Application/data
Classified data?
  Yes: Categorical labels (classification or quantitative (regression)
  no: Similarity (clustering) or partitioning (dimensionality reduction)

4 levels of abstraction

Model/method                                              CS189
e.g.
    - decision function: linear, polynomials, logistic, neural net...

    - nearest neighbors, decision trees

    - features: edge detectors, embeddings, raw pixels

    - low vs. high capacity model: how sinuous/complicated?
         ↳ affects overfitting/underfitting, inference

Optimization Problem
- variables/parameters, objective function, constraints
    e.g. unconstrained, convex programs, least squares linear regression, PCA

Optimization Algorithms
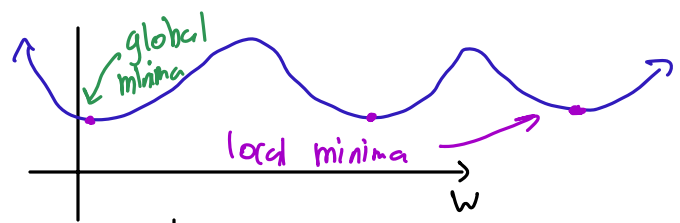e.g. gradient descent, SGD, simplex, SVD (for PCA)

Optimization Problems

Unconstrained: given continuous objective function f, find $\vec{w}$ that minimizes/maximizes f
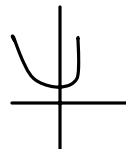
    - f is smooth if $\frac{df}{dx}$ is also continuous

    - global minimum of f: value $\vec{w}$ s.t. $f(\vec{w}) \leq f(\vec{v}) \; \forall \vec{v}$
         ↳ Sometimes need to settle for local minimum: $f(\vec{w}) \leq f(\vec{v}_i)$,
            $\vec{v}_i \in$ "tiny ball centered @ $\vec{w}$



global minima

local minima

W

Usually, finding local minima is easy, but
finding global minima is hard or impossible
    ↳ exception: convex functions: $\forall x, y \in \mathbb{R}^d$, line connecting $(x, f(x)) \& (y, f(x))$
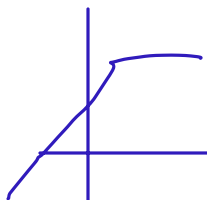       doesn't go below $f(\cdot)$

e.g. Perceptron risk function is convex & nonsmooth

       ↳ convex since it's a sum of convex parts

Consider a continuous, convex function over a closed domain:

  3 possibilities:

    1) no minima (goes to $-\infty$)

    2) just one local minimum

Perceptron { 3) connected set of local minima
risk function {
                   (that are all global minima
                   w/ same objective values)

     ↳ risk = 0 in pie slice

} Walk downhill

## Algos for smooth functions:

- gradient descent

    $w = w - \epsilon \nabla f(w)$

    - blind gradient descent: simply steps down, doesn't know how far to go down

    - SGD (also blind)

    - line search: looks ahead & tries to find minimum along the direction

- Newton's Method (needs Hessian matrix of f)

    ↳ can be expensive, e.g. for NN's

- Nonlinear Conjugate gradient method

    - has line search as subroutine

these algos find local minimum, not global

## Algos for nonsmooth f:

- can still do gradient descent

Won't work: Newton's Method (needs 2nd order information)

- BFGS

**line search:** repeat
1) pick direction
2) find local minimum along that dimension/line by solving an optimization problem in 1D

**Some examples:**
1) secant method (smooth only)

2) 2nd derivatives (must be smooth)

3) direct line search
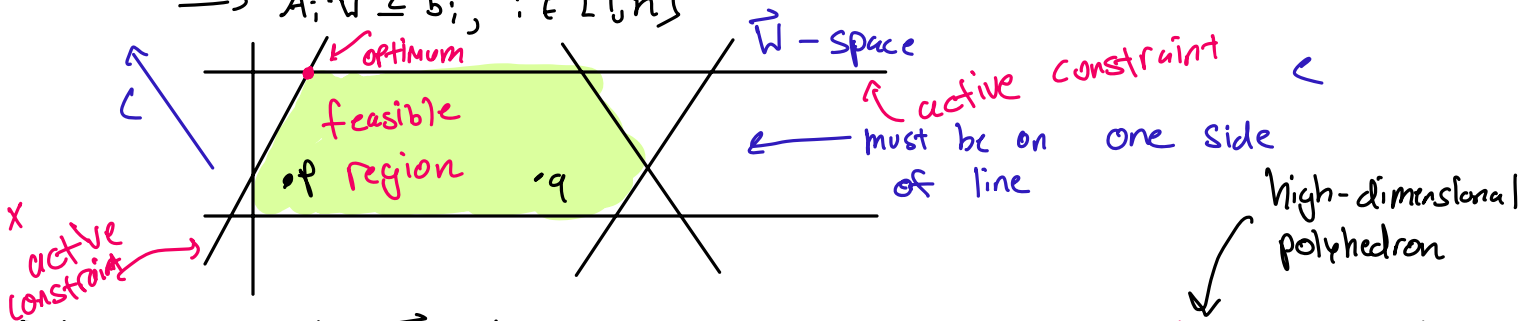   ↳ golden section search (nonsmooth functions)

**Constrained Optimization** (smooth equality constraints)

goal: find $\vec{u}$ that minimizes $f(\vec{u})$ s.t. $g(\vec{u}) = 0$, $g$ is smooth

:isosurface
↳ find point on surface s.t. $g(u) = 0$

algorithm: lagrange multipliers ( constrained smooth optimization problem → Smooth <u>unconstrained</u> optimization problem

**Linear Program:** linear objective func. & linear inequality constraints

goal: find $\vec{u}$ that maximizes $c \cdot u$ s.t. $A\vec{u} \leq b$

$A \in \mathbb{R}^{n,d}$, $b \in \mathbb{R}^n$ ⟹ $n$ linear constraints ← component-wise constraints

⟹ $A_i \cdot u \leq b_i$, $i \in [1, n]$



$\vec{W}$ - space

optimum

feasible region

active constraint
← must be on one side of line

high-dimensional polyhedron

x active constraint

The set of points $\vec{u}$ that satisfy all constraints is a **polytope** called the **feasible region F**

The **optimum** is point in F that maximizes $c \cdot W$ (furthest in direction of c)

A point set P is **convex** if for every $p, q \in P$, line connecting $p$ & $q$ is entirely in the point set

The optimum generally achieves equality for some constraints ( but not most of them)
↳ **active constraints** of the optimum

example: <u>every</u> feasible point (W, d) gives a linear classifier)

            ↳ not necessarily the best in test time

              ↳ 100% training accuracy

→ find $W, d$ that maximizes $0$    s.t. $y_i(x_i \cdot v + d) \geq 1$   $\forall i \in [1...n]$

IMPORTANT: linearly separable data iff feasible region $\neq$ empty set

      ⟶ also true for maximum margin classifier (quadratic program)

<u>Algos for solving linear programs:</u>

    – Simplex ( George Dantzig, 1947)

        – Walks from vertex to vertex in polytope

    – Interior Point methods

Quadratic Program: quadratic, convex objective function

      ⟹ hessian is PSD

goal: find $\vec{W}$ that minimizes $f(v) = W^T Q v + c^T v$   ←— quadratic term

    s.t. same linear constraints as linear program    ←— Symmetric Q, Positive definite

      ↳ only one local minimum (global minimum)

example: find maximum margin classifier