

CPSC 540: Machine Learning

Convex Functions, Gradient Descent, Convergence Rates

Winter 2016

Admin

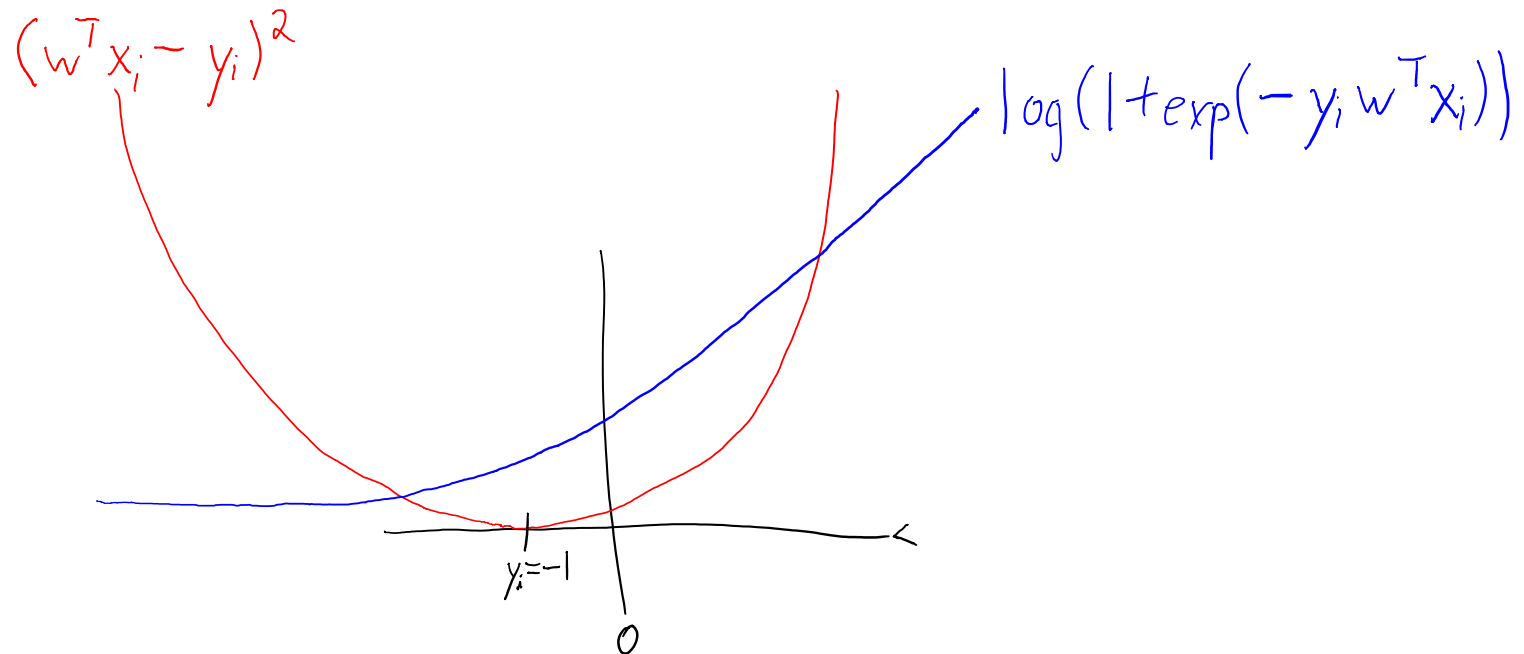
- Auditing/enrollment forms:
 - Drop-off/pickup your forms at the end of class.
 - It will be easier to argue for larger classroom if people are officially enrolled/auditing.
 - Remaining forms can be picked up at the **tutorials tomorrow**.
- CPSC and EECE graduate students: **prereq forms due now**.
- Assignment 1: **due Tuesday**.
 - Hand in **one assignment for the group** (of 1-3).
- Add/Drop deadline: **Monday**.
 - Last chance before you are locked in/out.

The 'Best' Machine Learning Model

- What is the 'best' machine learning model?
 - SVMs? Random forests? Deep learning?
- No free lunch theorem:
 - There is no 'best' model that achieves the best test error for every problem.
 - If model A works better than model B on one dataset, there is another dataset where model B works better.
- Asking what is the 'best' machine learning model is like asking which is 'best' among "rock", "paper", and "scissors".
- Caveat of no free lunch (NFL) theorem:
 - The world is very structured, some datasets are more likely than others.
 - Model A could be better than model B on a huge variety of practical applications.
- Machine learning emphasizes models useful across applications.

Last Time: Logistic Regression

- We considered **binary labels** y_i , and classifying with $\text{sign}(w^T x_i)$.
 - Squared error $(w^T x_i - y_i)^2$ is not ideal: **penalizes model for “too right”**.
 - Minimizing number of errors is also not ideal: **NP-hard**.
 - Tractable upper bounds are **hinge loss and logistic loss**.



Last Time: Maximum Likelihood and MAP

- Minimizing a loss function often equivalent to **maximum likelihood**.
 - For example, least squares is equivalent to using a Gaussian likelihood:

$$\text{If } y_i \sim \mathcal{N}(w^T x_i, \sigma^2), \quad \underset{w \in \mathbb{R}^d}{\operatorname{argmax}} \underbrace{p(y|w, X)}_{\text{likelihood}} \iff \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \frac{1}{2} \|Xw - y\|^2$$

- With a regularizer, often equivalent to MAP estimation:
 - For example, L2-regularization is equivalent to using a Gaussian prior:

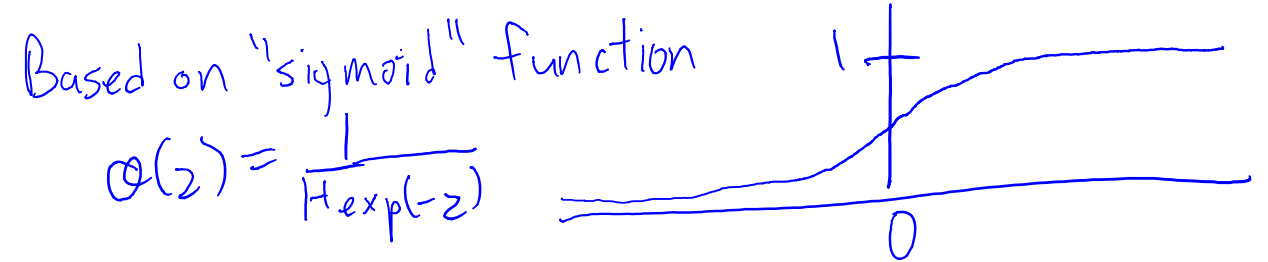
$$\text{If } y_i \sim \mathcal{N}(w^T x_i, \sigma^2) \quad \text{and} \quad w_j \sim \mathcal{N}(0, \lambda^{-1}) \quad \underset{w \in \mathbb{R}^d}{\operatorname{argmax}} \underbrace{p(w|y, X)}_{\text{posterior}} \iff \underset{w \in \mathbb{R}^d}{\operatorname{argmax}} \underbrace{p(y|w, X)}_{\text{likelihood}} \underbrace{p(w)}_{\text{prior}} \iff \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \frac{1}{2} \|Xw - y\|^2 + \frac{\lambda}{2} \|w\|^2$$

- Gives **probabilistic perspective on regularization**: prior on 'w'.

Last Time: Maximum Likelihood and MAP

- Logistic loss is equivalent to **maximum likelihood** logistic regression:

$$p(y_i | w, x_i) = \frac{1}{1 + \exp(-y_i w^T x_i)}$$



- L2-regularized logistic is **MAP estimate** with Gaussian prior:

$$p(y_i | w, x_i) = \frac{1}{1 + \exp(-y_i w^T x_i)} \quad p(w_j | \lambda) \propto \exp\left(-\frac{\lambda}{2} w_j^2\right)$$

- Advantage of likelihood/MAP perspective:
 - Allows us to **define objectives for other distributions** of y_i .

Multi-Class Logistic Regression

- Supposed y_i takes values from an **unordered discrete set** of classes.

$$x_i \in \mathbb{R}^d$$

$$y_i \in \{1, 2, 3, 4, \dots, k\}$$



→ "itchy"
(class 1)



→ "scratchy"
(class 2)



→ "poochie"
(class 3)

- Standard model:
 - Use a 'd'-dimensional weight **vector** ' w_c ' for each class 'c'.
 - Try to make **inner-product** $w_c^T x_i$ **big** when 'c' is the true label ' y_i '.
 - Classify by finding largest inner-product: $\hat{y}_i = \operatorname{argmax}_c \{w_c^T x_i\}$

Multi-Class Logistic Regression

We have a parameter matrix $W = \begin{bmatrix} | & | & | & \dots & | \\ w_1 & w_2 & w_3 & \dots & w_k \\ | & | & | & \dots & | \end{bmatrix}$

To make a prediction, compute $W^T x_i = \begin{bmatrix} w_1^T x_i \\ w_2^T x_i \\ \vdots \\ w_k^T x_i \end{bmatrix}$ and compute maximum.
prediction "2"

We want a loss function that will make $w_c^T x_i$ big when c is the true label y_i and will otherwise make $w_c^T x_i$ small.

We can define probability using softmax function:

$$p(y_i = c | W, x_i) = \frac{\exp(w_c^T x_i)}{\sum_{c'=1}^k \exp(w_{c'}^T x_i)} \propto \exp(w_c^T x_i)$$

To fit model, use

$$-\log p(y_i | W, x_i) = -w_{y_i}^T x_i + \log \left(\sum_{c'=1}^k \exp(w_{c'}^T x_i) \right)$$

If $k=3$:

$$p(y_i | W, x_i) = \frac{\exp(w_{y_i}^T x_i)}{\exp(w_1^T x_i) + \exp(w_2^T x_i) + \exp(w_3^T x_i)}$$

↑ true

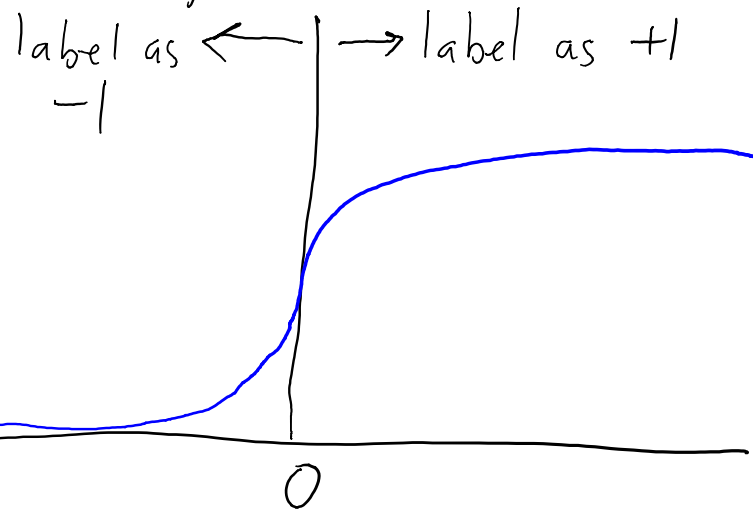
Generalizes sigmoid:

- special case of $k=2$ and $w_2 = 0$

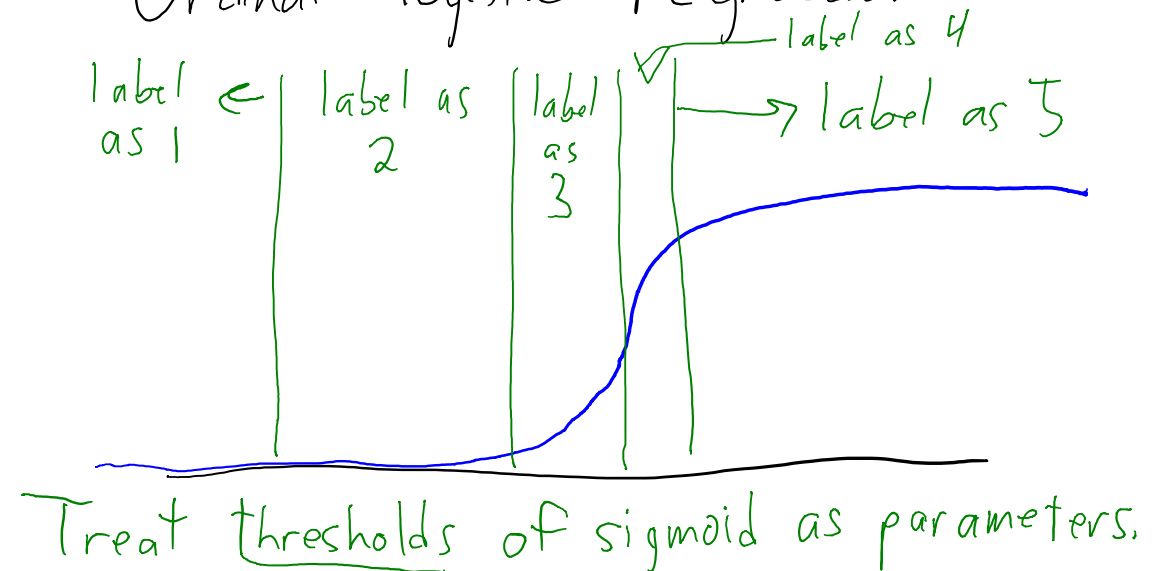
Ordinal Labels

- **Ordinal data**: categorical data where the **order matters**:
 - Rating hotels as {'1 star', '2 stars', '3 stars', '4 stars', '5 stars'}.
 - **Softmax would ignore order**.
- **'Proportional odds' or 'ordinal logistic regression'**:

Logistic regression



Ordinal logistic regression

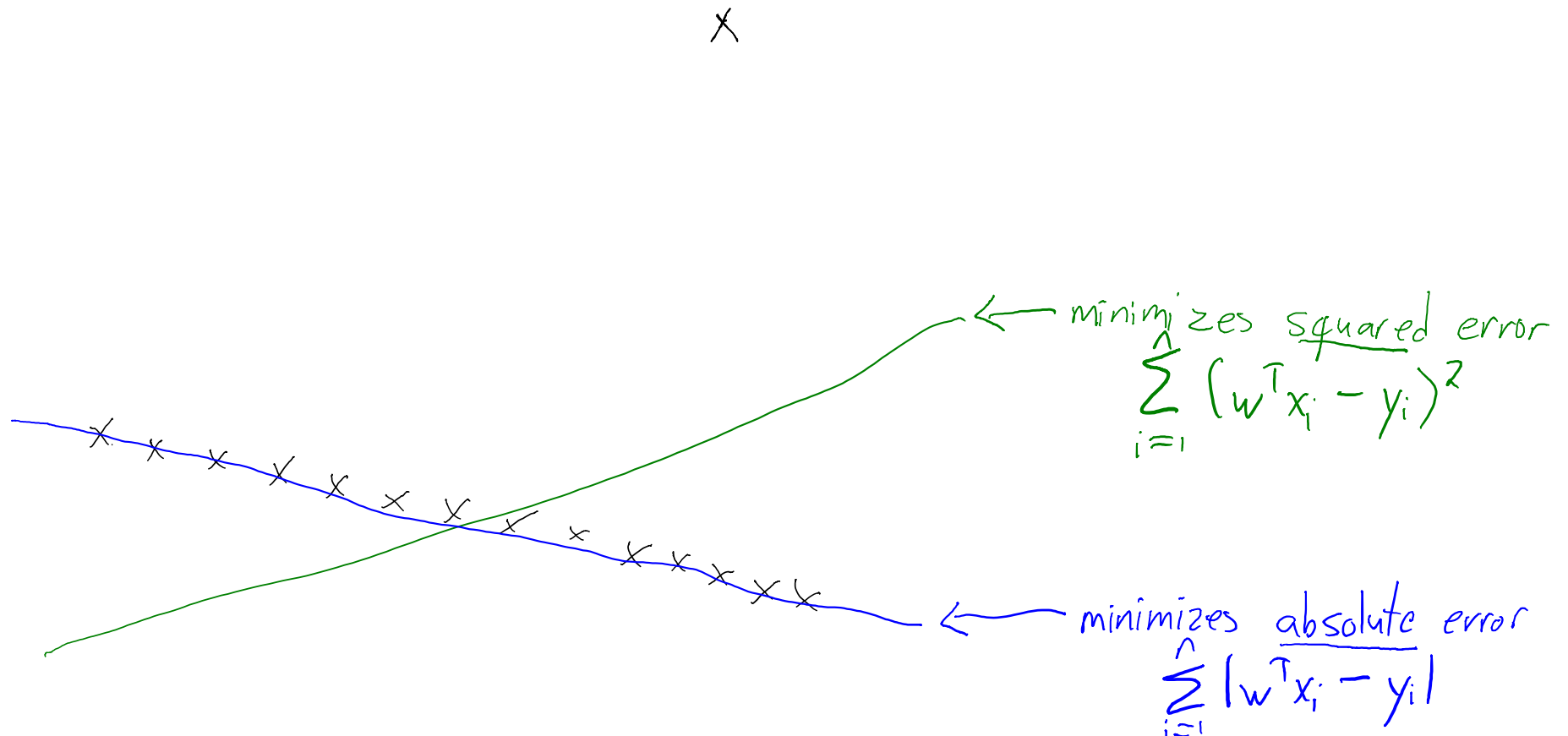


Count Labels

- **Count data**: predict the **number of times** something happens.
 - For example, $y_i = "602"$ Facebook likes.
- Softmax/ordinal **require finite number of categories**.
- We probably don't want separate parameter for '654' and '655'.
- **Poisson regression**: use probability from Poisson count distribution.
 - Many variations exist.

Last Time: Robust Regression

- We said that squared error is **sensitive to outliers**:
 - **Absolute error** is less sensitive: can be solved as a linear program.



'Brittle' Regression

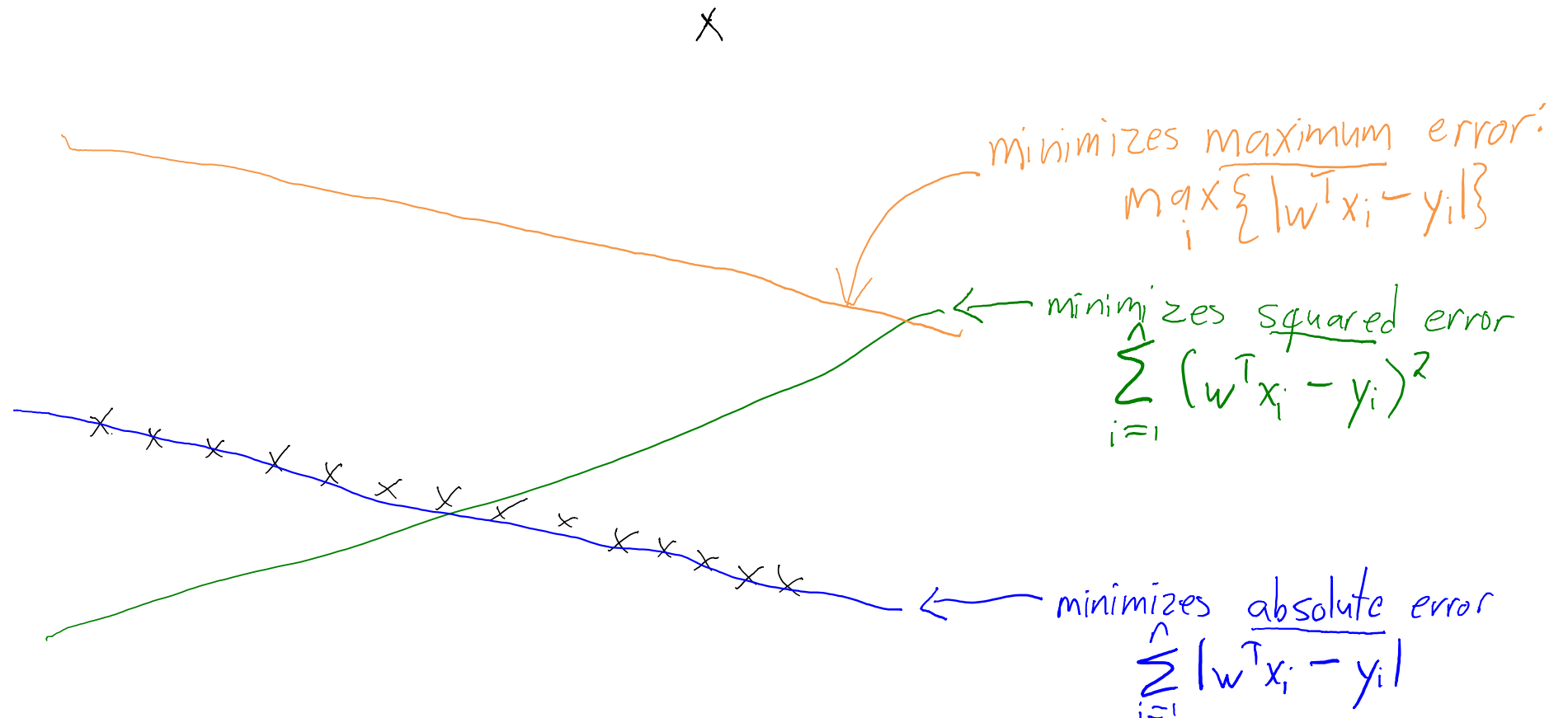
- What if you really care about getting the outliers right?
 - You want **best performance on worst training example**.
 - For example, if in worst case the plane can crash.
- In this case you can use something like the **infinity-norm**:

$$\arg \min_{x \in \mathbb{R}^d} \|Xw - y\|_{\infty} \quad \|z\|_{\infty} = \max_i \{|z_i|\}$$

- Very sensitive to outliers (brittle), but worst case will be better.

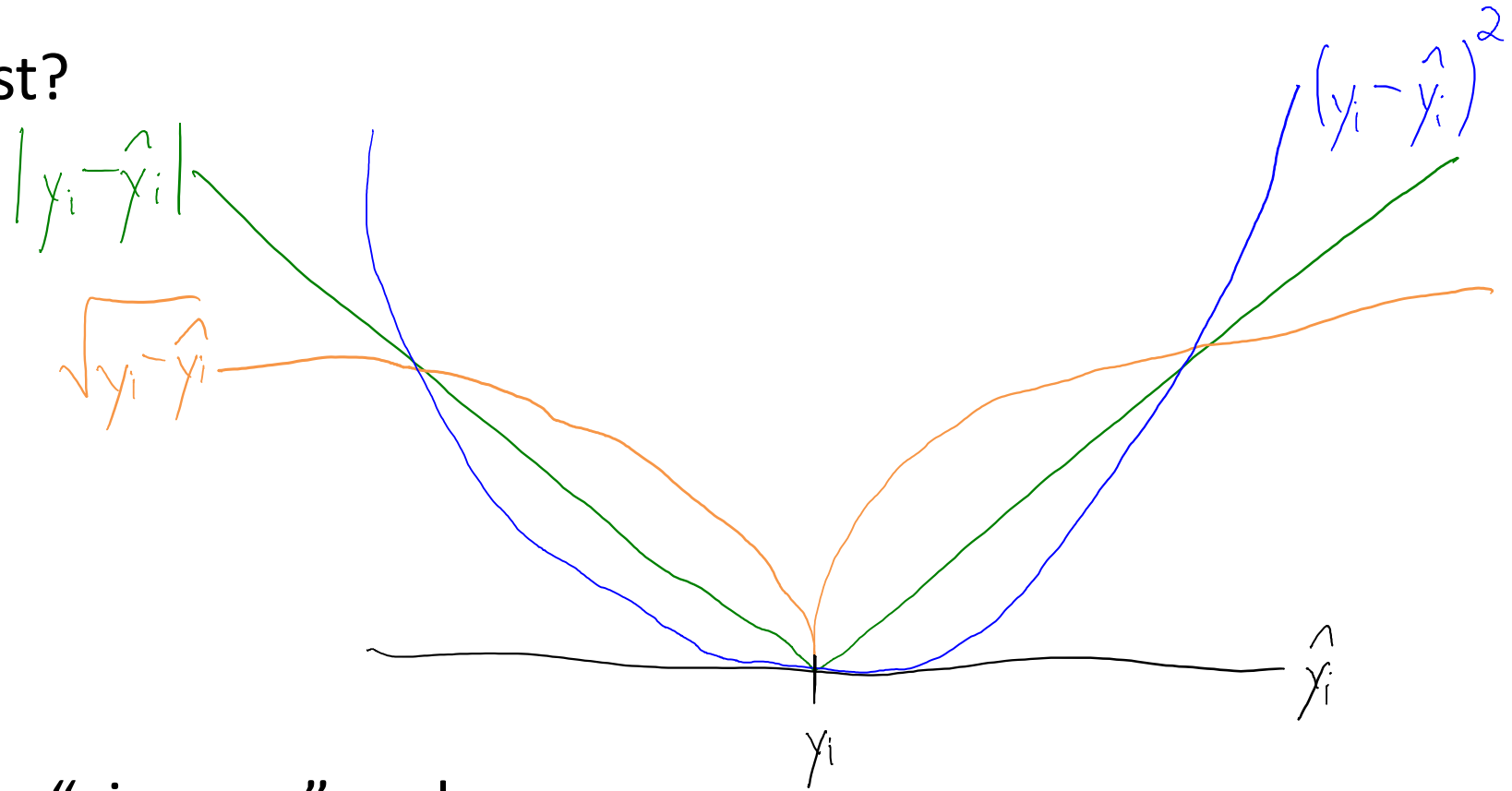
Last Time: Robust Regression

- We said that squared error is **sensitive to outliers**:
 - **Absolute error** is less sensitive: can be solved as a linear program.
 - **Maximum error** is more sensitive: can also be solved as linear program.



Very Robust Regression

- Can we be more robust?



- **Very robust:** eventually “gives up” on large errors.
- But **finding optimal ‘w’ is NP-hard.**
 - Absolute value is the most robust that is not NP-hard.

Course Roadmap

- Topics we discussed in part 1:
 - **Linear models**: change of basis, regularization, loss functions.
 - **Basics of learning theory**: Training vs. test error, bias-variance, fundamental trade-off, no free lunch.
 - **Probabilistic learning principles**: Maximum likelihood, MAP estimation.
- Part 2: **Large-scale machine learning**.
 - Why are SVMs/logistic **easy** while minimizing number of errors is **hard**?
 - How do we fit these models to **huge datasets**?

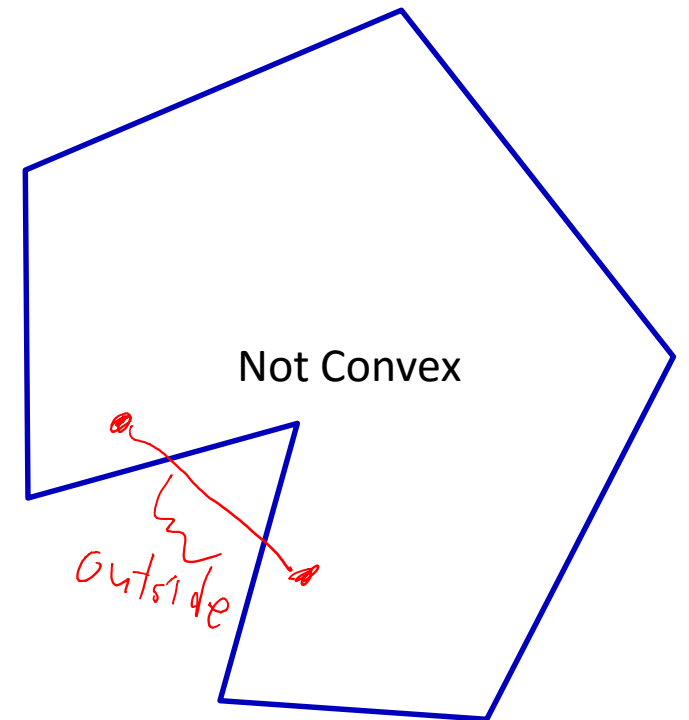
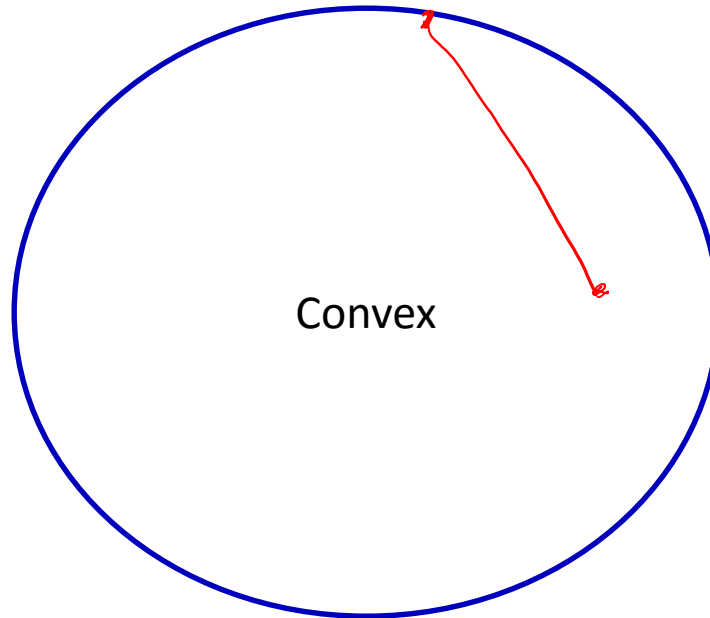
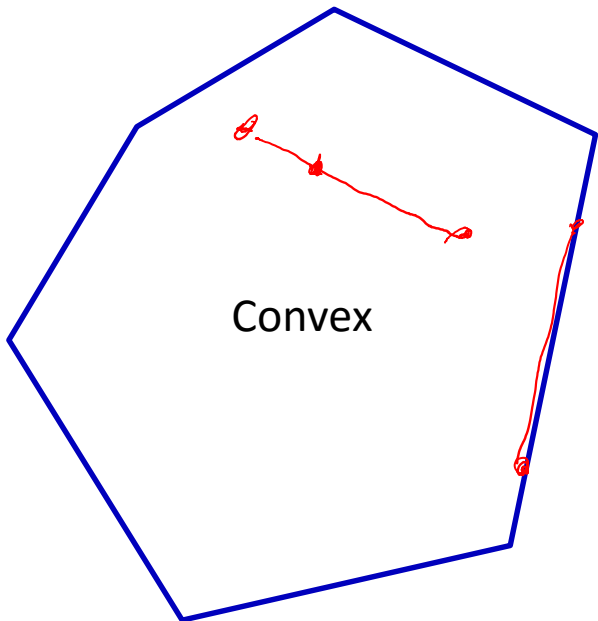
Convex Functions

- We are first going to discuss **convex functions**:
 - Minimizing convex functions is usually easy.
 - Minimizing non-convex functions is usually hard.
- The ‘easy’ problems we have discussed are **convex**:
 - Least squares, robust regression, logistic regression, support vector machines, multi-class logistic, brittle regression, Poisson regression.
 - All of the above with L2-regularization.
- The ‘hard’ problems we have discussed are **non-convex**:
 - 0-1 loss, “very robust” regression.

Convex Sets

- First we need to define a **convex set**:
 - A set is **convex** if the line between any two points stays in the set.

For all $x \in C$ and $y \in C$ we have $\theta x + (1 - \theta)y \in C$ for $0 \leq \theta \leq 1$



Convex Sets

- Examples:

Real-space: \mathbb{R}^d

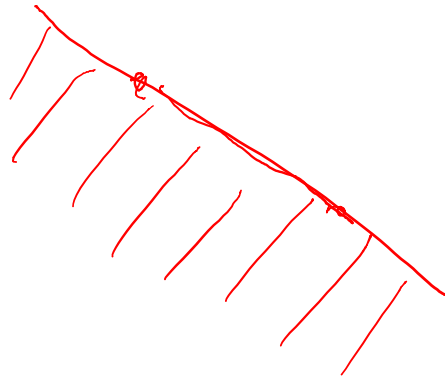
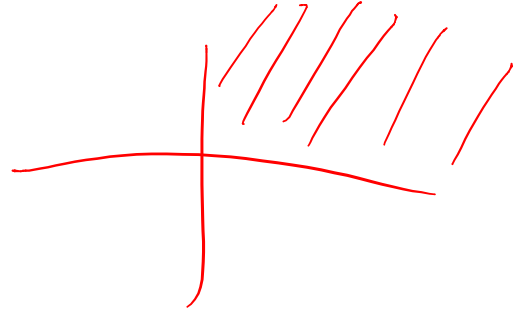
Positive orthant \mathbb{R}_+^d : $\{x \mid x \geq 0\}$

Hyper-plane: $\{x \mid a^T x = b\}$

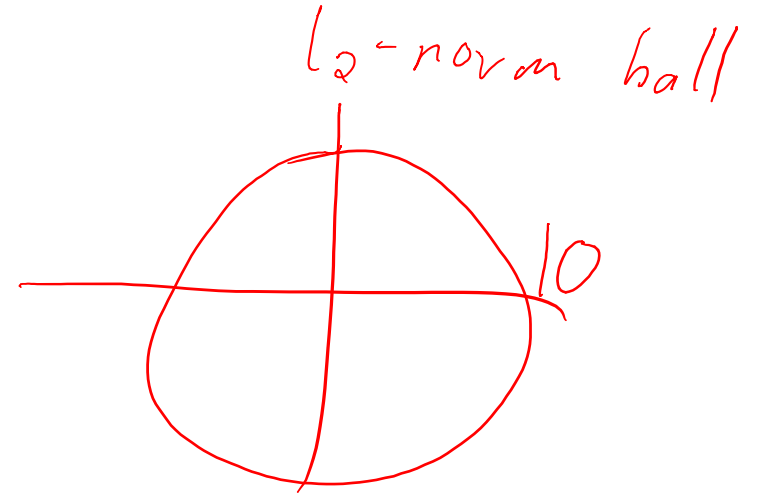
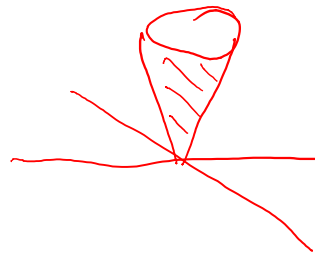
Half-space: $\{x \mid a^T x \leq b\}$

Norm-ball: $\{x \mid \|x\| \leq r\}$

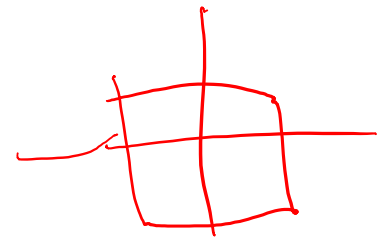
Norm-cone: $\{(x, r) \mid \|x\| \leq r\}$



L_2 -norm cone



L_∞ -norm "ball"



Showing a Set is Convex

How to prove a set is convex?

- One way: choose two generic x and y in the set, show that generic z between them is also in the set.

- Another way: show that set is intersection of sets that you know are convex.

E.g. if $C = \{x \mid a^T x = b\}$

then for $x \in C$ and $y \in C$ and $0 \leq \theta \leq 1$

$$\begin{aligned} \text{we have } a^T(\theta x + (1-\theta)y) &= \theta(a^T x) + (1-\theta)(a^T y) \\ &= \theta b + (1-\theta)b = b \end{aligned}$$

E.g., if $C = \{x \mid \|x\| \leq 10\}$

then for $x \in C$ and $y \in C$ and $0 \leq \theta \leq 1$

$$\text{we have } \|\theta x + (1-\theta)y\|$$

$$\leq \|\theta x\| + \|(1-\theta)y\| \quad (\text{triangle inequality})$$

$$= |\theta| \|x\| + |1-\theta| \|y\| \quad (\text{homogeneity})$$

$$= \theta \|x\| + (1-\theta) \|y\| \quad \theta \geq 0$$

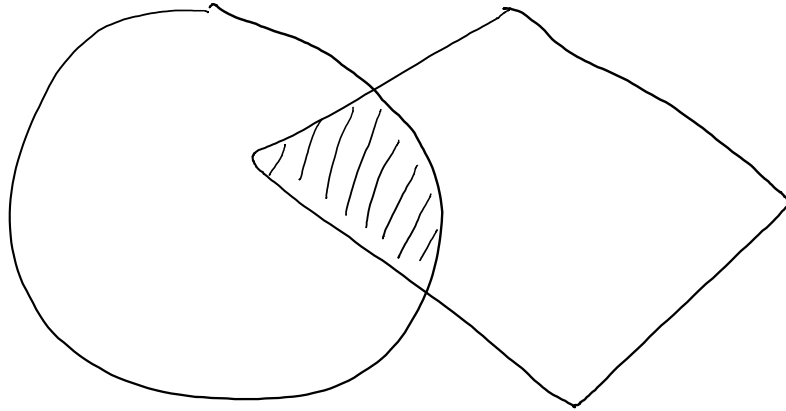
$$\leq \theta \max\{\|x\|, \|y\|\} + (1-\theta) \max\{\|x\|, \|y\|\}$$

$$= \max\{\|x\|, \|y\|\} \leq 10$$

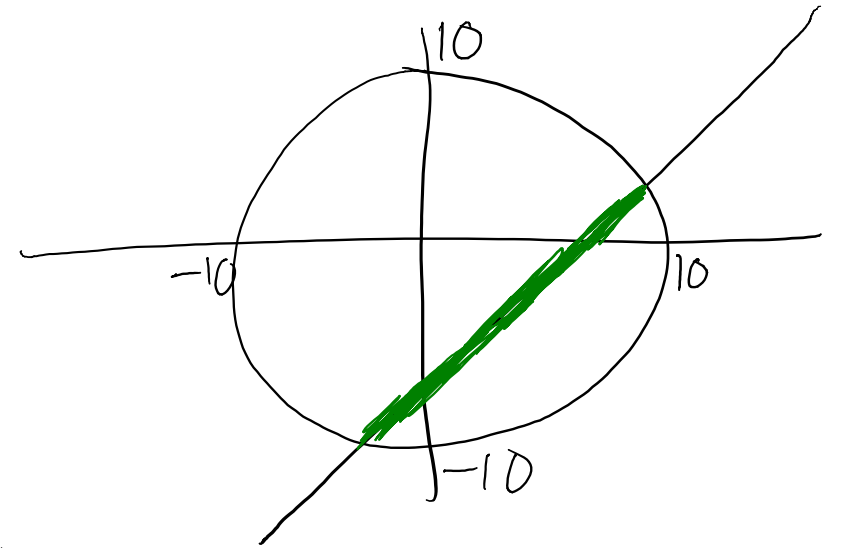
$$a \leq \max\{a, b\}$$

Intersection of Convex Sets

- Intersection of convex sets is convex:



E.g., $\{x \mid a^T x = b\} \cap \{x \mid \|x\| \leq 10\}$
is a convex set.



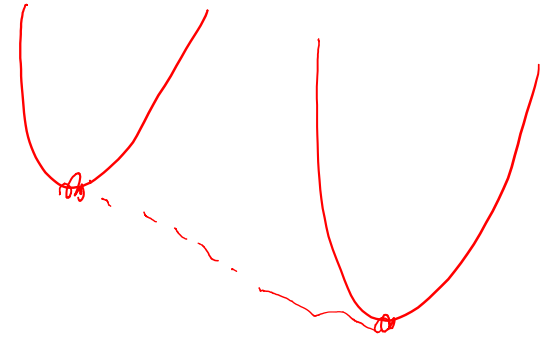
For example, x satisfying linear program constraints are a convex set:

$$Ax \leq b$$

$$A_{eq} x = b_{eq}$$

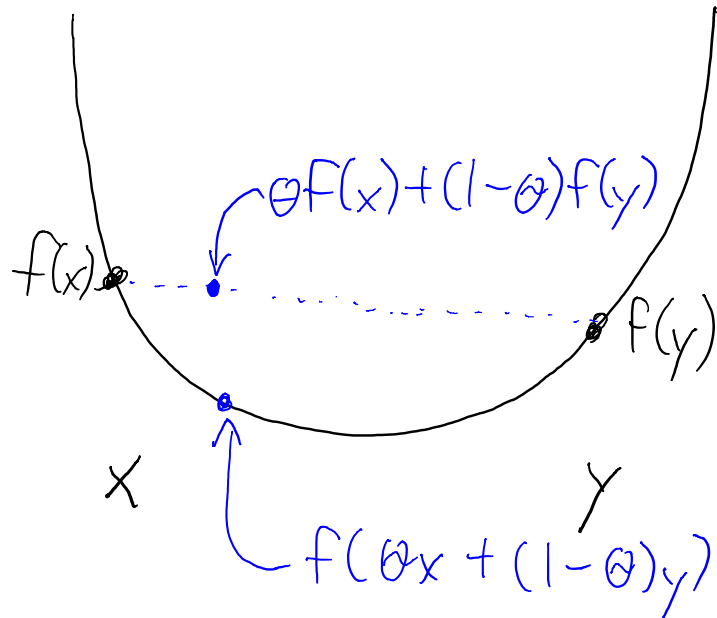
$$LB \leq x \leq UB$$

Convex Functions



- A function 'f' is convex if:
 1. The domain of 'f' is a convex set.
 2. The function is always below 'chord' between two points.

$$f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y) \quad \text{for all } x \in C, y \in C, \text{ and } 0 \leq \theta \leq 1$$



Implication: all local minima
are global minima.

We can minimize a convex function
by finding any stationary point.

Convex Functions

- **Examples:**

Quadratic functions: $f(x) = ax^2 + bx + c, a > 0$.

Linear functions $f(x) = a^T x + b$

Exponential: $f(x) = \exp(ax)$

Negative logarithm: $f(x) = -\log(x)$

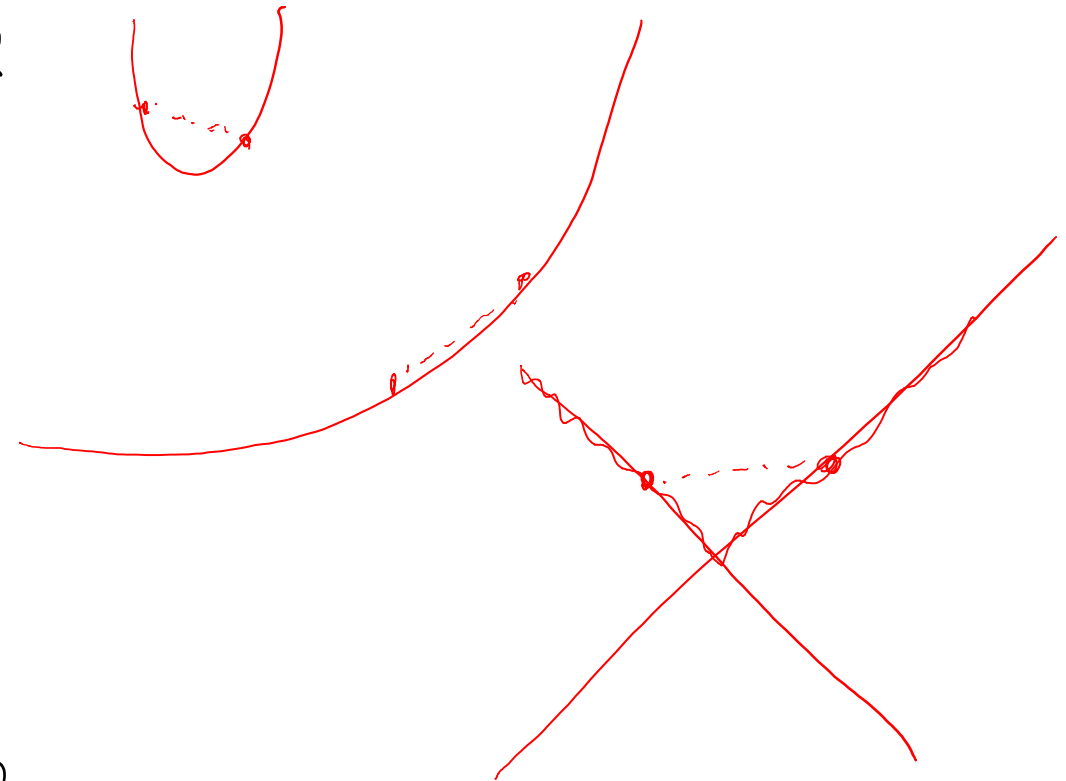
Absolute value: $f(x) = |x|$

Max function: $f(x) = \max_i \{x_i\}$

Negative entropy: $f(x) = x \log(x), x > 0$

Logistic loss: $f(x) = \log(1 + \exp(-z))$

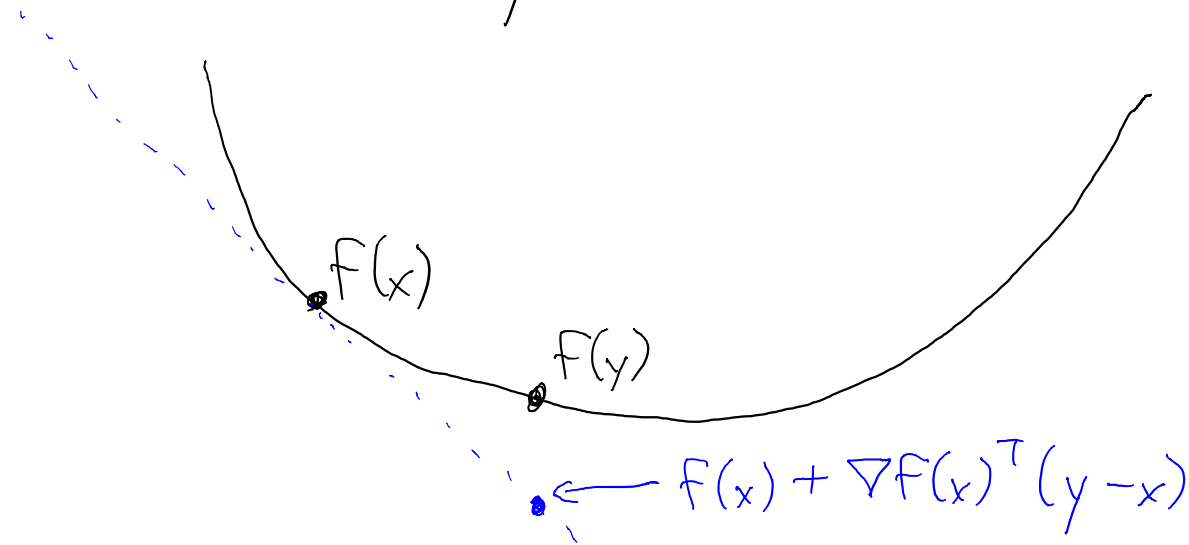
Log-sum-exp: $f(x) = \log\left(\sum_{i=1}^d \exp(x_i)\right)$



Differentiable Convex Functions

- A differentiable 'f' is **convex** iff 'f' is always above tangent:

$$f(y) \geq f(x) + \nabla f(x)^T (y - x) \quad \text{for all } x \in C \text{ and } y \in C$$



If $\nabla f(x) = 0$, this implies $f(y) \geq f(x)$ for all y so x is a global minimizer.

Twice-Differentiable Convex Functions



- A twice-differentiable 'f' is convex iff it's **curved upwards everywhere**.

For one-dimensional functions, reduces to $f''(x) \geq 0$.



— usually, this is the easiest way to show a function is convex.

For multivariate functions, generalization is

$$\nabla^2 f(x) \succeq 0 \text{ for all } x \in C.$$

$\succeq 0$ means A is symmetric and positive semi-definite: $y^T A y \geq 0$ for all y

Showing Functions are Convex

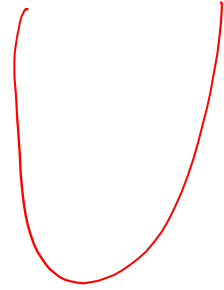
- Examples:

$$\text{If } f(x) = x^2$$

$$\text{then } f'(x) = 2x$$

$$\text{and } f''(x) = 2.$$

Since $2 \geq 0$ we've shown x^2 is convex.



$$\text{If } f(x) = \frac{1}{2} x^T A x + b^T x + c \text{ with } A \succeq 0$$

$$\text{then } \nabla f(x) = A x + b$$

$$\text{and } \nabla^2 f(x) = A$$

Since $\nabla^2 f(x) \succeq 0$ we've shown $f(x)$ is convex.



Showing Functions are Convex

- Examples:

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

Want to show that $\nabla^2 f(w) \succeq 0$,
or equivalently $y^T \nabla^2 f(w) y \geq 0$.

We have $y^T \nabla^2 f(w) y = y^T X^T X y$

$$= (Xy)^T (Xy)$$

$$= \|Xy\|^2 \geq 0.$$

So least squares is

convex and setting

$\nabla f(w) = 0$ gives global
minimum.

Strictly-Convex Functions

- A function is **strictly-convex** if these inequalities strictly hold:

$$f(\theta x + (1-\theta)y) < \theta f(x) + (1-\theta)f(y) \quad \text{for } 0 < \theta < 1.$$

$$f(y) > f(x) + \nabla f(x)^T (y - x)$$

$$\nabla^2 f(x) \succ 0 \quad (y^T \nabla^2 f(x) y > 0 \text{ for all } y \neq 0)$$

- Strict convexity implies **at most one global minimum**:

Points 'x' and 'y' can't both be global minima if $x \neq y$, since

This would imply $f(\theta x + (1-\theta)y)$ is below global min.

for $y \neq 0$

- This implies **L2-regularized least squares has unique solution**:

$$y^T \nabla^2 f(w) y = y^T (X^T X + \lambda I) y = y^T X^T X y + y^T (\lambda I) y = (Xy)^T (Xy) + \lambda y^T y = \|Xy\|^2 + \lambda \|y\|^2 > 0.$$

Operations that Preserve Convexity

- There are a few **operations preserve convexity**.
 - Often lets us avoid calculating Hessian.
 - Often lets us prove convexity of non-smooth functions.
- If f_1 and f_2 are convex, then convexity is preserved under:

1. Non-negative weighted sum:

$$f(x) = z_1 f_1(x) + z_2 f_2(x) \text{ is convex if } z_1 \geq 0 \text{ and } z_2 \geq 0$$

2. Composition with affine function:

$$f(x) = f_1(Ax + b) \text{ is convex.}$$

3. Pointwise maximum:

$$f(x) = \max\{f_1(x), f_2(x)\} \text{ is convex.}$$

Example: SVMs

$$f(x) = \sum_{i=1}^n \max\{0, | -y_i w^T x_i | \} + \frac{\lambda}{2} \|w\|^2$$

linear (convex) linear (convex)

max (convex)

so convex

non-negative sum of convex

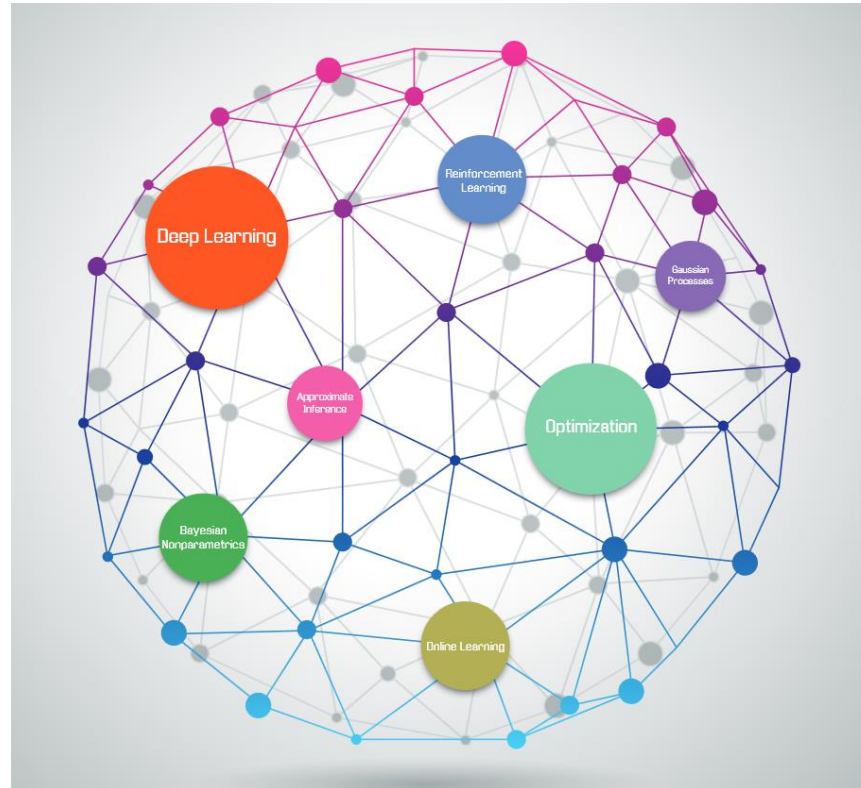
$$\nabla^2 \left[\frac{\lambda}{2} \|w\|^2 \right] = \lambda I \succeq 0$$

so convex

(pause)

Current Hot Topics in Machine Learning

- Graph of most common keywords among ICML papers last year:



- Why is there so much focus on **deep learning and optimization**?

Why Study Optimization in CPSC 540?

- In machine learning, **training is typically written as optimization**:
 - Numerically optimize parameters of model, given data.
- There are some exceptions:
 1. Counting- and distance-based methods (random forests, KNN).
 - See CPSC 340.
 2. Integration-based methods (Bayesian learning).
 - Later in course.

Although you still need to tune parameters in those models.
- But why study optimization? Can't I just use Matlab functions?
 - '\', linprog, quadprog, fmincon, CVX,...

The Effect of Big Data and Big Models

- **Datasets are getting huge**, we might want to train on:
 - Entire medical image databases.
 - Every webpage on the internet.
 - Every product on Amazon.
 - Every rating on Netflix.
 - All flight data in history.
- With bigger datasets, we can build bigger models:
 - This is where **deep learning** comes in.
 - Complicated models can address complicated problems.
- **Now optimization becomes a problem because of time/memory:**
 - We can't afford $O(d^2)$ memory, or an $O(d^2)$ operation.
 - Going through huge datasets 100s of times is too slow.
 - Evaluating huge models too many times is too slow.

Fitting Logistic Regression Models

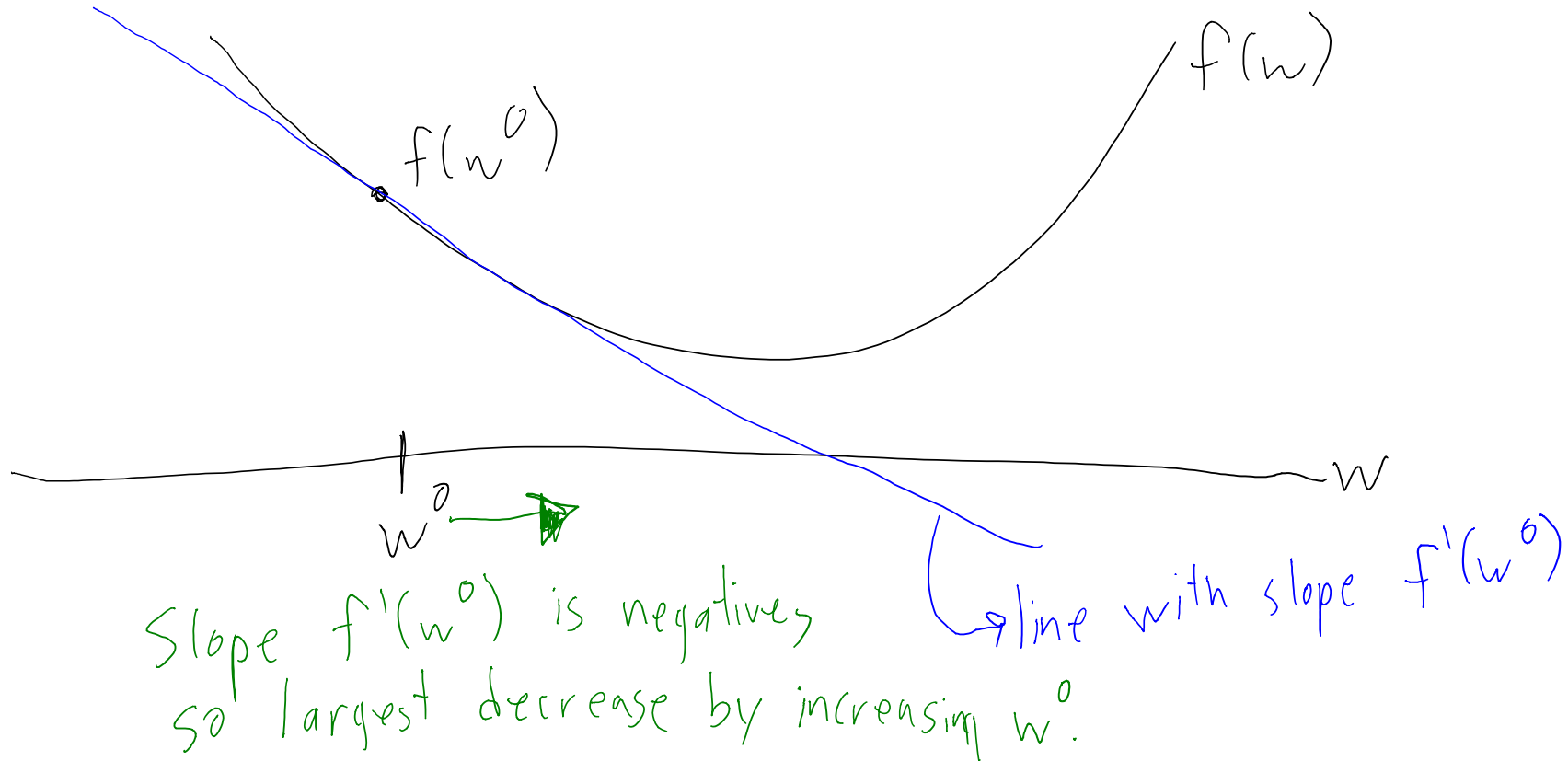
- Recall the **L2-regularized logistic regression** objective function:

$$\operatorname{arg\,min}_{w \in \mathbb{R}^d} \sum_{i=1}^d \log(1 + \exp(-y_i w^\top x_i)) + \frac{\lambda}{2} \|w\|^2$$

- This objective function is **strictly-convex and differentiable**.
- But we **can't formulate as linear system or linear program**.
- Nevertheless, we can efficiently solve this problem.
- There are many ways to do this, but we focus on **gradient descent**:
 - **Iteration cost is linear in 'd'** (not true of IRLS/Newton's method).
 - We can prove that **we don't need too many iterations**:
 - Number of iterations does not directly depend on 'd'.

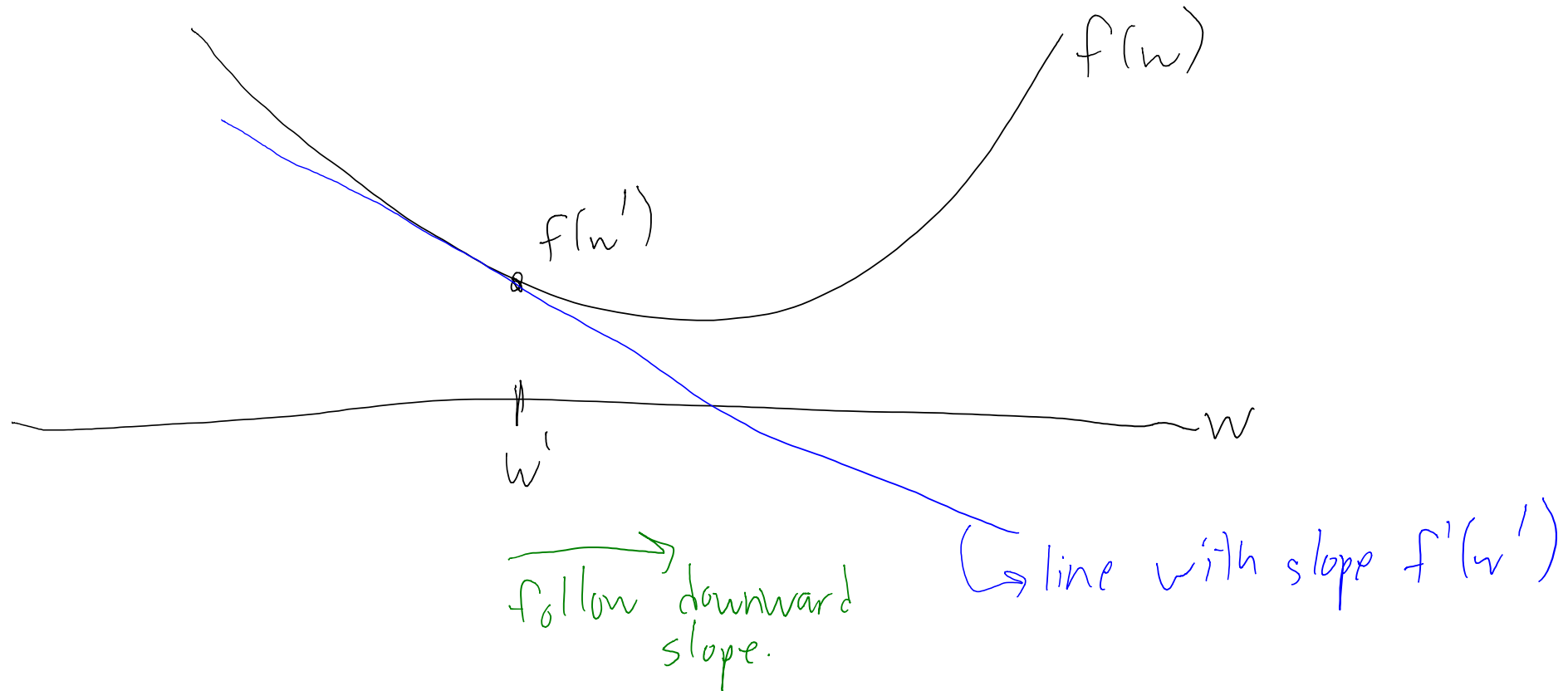
Gradient Descent

- Gradient descent is based on a simple observation:
 - Given parameters ' w^0 ', direction of largest decrease is $-\nabla f(w^0)$.



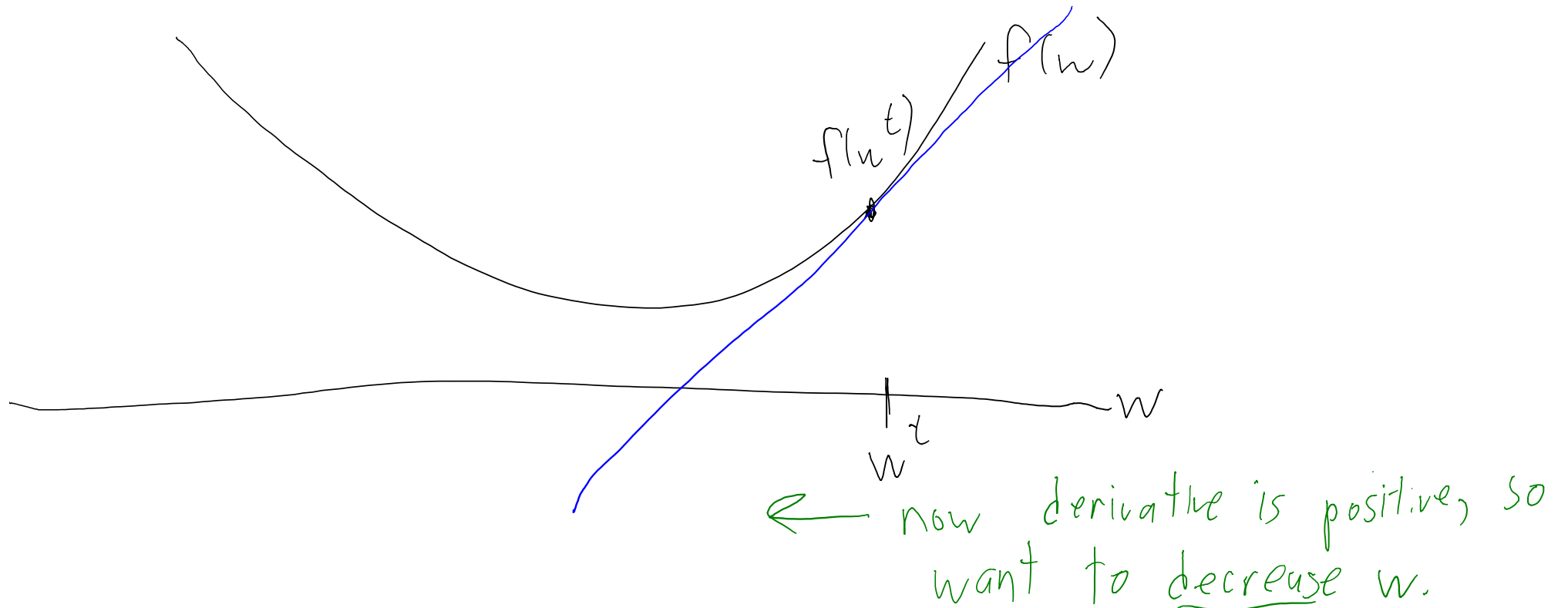
Gradient Descent

- Gradient descent is based on a simple observation:
 - Given parameters w^0 , direction of largest decrease is $-\nabla f(w^0)$.



Gradient Descent

- **Gradient descent** is based on a simple observation:
 - Given parameters ' w^0 ', **direction of largest decrease is $-\nabla f(w^0)$** .



Gradient Descent

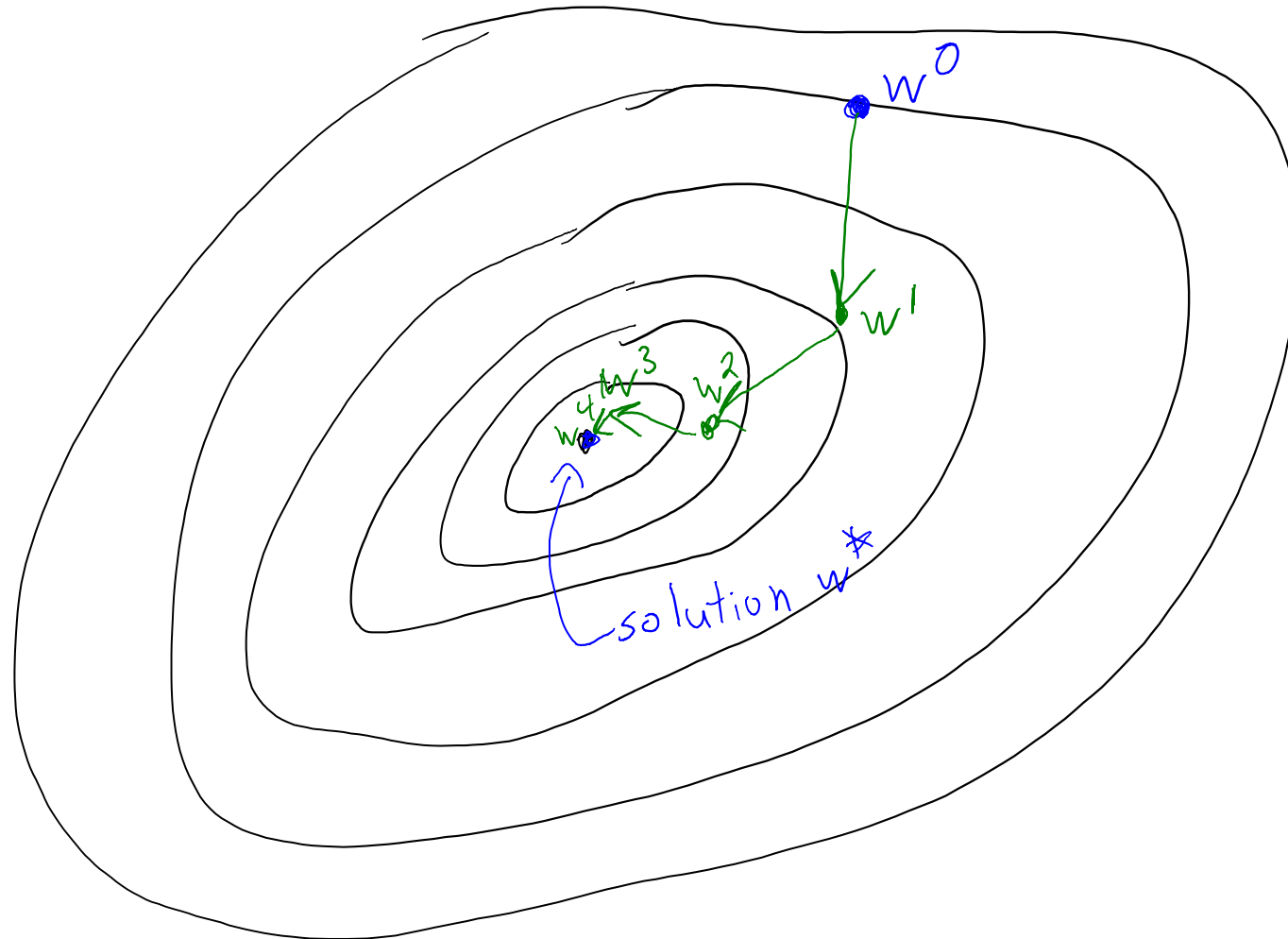
- Gradient descent is an iterative algorithm:
 - We start with some initial guess, w^0 .
 - Generate new guess by moving in the negative gradient direction:

$$w^1 = w^0 - \alpha_0 \nabla f(w^0).$$

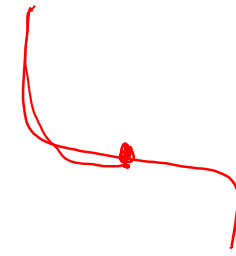
(The scalar α_0 is the 'step size'.)

- Repeat to successively refine the guess: $w^{t+1} = w^t - \alpha_t \nabla f(w^t)$.
- Generate w^2, w^3, w^4, \dots
- Stop if not making progress or $\|\nabla f(w^t)\| \leq \delta$ (some small number)

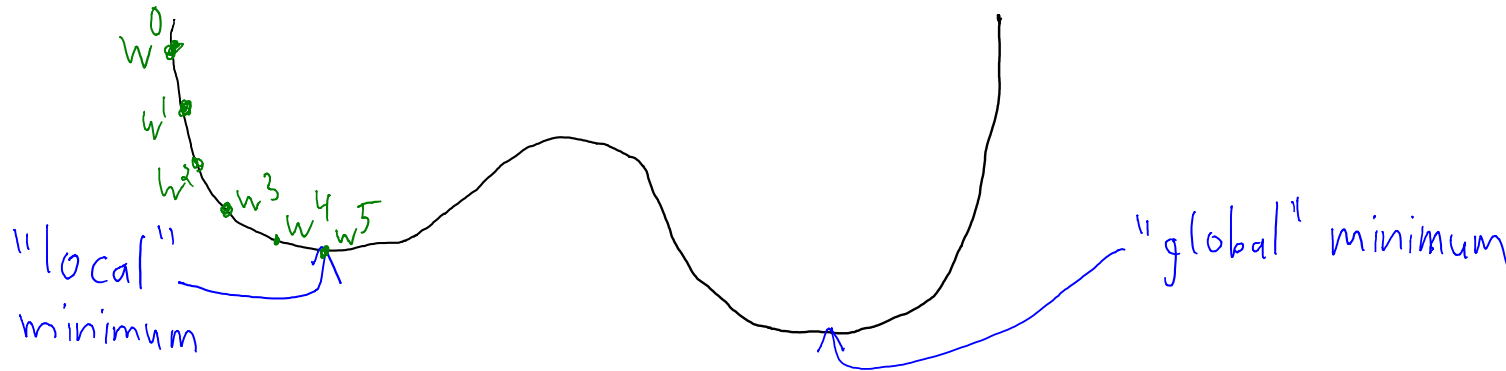
Gradient Descent in 2D



Gradient Descent



- If α_t is small enough and $\nabla f(w^t) \neq 0$, guaranteed to decrease 'f':
$$f(w^{t+1}) < f(w^t)$$
- Under weak conditions, procedure converges to a stationary point.



If 'f' is convex, converges to global minimum.

- **Least squares via linear system vs. gradient descent:**
 - Solving linear system cost $O(nd^2 + d^3)$.
 - Gradient descent costs $O(ndt)$ to run for 't' iterations.
 - Will be faster if $t < d$.

$$X^T(Xw) + X^T y$$

$O(nd)$

Convergence Rate of Gradient Descent

- How many iterations do we need?

- Let x^* be the optimal solution and ϵ be the accuracy we want.
- What is the smallest number of iterations 't' such that:

$$f(x^t) - f(x^*) \leq \epsilon$$

Notation:

- In optimization, we usually talk about optimizing x.

- To answer this question, **need assumptions**:

Let's assume

$$\underbrace{\mu I \preceq \nabla^2 f(x)}_{\text{"strongly convex"}} \preceq \underbrace{L I}_{\text{"strongly smooth"}} \quad \text{for all } x \text{ and some } L < \infty \text{ and } \mu > 0.$$

Strongly-convex
 \Rightarrow strictly-convex
 \Rightarrow convex.

"strongly convex"

$$y^T \nabla^2 f(x) y \geq \mu \|y\|^2$$

"strongly smooth"

$$y^T \nabla^2 f(x) y \leq L \|y\|^2$$

Example:

If $f(w) = \frac{1}{2} \|Xw - y\|^2 + \frac{\lambda}{2} \|w\|^2$, then $\mu = \min \text{eig}(X^T X) + \lambda \geq \lambda$ and $L = \max \text{eig}(X^T X) + \lambda$.

Bonus Slide: Constants for Least Squares

- Consider **least squares**: $f(x) = \frac{1}{2} \|Ax - b\|^2$

What are 'L' and 'u' such that $\mu I \preceq \nabla^2 f(x) \preceq L I$?

Note that $\nabla^2 f(x) = A^T A$, and since it's symmetric we can ^{use} Spectral decomposition:

$$A^T A = \sum_{j=1}^d \lambda_j q_j q_j^T \text{ where } q_j^T q_j = 1 \text{ and } q_i^T q_j = 0 \text{ for } i \neq j. \text{ (Assume } \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d)$$

We can write any y as linear combination of orthogonal basis, $y = \alpha_1 q_1 + \alpha_2 q_2 + \dots + \alpha_d q_d$.

$$\text{So we have } y^T \nabla^2 f(x) y = y^T A^T A y = y^T \left(\sum_{j=1}^d \lambda_j q_j q_j^T \right) y = \sum_{j=1}^d \lambda_j \underbrace{y^T q_j}_{=\alpha_j} \underbrace{q_j^T y}_{=\alpha_j} = \sum_{j=1}^d \lambda_j \alpha_j^2$$

Note that we can assume $\|y\|=1$

$$\text{or } y^T y = \sum_{j=1}^d \alpha_j^2 = 1.$$

So $y^T \nabla^2 f(x) y$ is maximized when $\alpha_1^2 = 1$ and minimized when $\alpha_d^2 = 1$, giving $L = \lambda_1 = \max(\text{eig}(A^T A))$ and $\mu = \lambda_n = \min(\text{eig}(A^T A))$

Convergence Rate of Gradient Descent

- The gradient descent iteration:

$$x^{t+1} = x^t - \alpha_t \nabla f(x^t)$$

- Assumptions:

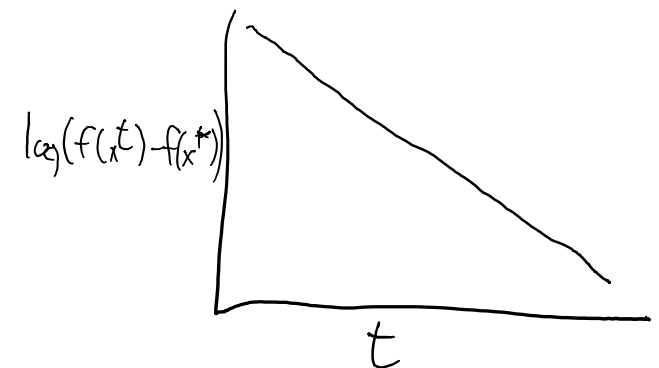
- Function 'f' is **L-strongly smooth** and **μ -strongly convex**.
- We set the step-size to $\alpha_t = 1/L$.

- Then gradient descent has a **linear convergence rate**:

$$f(x^t) - f(x^*) \leq O(\rho^t) \text{ for } \rho < 1.$$

- It follows that **we need $t = O(\log(1/\epsilon))$ iterations**.
 - This is good! We **want 't' to grow slowly** in accuracy $1/\epsilon$.
- Also called 'exponential' convergence rate.

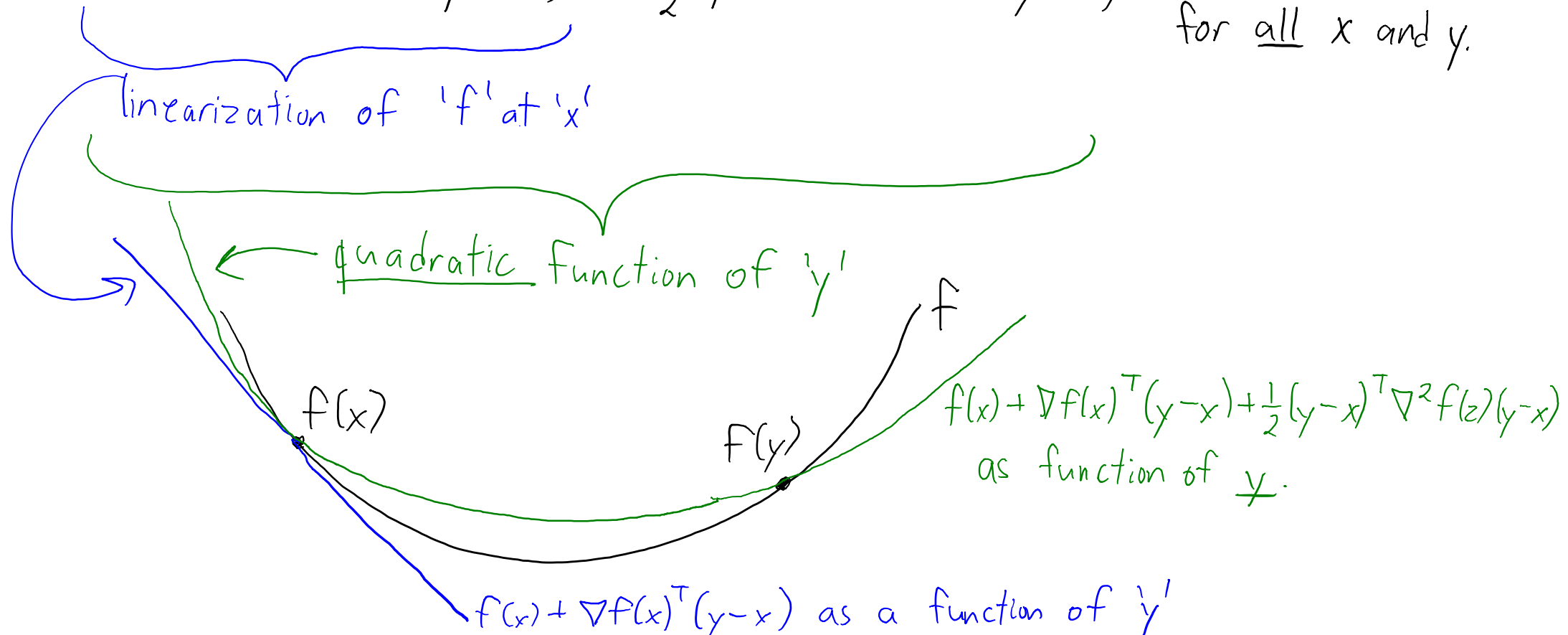
$$\begin{aligned} f(x^t) - f(x^*) = \epsilon &\leq O(\rho^t) \\ \text{means } \epsilon &\leq c \rho^t \text{ for } t \text{ large} \\ \text{or } \log(\epsilon) &\leq \log(c \rho^t) \\ &= \log(c) + t \log(\rho) \\ \text{or } t &\geq \frac{\log(\epsilon)}{\log(\rho)} - \text{const} \\ \text{or } t &\geq O(\log(1/\epsilon)) \\ &\text{(since } \rho < 1) \end{aligned}$$



Convergence Rate of Gradient Descent

- One version of **Taylor expansion**:

$$f(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(z) (y-x) \quad \text{for some } z \text{ for all } x \text{ and } y.$$



Using Strong-Smoothness

- One version of **Taylor expansion**:

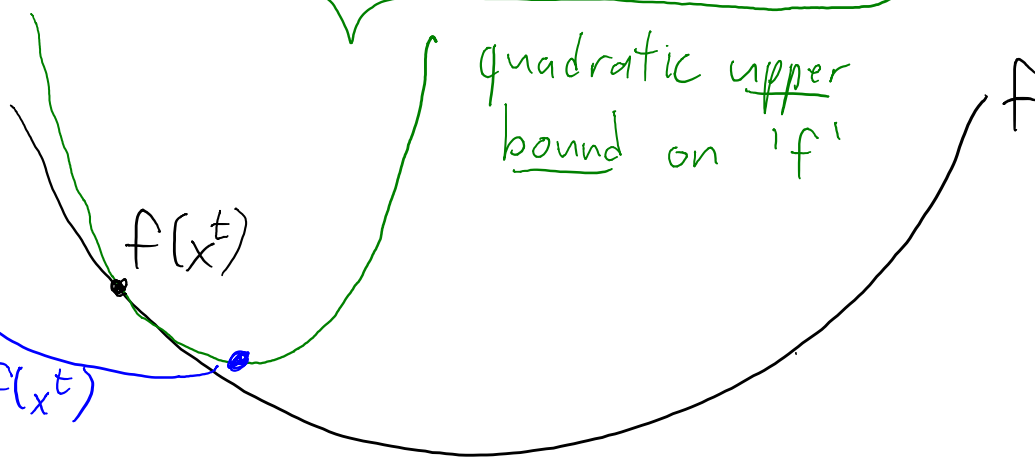
$$f(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(z) (y-x) \quad \text{for some } z$$

for all x and y .

From strong-smoothness we have: $v^T \nabla^2 f(z) v \leq L \|v\|^2$ for any z and v .

$$f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|y-x\|^2 \quad \text{for all } x \text{ and } y.$$

If we set x^{t+1}
to minimize bound
we get $x^{t+1} = x^t - \frac{1}{L} \nabla f(x^t)$



Let's find min of quadratic upper bound:
Let $q(y) = f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|y-x\|^2$
 $\nabla q(y) = 0 + \nabla f(x) - 0 + L(y-x)$
or
 $y = x - \frac{1}{L} \nabla f(x).$

Using Strong-Smoothness

- One version of **Taylor expansion**:

$$f(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(z) (y-x) \quad \text{for some } z$$

for all x and y .

From strong-smoothness we have: $v^T \nabla^2 f(z) v \leq L \|v\|^2$ for any z and v .

$$f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|y-x\|^2 \quad \text{for all } x \text{ and } y.$$

Set $x = x^t$ and $y = x^{t+1}$:

$$\begin{aligned} f(x^{t+1}) &\leq f(x^t) + \nabla f(x^t)^T (x^{t+1} - x^t) + \frac{L}{2} \|x^{t+1} - x^t\|^2 \\ &= f(x^t) + \nabla f(x^t)^T \left(-\frac{1}{L} \nabla f(x^t)\right) + \frac{L}{2} \left\|-\frac{1}{L} \nabla f(x^t)\right\|^2 \\ &= f(x^t) - \frac{1}{L} \nabla f(x^t)^T \nabla f(x^t) + \frac{1}{2L} \|\nabla f(x^t)\|^2 \\ &= f(x^t) - \frac{1}{2L} \|\nabla f(x^t)\|^2 \end{aligned}$$

$$x^{t+1} = x^t - \frac{1}{L} \nabla f(x^t)$$

(minimum of upper bound)

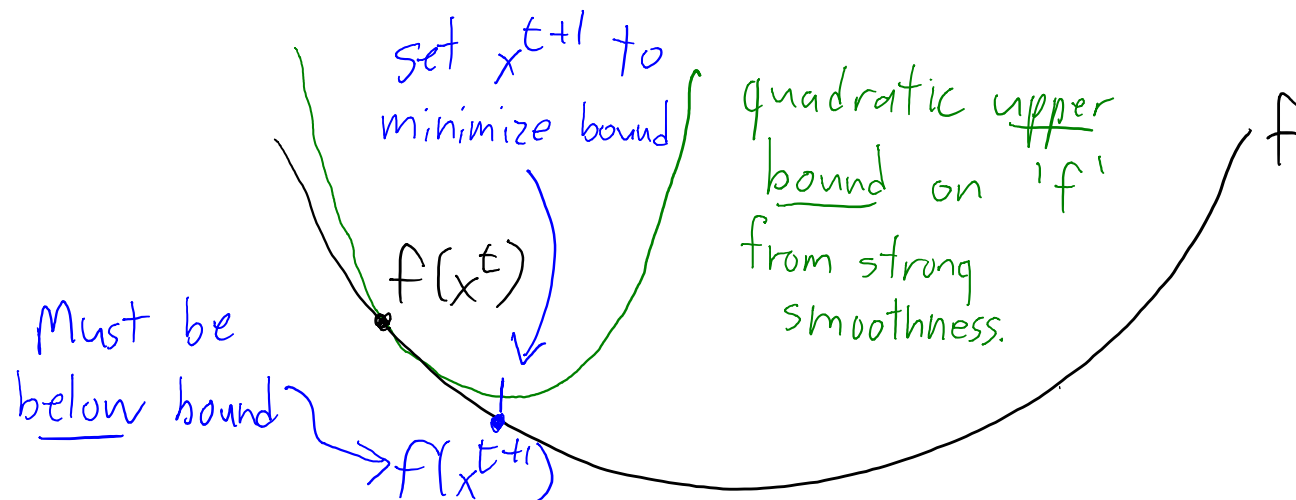
$$\nabla f(x)^T \nabla f(x) = \|\nabla f(x)\|^2$$

Using Strong-Smoothness

- We've derived a **bound on guaranteed progress** at iteration 't':

$$f(x^{t+1}) \leq f(x^t) - \frac{1}{2L} \|\nabla f(x^t)\|^2$$

- If gradient is non-zero, **guaranteed to decrease objective**.
- Amount we decrease grows with the size of the gradient.
- Note: bound involves for any strongly-smooth function (e.g., non-convex)



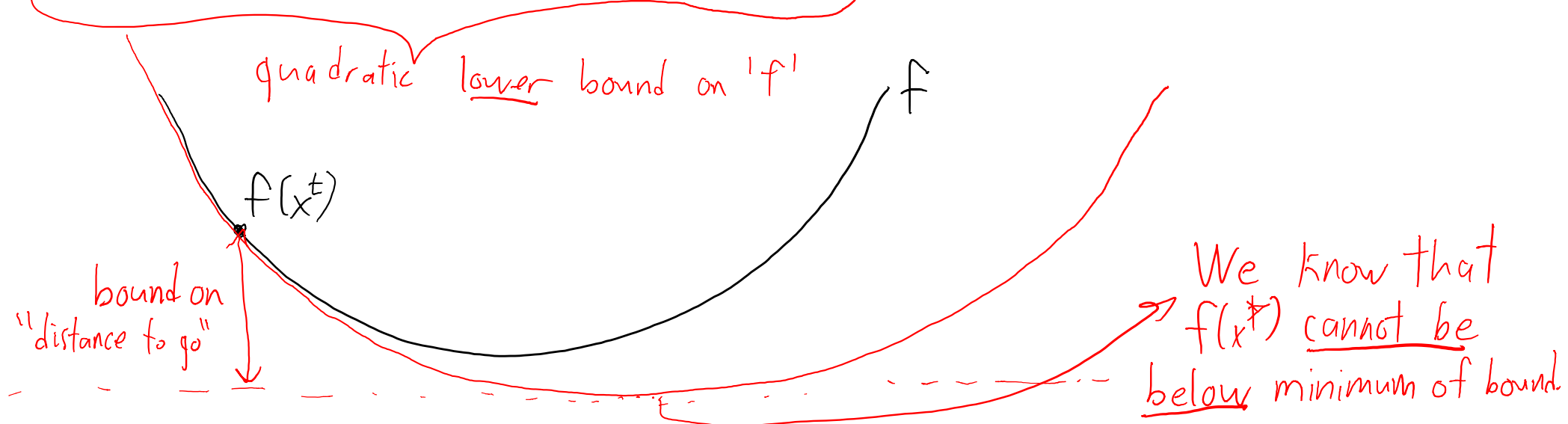
Using Strong-Convexity

- One version of **Taylor expansion**:

$$f(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(z) (y-x) \quad \text{for some } z \text{ for all } x \text{ and } y.$$

By strong-convexity we have $v^T \nabla^2 f(z) v \geq \mu$ for all v and z .

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|^2$$



Using Strong-Convexity

- One version of **Taylor expansion**:

$$f(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(z) (y-x) \quad \text{for some } z$$

for all x and y .

By strong-convexity we have $v^T \nabla^2 f(z) v \geq \mu$ for all v and z .

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|^2$$

Minimize both sides with respect to y :

$$f(x^*) \geq f(x) - \frac{1}{2\mu} \|\nabla f(x)\|^2$$

Combining Strong-Smoothness and Convexity

- Our bound on **guaranteed progress**:

$$f(x^{t+1}) \leq f(x^t) - \frac{1}{2L} \|\nabla f(x^t)\|^2$$

- Our bound on 'distance to go':

$$f(x^*) \geq f(x^t) - \frac{1}{2\mu} \|\nabla f(x^t)\|^2 \Leftrightarrow -\frac{1}{2} \|\nabla f(x^t)\|^2 \leq -\mu(f(x^t) - f(x^*))$$

- Use 'distance to go' bound in guaranteed progress bound:

$$f(x^{t+1}) \leq f(x^t) - \frac{1}{L} (-\mu(f(x^t) - f(x^*)))$$

- Subtract $f(x^*)$ from both sides and simplify:

$$\begin{aligned} f(x^{t+1}) - f(x^*) &\leq f(x^t) - f(x^*) - \frac{\mu}{L} (f(x^t) - f(x^*)) \\ &= \left(1 - \frac{\mu}{L}\right) [f(x^t) - f(x^*)] \end{aligned}$$

Combining Strong-Smoothness and Convexity

- We've shown that:

$$f(x^t) - f(x^*) \leq \left(1 - \frac{\mu}{L}\right) [f(x^{t-1}) - f(x^*)]$$

- Applying this recursively:

$$\begin{aligned} f(x^t) - f(x^*) &\leq \left(1 - \frac{\mu}{L}\right) \left[\left(1 - \frac{\mu}{L}\right) [f(x^{t-2}) - f(x^*)] \right] \\ &= \left(1 - \frac{\mu}{L}\right)^2 [f(x^{t-1}) - f(x^*)] \\ &= \left(1 - \frac{\mu}{L}\right)^3 [f(x^{t-2}) - f(x^*)] \\ &\vdots \\ &= \left(1 - \frac{\mu}{L}\right)^t [f(x^0) - f(x^*)] \end{aligned}$$

$\rightarrow O(p^t)$

- Since $\mu \leq L$, we've shown **linear convergence rate**.

Discussion of Linear Convergence Rate

- We've shown that gradient descent under certain settings has:

$$f(x^t) - f(x^*) \leq \left(1 - \frac{\mu}{L}\right)^t [f(x^0) - f(x^*)]$$

- The number L/μ is called the 'condition number' of 'f'.
- Connection to matrix condition number:
 - For least squares, condition number of 'f' is condition number of $X^T X$.
- This rate is **dimension-independent**:
 - It does not directly depend on dimensions 'd'.
 - In principle, applies to infinite-dimensional problems.
 - But, L and μ may be larger in high-dimensional spaces.
- In practice, typically **you don't have 'L'**.
 - We'll discuss practical issues next time.

Summary

- **No free lunch**: there is no ‘best’ machine learning model.
- **Softmax loss** to model discrete y_i , other losses can be derived.
- **Convex functions**: all stationary points are global minima.
- **Show functions are convex**.
- **Gradient descent** finds stationary point of differentiable function.
- **Rate of convergence** of gradient descent is linear.

- **Next time**:
 - What if we don’t know which features are relevant or which basis to use?