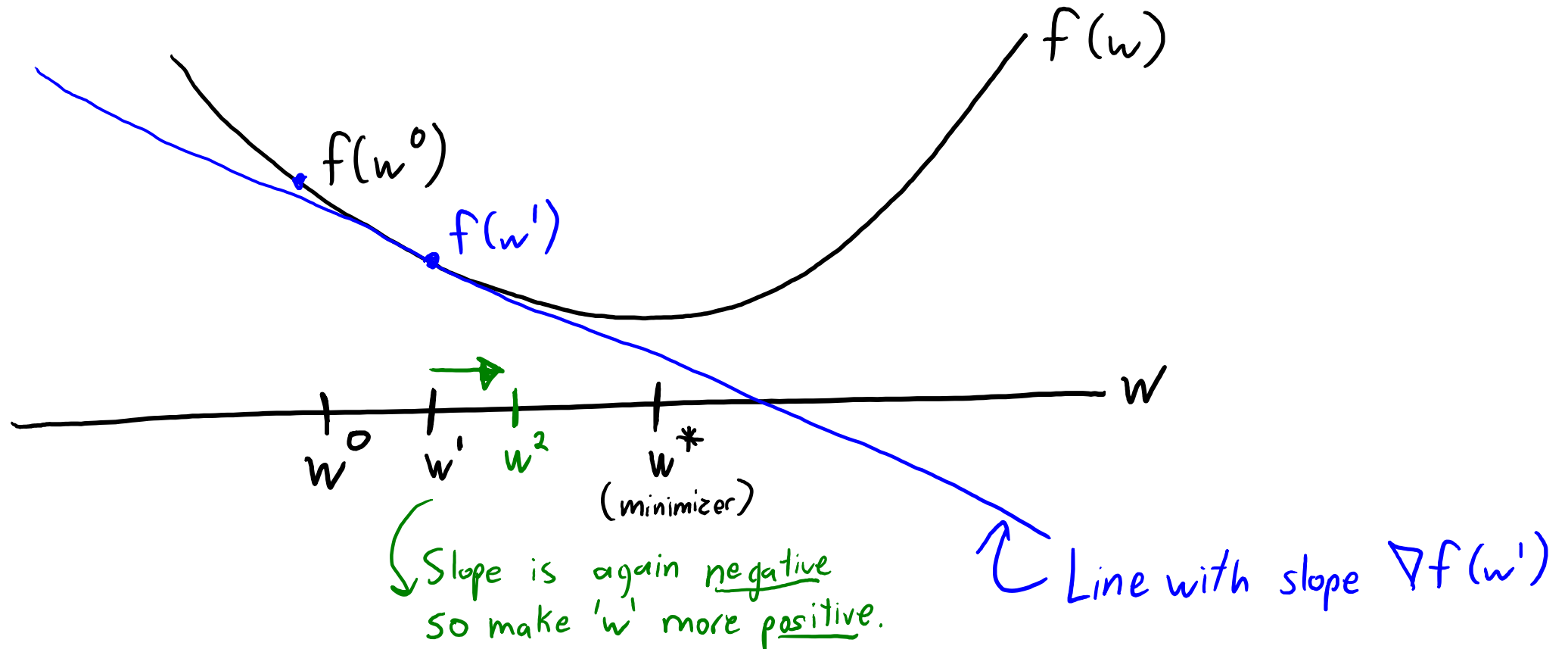# Gradient Descent for Finding a Local Minimum

- Gradient descent is based on a simple observation:
  - Given parameters 'w', the direction of largest decrease is $-\nabla f(w)$.

# Gradient Descent for Finding a Local Minimum

- Gradient descent is based on a simple observation:
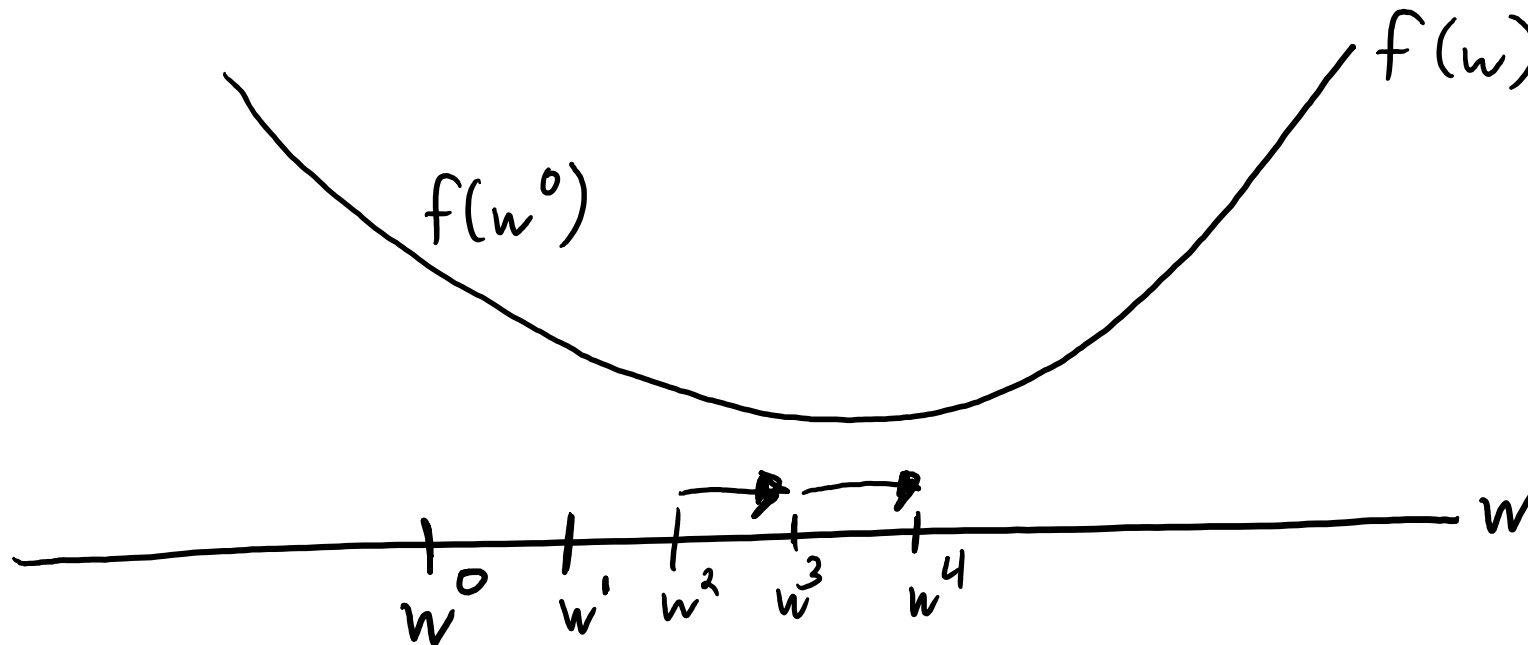  - Given parameters 'w', the direction of largest decrease is $-\nabla f(w)$.

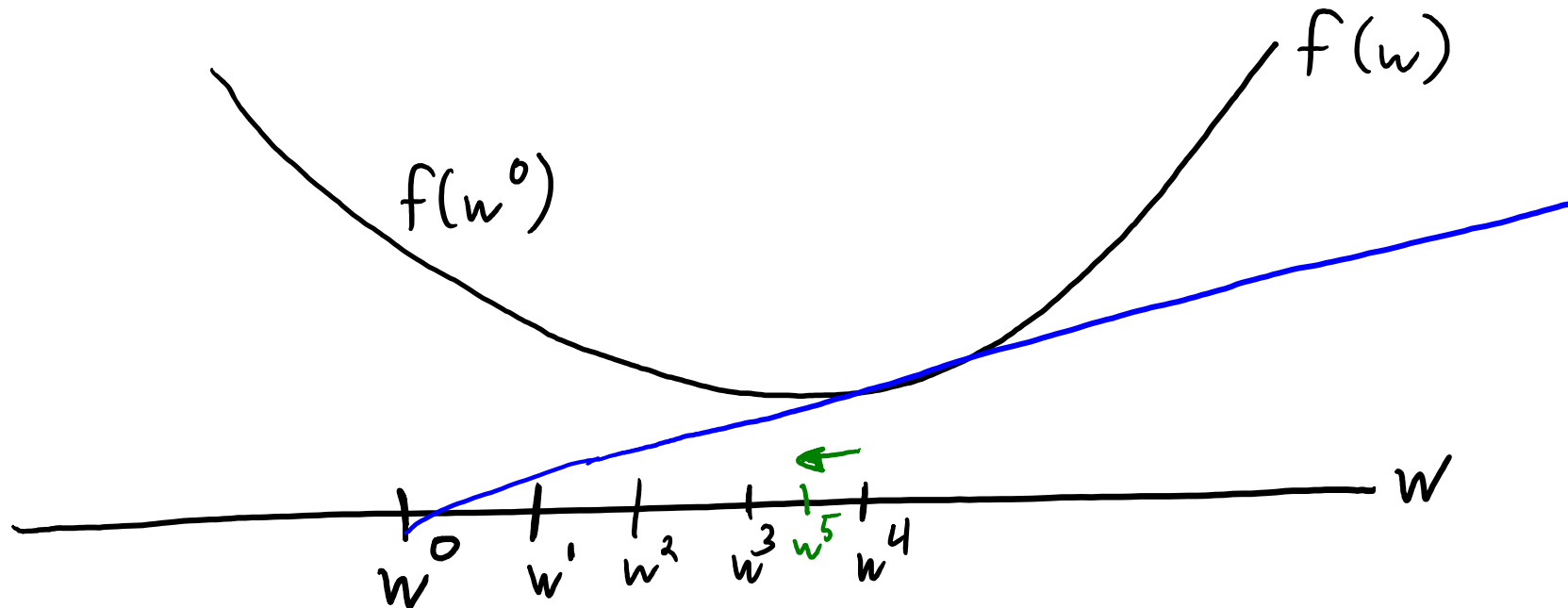# Gradient Descent for Finding a Local Minimum

- Gradient descent is based on a simple observation:
  - Given parameters 'w', the direction of largest decrease is $-\nabla f(w)$.

# Gradient Descent for Finding a Local Minimum

- Gradient descent is based on a simple observation:
  - Given parameters 'w', the direction of largest decrease is $-\nabla f(w)$.



Now the slope $\nabla f(w^4)$ is positive so we move in the negative direction.

# Gradient Descent for Finding a Local Minimum

- We start with some initial guess, $w^0$.

- Generate new guess by moving in the negative gradient direction:

$$w^1 = w^0 - \alpha^0 \nabla f(w^0)$$

  - This decreases 'f' if the "step size" (aka "learning rate") $\alpha^0$ is small enough.
  - Usually, we decrease $\alpha^0$ if it increases 'f'

"step size" is not a great name since it is not sign(grad)*alpha

- Repeat to successively refine the guess:

$$w^{t+1} = w^t - \alpha^t \nabla f(w^t) \quad \text{for } t = 1, 2, 3, \ldots$$

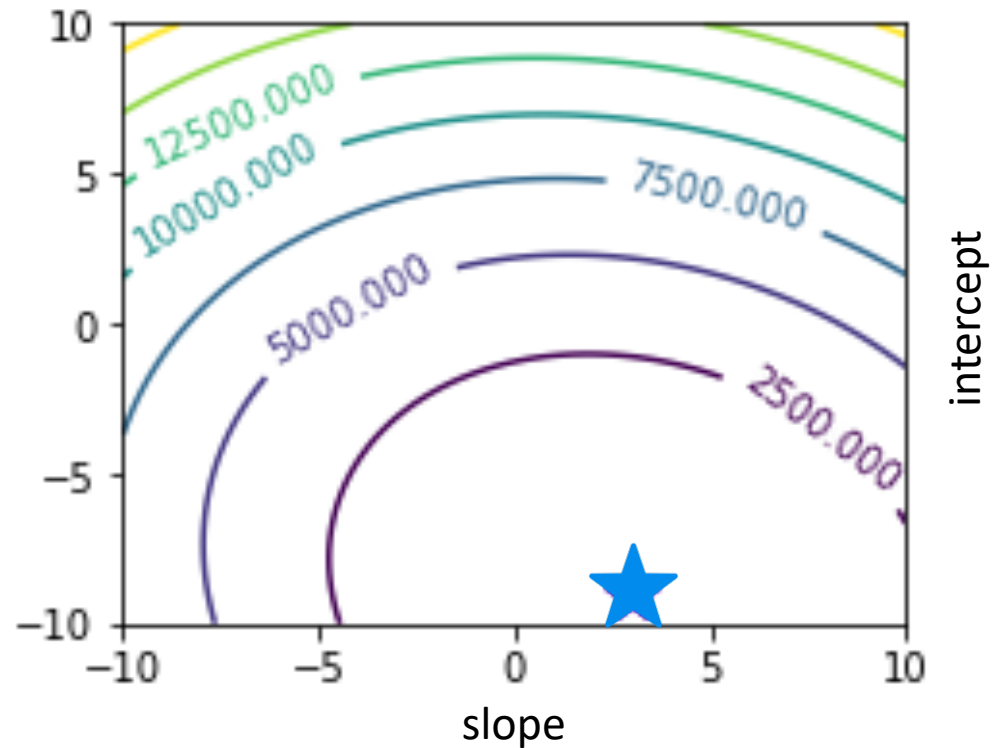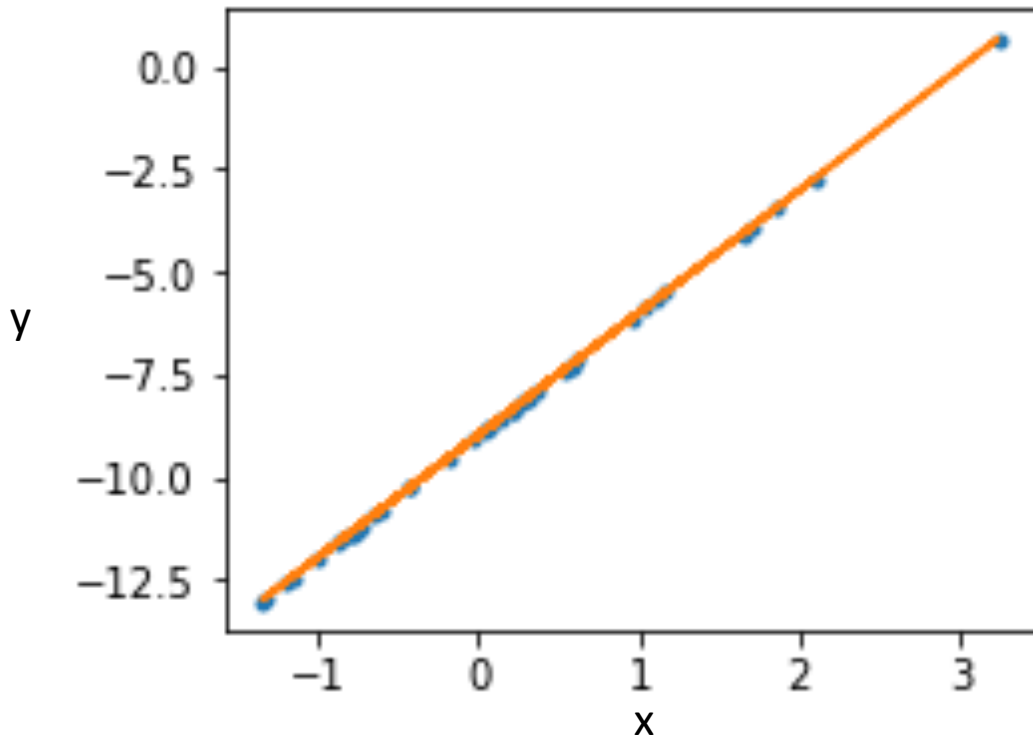- Stop if not making progress or $\|\nabla f(w^t)\| \leq \varepsilon$

Approximate local minimum

$\longrightarrow$ Some small scalar.

# Data Space vs. Parameter Space

- Left: Usual regression plot in the "x vs. y" data space



- Right: plot of "slope vs. intercept" in parameter space
  - Points in parameter space correspond to models (★ is the least squares parameters, which is the line shown in the left plot)

# Gradient Descent in Data Space vs. Parameter Space

- Gradient descent starts with an initial guess in parameter space:



– And each iteration tries to move guess closer to solution.

# Gradient Descent in Data Space vs. Parameter Space

- Gradient descent starts with an initial guess in parameter space:



– And each iteration tries to move guess closer to solution.

# Gradient Descent in Data Space vs. Parameter Space

- Gradient descent starts with an initial guess in parameter space:



– And each iteration tries to move guess closer to solution.

# Gradient Descent in Data Space vs. Parameter Space
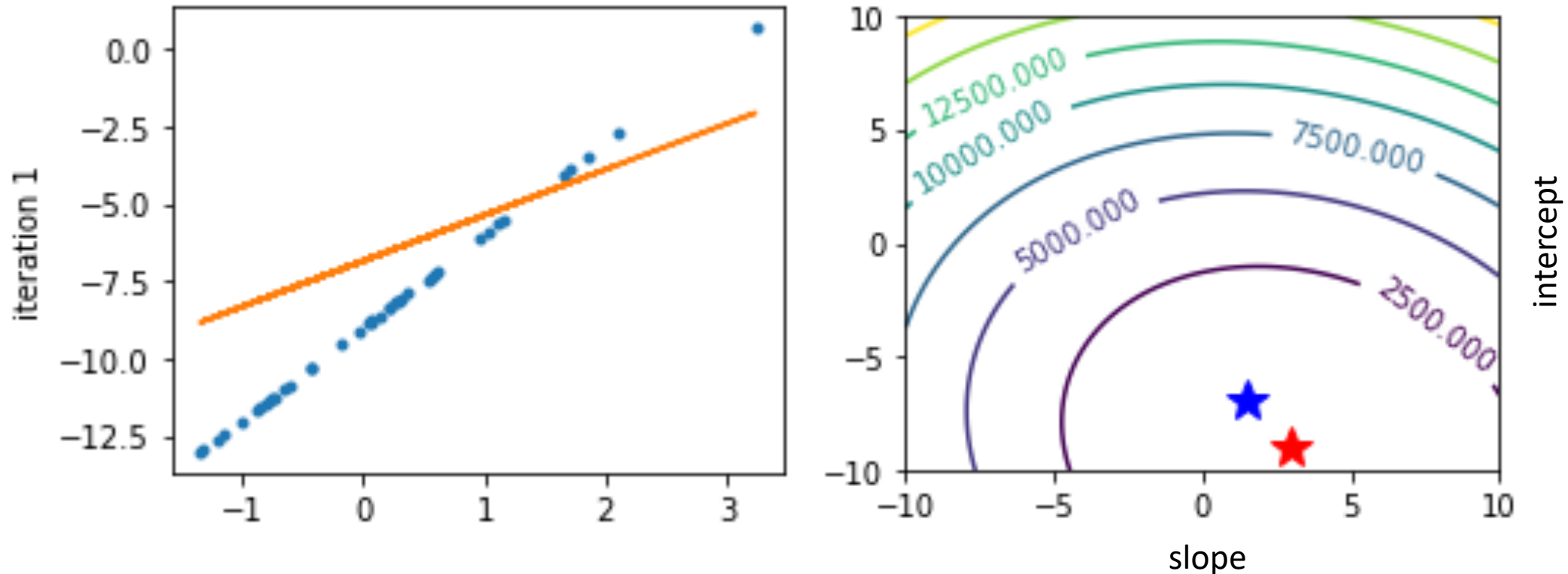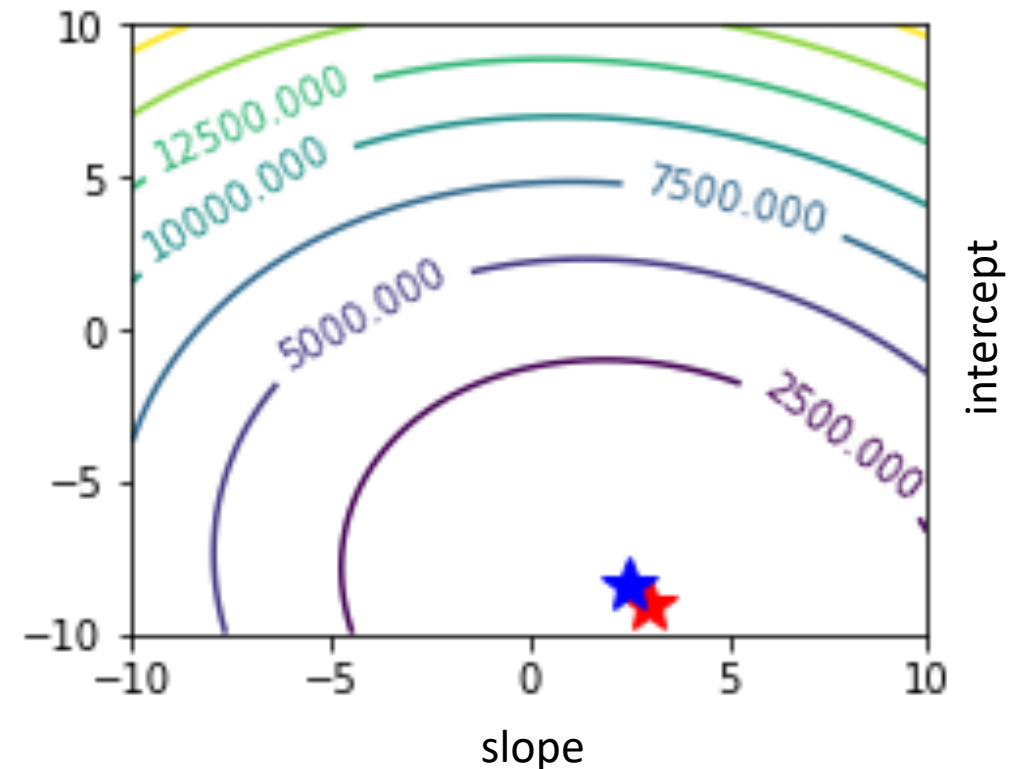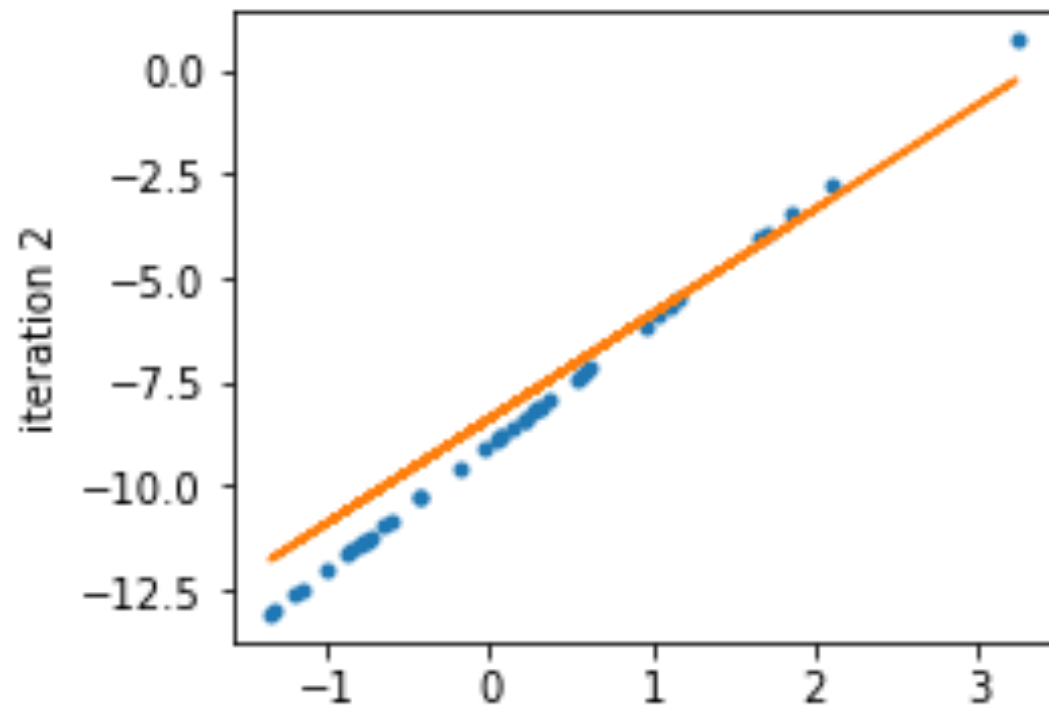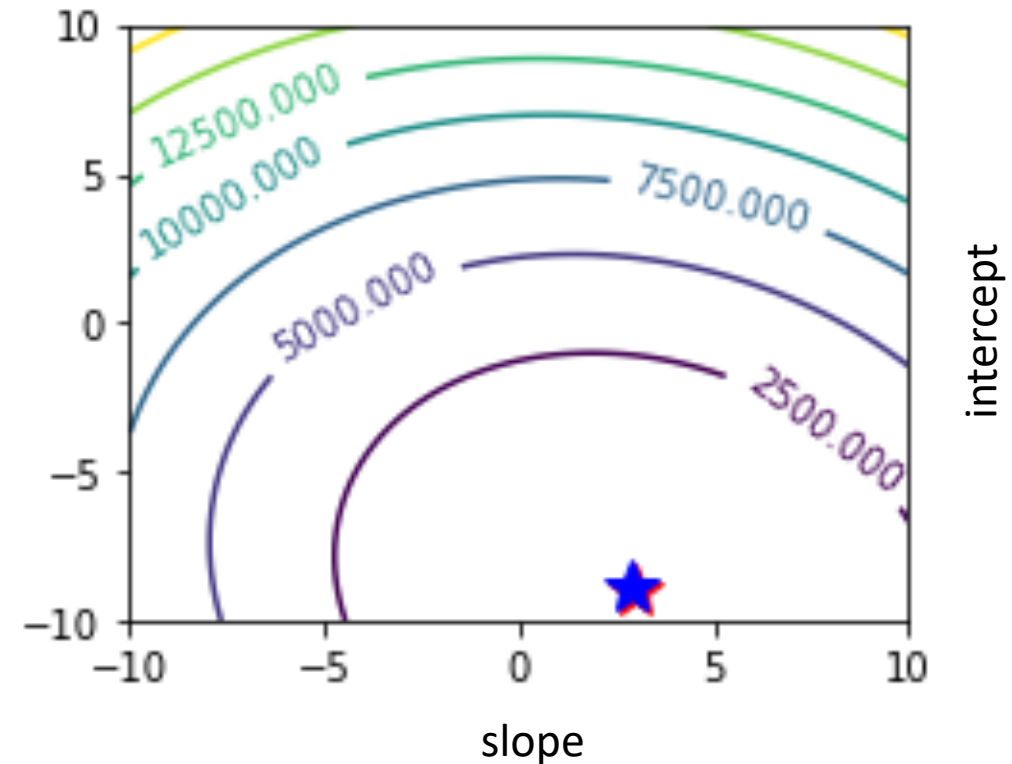
- Gradient descent starts with an initial guess in parameter space:



  – And each iteration tries to move guess closer to solution.

# Gradient Descent in 2D

- Under weak conditions, algorithm converges to a 'w' with $\nabla f(w) = 0$.
  - 'f' is bounded below, $\nabla f$ can't change arbitrarily fast, small-enough constant $\alpha^t$.

# Gradient Descent for Least Squares

- The least squares objective and gradient:

$$f(w) = \frac{1}{2}\|Xw - y\|^2 \qquad \nabla f(w) = X^T(Xw - y)$$

- Gradient descent iterations for least squares:

$$w^{t+1} = w^t - \alpha^t \underbrace{X^T(Xw^t - y)}_{\nabla f(w^t)}$$

lowercase t means iteration step

- Cost of gradient descent iteration is O(nd) (no need to form $X^TX$).

$$\text{Bottleneck is computing } \nabla f(w^t) = X^T(\underbrace{Xw^t - y}_{\substack{O(nd) \\ O(n)}})$$

$$O(nd)$$

# Normal Equations vs. Gradient Descent

- Least squares via normal equations vs. gradient descent:
  - Normal equations cost $O(nd^2 + d^3)$.
  - Gradient descent costs $O(ndt)$ to run for 't' iterations.
    - Each of the 't' iterations costs $O(nd)$.

  - Gradient descent can be faster when 'd' is very large:
    - If solution is "good enough" for a 't' less than $\text{minimum}(d, d^2/n)$.
    - CPSC 5XX: 't' proportional to "condition number" of $X^TX$ (no direct 'd' dependence).

  - Normal equations only solve linear least squares problems.
    - Gradient descent solves many other problems.

# Step Size Considerations

- The step size $\alpha^t$ must be set carefully for gradient descent to work.
  - If $\alpha^t$ is too large gradient descent may not converge.
  - If $\alpha^t$ is too small gradient descent may be too slow.

# Beyond Gradient Descent

- There are many variations on gradient descent.
  - Methods employing a "line search" to choose the step-size.
  - "Conjugate" gradient and "accelerated" gradient methods.
  - Newton's method (which uses second derivatives).
  - Quasi-Newton and Hessian-free Newton methods.
  - Stochastic gradient descent (later in course).

- This course focuses on gradient descent and stochastic gradient:
  - They're simple and give reasonable solutions to most ML problems.
  - But the above can be faster for some applications.

# Next Topic: Convex Functions

# Convex Functions

- Is finding a 'w' with $\nabla f(w) = 0$ good enough?
  - Yes, for convex functions.



- A function is convex if the area above the function is a convex set.
  - All values between any two points above function stay above function.

# Convex Functions

- Example of non-convex function:
  - In this case, we have $\nabla f(w)=0$ at a maximum.

- Example of two-dimensional convex function:

# Convex Functions

- All 'w' with $\nabla f(w) = 0$ for convex functions are global minima.

Proof by contradiction:

Consider a local minimum

If this is not global minimum, there must a smaller value.

But this contradicts that we are at a local minimum.

↑
By convexity we can move along line to global minimum and decrease objective.

— Normal equations find a global minimum because least squares is convex.

# How do we know if a function is convex?

- Some useful tricks for showing a function is convex:
  - 1-variable, twice-differentiable function is convex iff f''(w) ≥ 0 for all 'w'.

Consider $f(w) = \frac{1}{2} a w^2$ for $a > 0$.

We have $f'(w) = aw$

and $f''(w) = a > 0$

By assumption

Consider $f(w) = e^w$

We have $f'(w) = e^w$

and $f''(w) = e^w > 0$

By definition of exponential function.

# How do we know if a function is convex?

- Some useful tricks for showing a function is convex:
  - 1-variable, twice-differentiable function is convex iff f''(w) ≥ 0 for all 'w'.
  - A convex function multiplied by non-negative constant is convex.

We showed that $f(w) = e^w$ is convex, so $f(w) = 10e^w$ is convex.

# How do we know if a function is convex?

- Some useful tricks for showing a function is convex:
  - 1-variable, twice-differentiable function is convex iff f''(w) ≥ 0 for all 'w'.
  - A convex function multiplied by non-negative constant is convex.
  - Linear functions are convex.

$$f(w) = w^T v \quad \text{for any vector 'v' is convex}$$

Area above a line is convex.

Area above a plane or hyper-plane is convex

# How do we know if a function is convex?

- Some useful tricks for showing a function is convex:
  - 1-variable, twice-differentiable function is convex iff $f''(w) \geq 0$ for all 'w'.
  - A convex function multiplied by non-negative constant is convex.
  - Linear functions are convex.
  - Norms and squared norms are convex.

$$\|w\|_2, \quad \|w\|_2^2, \quad \|w\|_1, \quad \|w\|_\infty, \quad \|w\|_1^2, \quad \text{and so on are all convex.}$$
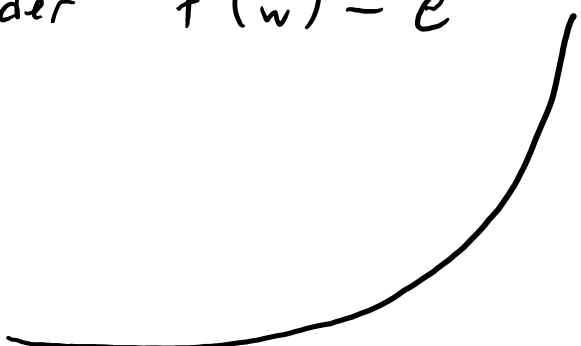
# How do we know if a function is convex?

- Some useful tricks for showing a function is convex:
  - 1-variable, twice-differentiable function is convex iff f''(w) ≥ 0 for all 'w'.
  - A convex function multiplied by non-negative constant is convex.
  - Linear functions are convex.
  - Norms and squared norms are convex.
  - The sum of convex functions is a convex function.

$$f(w) = 10e^w + \frac{\lambda}{2}\|w\|^2 \quad \text{is convex}$$

From earlier

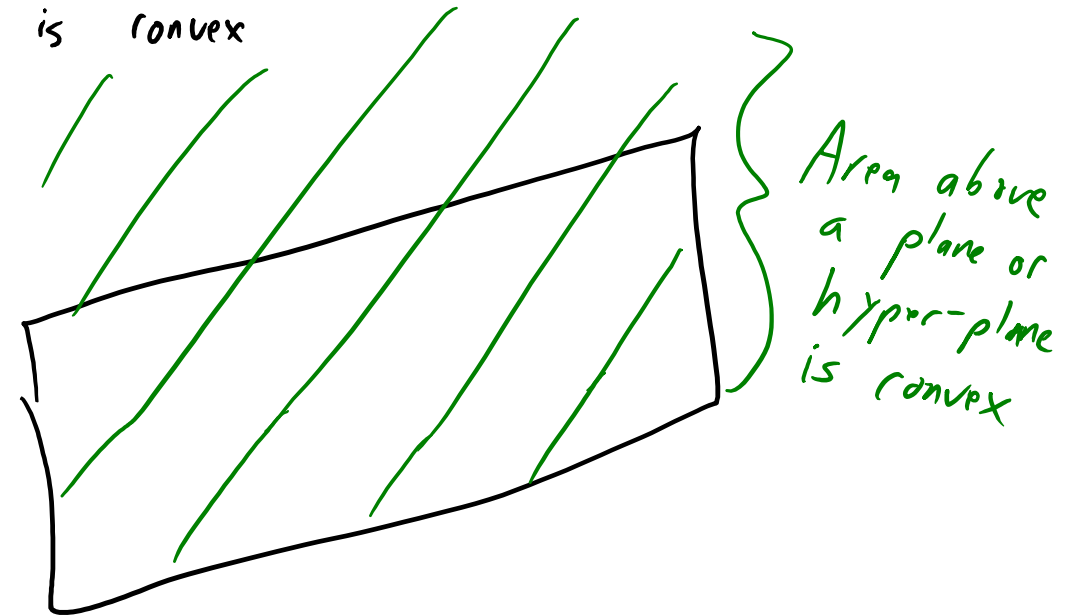constant norm squared

# How do we know if a function is convex?

- Some useful tricks for showing a function is convex:
  - 1-variable, twice-differentiable function is convex iff $f''(w) \geq 0$ for all 'w'.
  - A convex function multiplied by non-negative constant is convex.
  - Linear functions are convex.
  - Norms and squared norms are convex.
  - The sum of convex functions is a convex function.
  - The max of convex functions is a convex function.

$$f(w) = \max\{1, 2w, w^2\} \quad \text{is convex.}$$

convex

$y = w^2$
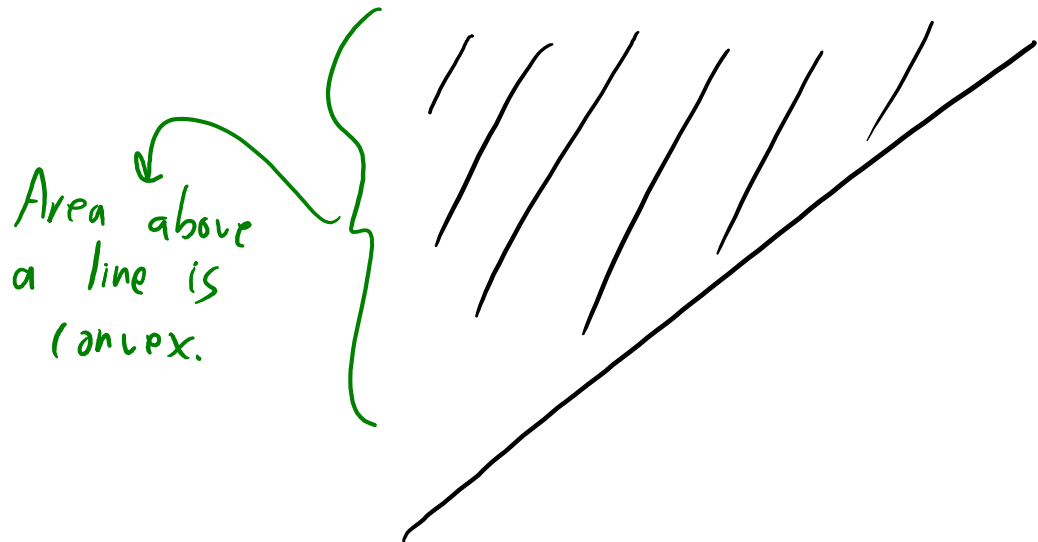
$y = 2w$

$y = 1$

# How do we know if a function is convex?

- Some useful tricks for showing a function is convex:
  - 1-variable, twice-differentiable function is convex iff f''(w) ≥ 0 for all 'w'.
  - A convex function multiplied by non-negative constant is convex.
  - Linear functions are convex.
  - Norms and squared norms are convex.
  - The sum of convex functions is a convex function.
  - The max of convex functions is a convex function.
  - Composition of a convex function and a linear function is convex.

$$\text{If } f(w) = g(Xw - y) \text{ then 'f' is convex if 'g' is convex.}$$

$Xw - y$ → linear function (from $\mathbb{R}^d$ to $\mathbb{R}^n$)

# How do we know if a function is convex?

- Some useful tricks for showing a function is convex:
  - 1-variable, twice-differentiable function is convex iff f''(w) ≥ 0 for all 'w'.
  - A convex function multiplied by non-negative constant is convex.
  - Linear functions are convex.
  - Norms and squared norms are convex.
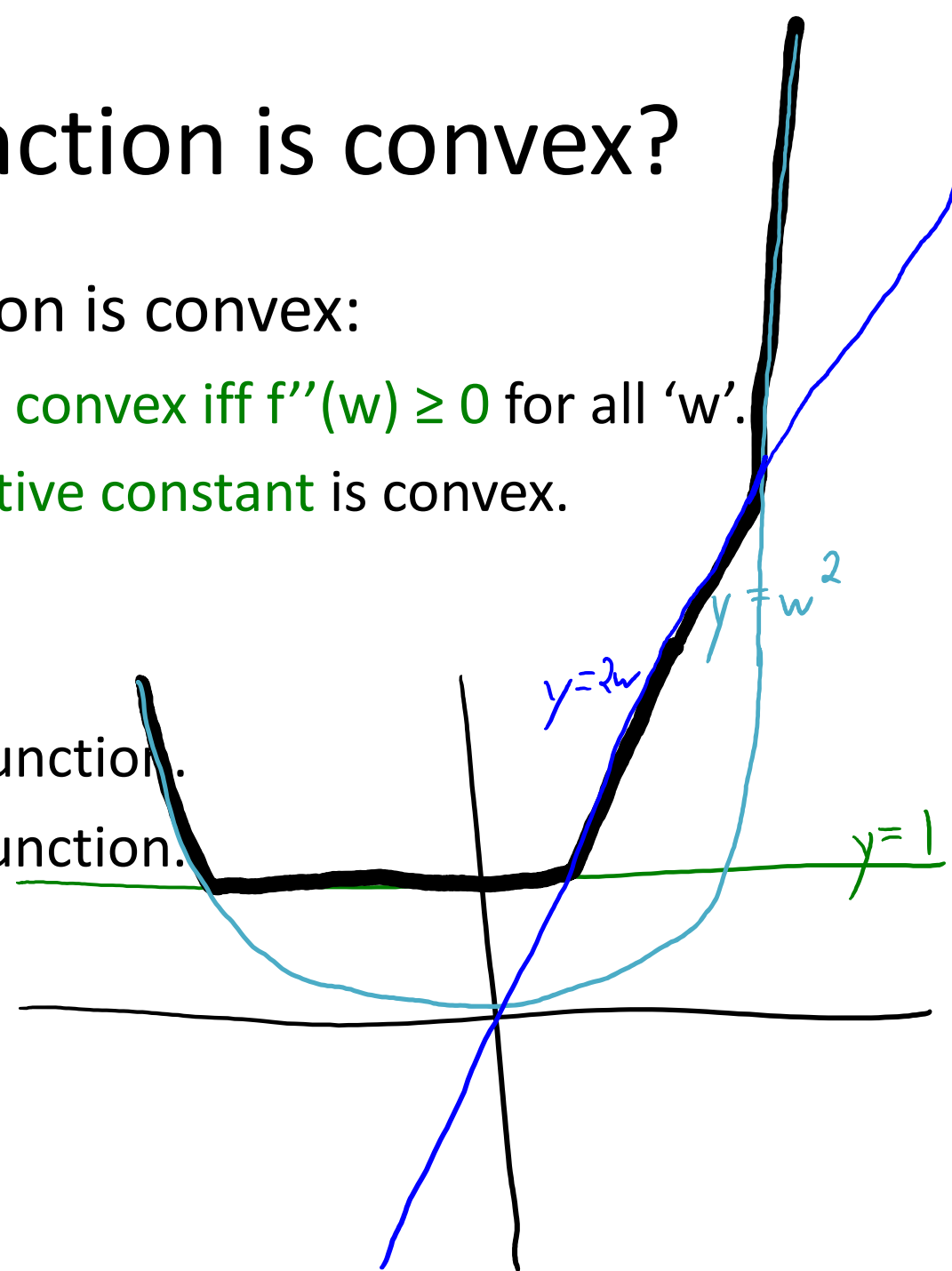  - The sum of convex functions is a convex function.
  - The max of convex functions is a convex function.
  - Composition of a convex function and a linear function is convex.
- But: not true that multiplication of convex functions is convex:
  - If $f(x)=x$ (convex) and $g(x)=x^2$ (convex), $f(x)g(x) = x^3$ (not convex).

# How do we know if a function is convex?

- Some useful tricks for showing a function is convex:
  - 1-variable, twice-differentiable function is convex iff $f''(w) \geq 0$ for all 'w'.
  - A convex function multiplied by non-negative constant is convex.
  - Linear functions are convex.
  - Norms and squared norms are convex.
  - The sum of convex functions is a convex function.
  - The max of convex functions is a convex function.
  - Composition of a convex function and a linear function is convex.
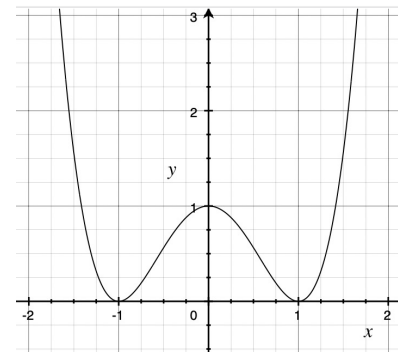
$(w^2 - 1)^2$

- Also not true that composition of convex with convex is convex:

Even if 'f' is convex and 'g' is convex, $f(g(w))$ might not be convex.

E.g., $w^2$ is convex and $(w-1)^2$ is convex, but $(w^2-1)^2$ is not convex.

# Example: Convexity of Linear Regression (Easy Way)

- Consider linear regression objective with squared error:

$$f(w) = \|Xw - y\|^2$$

- We can use that this is a convex function composed with linear:

Let $h(w) = Xw - y$, which is a linear function ('d' inputs, 'n' outputs)

Let $g(r) = \|r\|^2$, which is convex because it's a squared norm.

Then $f(w) = g(h(w))$, which is convex because it's a convex function composed with a linear function

# Convexity in Higher Dimensions



- Twice-differentiable 'd'-variable function is convex iff:
  – Eigenvalues of Hessian $\nabla^2 f(w)$ are non-negative for all 'w'.

- True for least squares where $\nabla^2 f(w) = X^T X$ for all 'w'.
  – See bonus slides for why $X^T X$ has non-negative eigenvalues.

- Unfortunately, sometimes it is hard to show convexity this way.
  – Usually easier to just use some of the rules as we did on the last slide.

# Summary

- Gradient descent finds critical point of differentiable function.
  - Can be faster than normal equations for large 'd' values.
  - Finds global optimum if function is convex.
- Convex functions:
  - Set of functions with property that $\nabla f(w) = 0$ implies 'w' is a global min.
  - Can (usually) be identified using a few simple rules.

- Next time:
  - Linear regression without the outlier sensitivity…

# Norms Norms Norms: Getting from Sums to Norms

I was going over the solutions for A3 and I am still a bit confused on how to get from a sum to a norm in some situations. I know the basic ones that give me $\|Xw - y\|^2$ and stuff, but when other things are thrown in the mix I get a bit confused. For example, $\sum_{i=1}^n v_i(w^T x_i - y_i)^2$ gives $\|V^{1/2}(Xw - y)\|^2$. From my understanding v is a vector and vi is the number at position I in that vector. How does the summation of these indices result in the diagonal matrix V and not just the vector v?

Furthermore, when we have a summation like $\sum_{j=1}^d \lambda_j|w_j|$, it is simplified to $\|\Lambda w\|_1$. How does the lambda end up inside the L1 norm? I thought that a summation could be simplified to a L1 norm if its terms are wrapped around the absolute value symbol. In this case the lambda is not, so how is it able to appear inside the norm like that?

hw3    midterm_exam

---

i **the instructors' answer,** *where instructors collectively construct a single answer*

I know that this notation seems intimidating if this is the first time you see it. Fortunately, there are really only a few "rules" you need to figure out, and you'll find that these are use all over the place.

For those particular questions you'll want to memorize the way that the three common norms appear:
$\sum_{i=1}^n |r_i| = \|r\|_1, \sum_{i=1}^n r_i^2 = \|r\|^2, \max_{i \in \{1,2,\ldots,n\}}\{|r_i|\} = \|r\|_\infty$. So when you see max, sum of non-negative values, or sum of squared values you should think of these norms.

Next, notice what multiplying by a diagonal matrix does: if you multiply a vector $w$ (for example) by a diagonal matrix then you multiply each element $w_i$ by the corresponding diagonal element. If you multiply matrix $X$ (for example) by a diagonal matrix then you multiply each row of $X$ by the corresponding diagonal element.

The $V^{1/2}$ shows up because we're multiplying the square.

The other really useful ones to know are $\sum_{i=1}^n v_i r_i = v^T r$ if the elements aren't necessarily non-negative, $\sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j = x^T A x$, and $\sum_{i=1}^n x_i r_i = X^T r$.

(All of the above follow from definitions, but it takes some practice to recognize these common forms. That's why we made you get some practice on the assignments, and why we covered this notation before the midterm so that you study it before we start using it a lot. It is incredibly common the ML world.)

---

edit  ·  good answer | 2        Updated 8 months ago by Mark Schmidt

# Constraints, Continuity, Smoothness

- Sometimes we need to optimize with constraints:
  - Later we'll see "non-negative least squares".

$$\min_{w \geq 0} \frac{1}{2} \sum_{i=1}^{n} \left( w^\top x_i - y_i \right)^2$$

  - A vector 'w' satisfying w ≥ 0 (element-wise) is said to be "feasible".
- Two factors affecting difficulty are continuity and smoothness.
  - Continuous functions tend to be easier than discontinuous functions.
  - Smooth/differentiable functions tend to be easier than non-smooth.
  - See the calculus review here if you haven't heard these words in a while.

# Convexity, min, and argmin

- If a function is convex, then all critical points are global optima.

- However, convex functions don't necessarily have critical points:
  – For example, $f(x) = a*x$, $f(x) = \exp(x)$, etc.

- Also, more than one 'x' can achieve the global optimum:
  – For example, $f(x) = c$ is minimized by any 'x'.

# Why use the negative gradient direction?

- For a twice-differentiable 'f', multivariable <span style="color:blue">Taylor expansion</span> gives:

$$f\left(w^{t+1}\right) = f(w^t) + \nabla f(w^t)^\top (w^{t+1} - w^t) + \frac{1}{2}(w^{t+1} - w^t)\nabla^2 f(v)(w^{t+1} - w^t)$$

for some 'v' between $w^{t+1}$ and $w^t$.

- If gradient can't change arbitrarily quickly, Hessian is bounded and:

$$f(w^{t+1}) = f(w^t) + \nabla f(w^t)^\top (w^{t+1} - w^t) + O\left(\|w^{t+1} - w^t\|^2\right)$$

becomes negligible as $w^{t+1}$ gets close to $w^t$

  - But which <span style="color:green">choice of $w^{t+1}$ decreases 'f' the most</span>?

    - As $\|w^{t+1} - w^t\|$ gets close to zero, the value of $w^{t+1}$ minimizing $f(w^{t+1})$ in this formula converges to $(w^{t+1} - w^t) = -\alpha^t \nabla f(w^t)$ for some scalar $\alpha^t$.
    - So if we're moving a small amount, the optimal $w^{t+1}$ is: $w^{t+1} = w^t - \alpha_t \nabla f(w^t)$ for some scalar $\alpha_t$.

# Normalized Steps

Question from class: "can we use $w^{t+1} = w^t - \frac{1}{\|\nabla f(w^t)\|} \nabla f(w^t)$ "

This will work for a while, but notice that

$$\|w^{t+1} - w^t\| = \|\frac{1}{\|\nabla f(w^t)\|} \nabla f(w^t)\|$$

$$= \frac{1}{\|\nabla f(w^t)\|} \|\nabla f(w^t)\|$$

$$= 1$$

So the algorithm <u>never converges</u>

# Optimizer "findMin" Details

## The minimizer function

Hi all,

I'm just curious how the minimizers given to us works. Are there any resources can give us more details about it?

**i the instructors' answer,** *where instructors collectively construct a single answer*

It's just a basic gradient descent implementation with some clever guesses for the step-size.

The step-size on each iteration is initialized using the method from this classic paper (which works surprisingly well but we don't really know why except in two dimensions):
http://pages.cs.wisc.edu/~swright/726/handouts/barzilai-borwein.pdf

That step-size is evaluated using a standard condition ("Armijo condition") and then it fits a polynomial regression model based on the function and directional derivative values and tries the step-size minimizing this polynomial. Both these tricks are described in Nocedal and Wright's "Numerical Optimization" book.

# 3.3.8

Probably a stupid question, but when we do the gradient descent, is it required to calculate the loss value at each step? Are we allowed to save some middle results and reused them if we calculate both gradient and the loss? Thank you!

hw1

**~ An instructor (Mark Schmidt) endorsed this question ~**

Edit    undo good question | 1                Updated 8 months ago by          (Anon. Scale to class

---

ℹ **the instructors' answer,** *where instructors collectively construct a single answer*

This is an important practical detail when you implement gradient descent, although there are few ways to interpret the question.

1. Technically, gradient descent never needs to compute the loss value. But you often compute it because it can help in checking whether you have a good step size.

2. Usually, some of the calculations involved in computing the gradient are also need to compute the function value. For example, a function of the form $f(Xw)$ has a gradient of the form $X^T f'(Xw)$, so if you evaluate the function and its gradient at the same 'w' you only want to compute $Xw$ once. The function *findMin* asks for the function and gradient together so you can do things like this.

3. When you do a line-search, you often compute $f(\tilde{w})$ and possibly $\nabla f(\tilde{w})$ for some values $\tilde{w}$ that may later become $w^{k+1}$. You can use these quantities in the next iteration, instead of re-computing them. In *findMin*, it assumes that the line-search from the previous iteration gives you the current function and gradient value, so these do not need to be computed.

4. For many problems, there are tricks available where you only compute the loss/gradient with respect to a subset of the examples on each iteration, but then you still do the full update:
https://arxiv.org/pdf/1309.2388v2.pdf

# Example: Convexity of Linear Regression (Hard Way)

- Consider linear regression objective with squared error:

$$f(w) = \|Xw - y\|^2$$

- Twice-differentiable 'f' is convex if $\nabla^2 f(x)$ has eigenvalues $\geq 0$.
  - This is equivalent to saying $v^T \nabla^2 f(x) v \geq 0$ for all vectors $v$.
- The Hessian for least squares is $\nabla^2 f(x) = X^T X$.
  - See notes on Gradients and Hessians of quadratics on webpage.
- We have:

$$v^T \nabla^2 f(w) v = v^T X^T X v = (Xv)^T (Xv) = \|Xv\|^2 \geq 0 \quad (\text{because norms are} \geq 0)$$

$$\text{So it's convex}$$