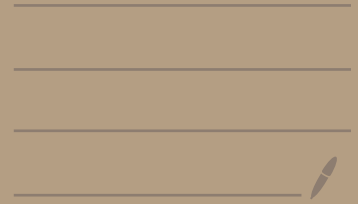
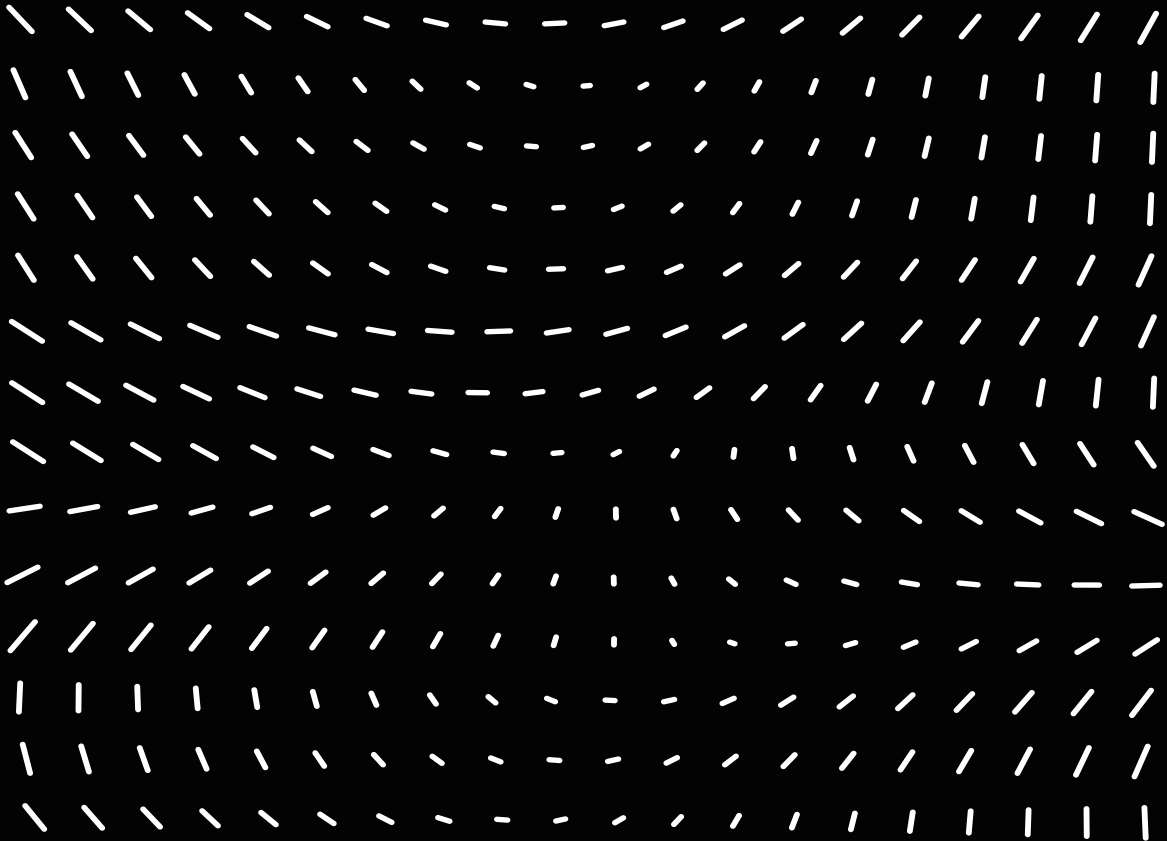


CISC 371



All Lectures



Class #0: Elementary Differential Calculus

Functions Overview

$$f: \mathbb{R} \rightarrow \mathbb{R}$$

A function f takes real number as input and gives real number as output

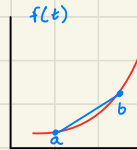
Limit

$$\lim_{x \rightarrow 0} f(x) = C$$

As x approaches 0, the value of the function $f(x)$ approaches C .

Continuity

- A function is continuous at a if $\lim_{t \rightarrow a} f(t) = f(a)$
- NO breaks, jump discontinuity
- Chord is a line segment for a function



Derivatives

Power rule: $t^x = x t^{x-1}$ $2^5 = 5 \cdot 2^4$

trig basics: $\sin(x) = \cos(x)$, $\cos(t) = -\sin(t)$

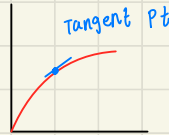
$$e^x = e^x, \ln(x) = \frac{1}{x}$$

constant rule: $c \cdot f(x) = c \left(\frac{d}{dt} f(x) \right)$

Power rule: $f(x)' \cdot g(x) + g(x)' \cdot f(x)$

Composition is $f(g(x))$ or $(f \circ g)(x)$

Composite function implies plugging one into another



Ex 1

$$f(x) = x^2$$

$$g(x) = t + 1$$

$$\begin{aligned} f(g(t)) &= f(t+1) \\ &= (t+1)^2 \\ &= t^2 + 2t + 2 \end{aligned}$$

Chain Rule:

$$(f \circ g)' = (f' \circ g) \cdot g'$$

Differentiate the outside, keep the inside, multiply by inner deriv

$$\begin{aligned} \text{Ex } f(g(t)) &= (t+1)^2 \\ &= 2(t+1)(1) \\ &= 2t + 2 \end{aligned}$$

Reciprocal Rule

$$h(t) = \frac{1}{g(t)}$$

Ex $h(t) = \frac{1}{t}, h'(t) = -\frac{1}{t^2}$

$$h'(t) = \frac{g'(t)}{(g(t))^2}$$

Quotient Rule

$$h(t) = \frac{f(t)}{g(t)}$$

$$h'(t) = \frac{g(t) f'(t) - f(t) g'(t)}{(g(t))^2}$$

Ex

$$h(t) = \frac{t^2}{t+1} \rightarrow h'(t) = \frac{(t+1)(2t) - (t^2)(1)}{(t+1)^2}$$

Taylor Series

- Taylor series is a way to approximate a complicated function using polynomials
- Expands around a point to

Ex $f(t) = f(t_0) + f'(t_0)(t-t_0) + \frac{f''(t_0)}{2!}(t-t_0)^2 + \dots$

- Analytic means a function is infinitely differentiable
- Zooming in on curve

Linear Approximation

$$f(t) \approx f(t_0) + f'(t_0)(t-t_0)$$

Quadratic Approximation

$$f(t) \approx f(t_0) + f'(t_0)(t-t_0) + \frac{f''(t_0)}{2}(t-t_0)^2$$

Lagrange Remainder

- No approximation is perfect, Lagrange remainder tells us how much error there is after we stop after k terms.
- gives worst case error

Class #01: Introduction to Optimization - PDF Notes

- Optimization is choosing the best element from a set according to a rule

EX $f(t) = \frac{-t}{t^2+1}$

This has max at $t = -1$, min at $t = 1$

- Minimum is the lowest value the function takes
- Minimizer is the input that gets you the minimum

Fermat's Problem

- Find the point in the plane that minimizes the sum of distances to all pts
 $f(w) = \sum_{j=1}^n \|w - a_j\|$
- For 3 anchors:
 - if all angles are $< 120^\circ$, minimizer lies inside triangle
 - if an angle is $\geq 120^\circ$, minimizer is that vertex
- For 5 anchors, there's no closed form solution

3 Types of Optimization Problems

- 1) Unconstrained, scalar (easy, one var)
- 2) Unconstrained, vector (harder, need direction + step sizes)
- 3) Constrained, vector (hardest, reduce to type 2)

Important Definitions

- **Open Set**, for every vector \vec{v} in an open space if a vector is near it, only **interior points**
- **Interior Point**, a point bounded by something, with some wiggle room.
- **Boundary Point**, its in the subset, but its not an interior point. Touching the edge

More Definitions

- **Global minimizer**: lowest overall, allows equality, \vec{w}^* of $f(\vec{w})$
- **Strict Global minimizer**: absolute lowest point, $f(t^*) < f(w)$, $\forall w \neq t^*$
- **Global minimum**: L is lowest level of some function f
- **Local minimizer**: lowest in neighborhood but allows equality

Class #02: Minimizing By Approximation - PDF Notes

Approximation

- Some objective functions are too complicated to minimize
 - ↳ Idea: Replace a hard function $f(t)$ with a simpler polynomial model
- Build quadratic approximation
- Solve minimizer
- Get better guess

amplitude = size of oscillations

3 Quadratic Interpolation Methods

1) 3 points (NO Derivatives)

- take 3 nearby points $(t_1, f_1), (t_2, f_2), (t_3, f_3)$

$$p(t) = a_1 t^2 + a_2 t + a_3$$

- Solve system Vandermonde matrix to find coefficients

• Minimizer: $\hat{t} = \frac{a_2}{2a_1}$

2) 2 points + First Derivative

- if you have derivative at one point
- Constraints $p(t_1) = f_1, p(t_2) = f_2, p'(t_1) = f_1'$
- Solve for quadratic coefficients

3) 1 point + First and Second Derivative

- $f(t_1)$
- $f'(t_1)$ enough to know quadratic
- $f''(t_1)$

Always check if new guess is improvement from old guess.

First Derivative \rightarrow descent

Second Derivative \rightarrow positive

Higher order Approximations

- Use quadratics (sweet spots)
- ↳ could use cubic models

Class 2: In person / Lecture video Notes

Orders of Approximation

$$p_0(t) = c \quad \text{Flat line}$$

$$p_1(t) = mt + b \quad \text{Linear}$$

$$p_2(t) = a_1 t^2 + a_2 t + a_3 \quad \text{Quadratic, can actually have minimum}$$

Vandermode Matrix

- Quadratic Interpolation using Vandermode matrices
- approximate using a quadratic model

$$p(t) = a_1 t^2 + a_2 t + a_3$$

To determine a_1, a_2, a_3 (coefficients) we need constraints (data points / derivatives)

Case 1) 3 Points

Plug in $(t_1, f_1), (t_2, f_2), (t_3, f_3)$

$$\begin{bmatrix} t_1^2 & t_1 & 1 \\ t_2^2 & t_2 & 1 \\ t_3^2 & t_3 & 1 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = \begin{bmatrix} f_1 \\ f_2 \\ f_3 \end{bmatrix}$$

Minimizer: $t = -\frac{a_2}{2a_1}$ (need coeffs)

Think like

This is Vandermode Matrix

$$\begin{bmatrix} x_1^2 & x_1 & 1 \\ x_2^2 & x_2 & 1 \\ x_3^2 & x_3 & 1 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}$$

Case 2: 2 points + 1 Derivative

Points: $(t_1, f_1), (t_2, f_2)$ but also

Slope at one point: $f'(t_1)$

$$\begin{bmatrix} t_1^2 & t_1 & 1 \\ t_2^2 & t_2 & 1 \\ 2t_1 & 1 & 0 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = \begin{bmatrix} f_1 \\ f_2 \\ f'_1 \end{bmatrix}$$

Case 3: 1 point, First and Second Deriv

$$\begin{matrix} f(t_1) \\ f'(t_1) \\ f''(t_1) \end{matrix} \begin{bmatrix} t_1^2 & t_1 & 1 \\ 2t_1 & 1 & 0 \\ 2 & 0 & 0 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = \begin{bmatrix} f_1 \\ f'_1 \\ f''_1 \end{bmatrix}$$

Once coefficients are found, can analytically solve minimizer $p'(t) = 0$

Class #03: Stationarity and Convexity - PDF Notes

Stationary Points

t^* could be a min/max or saddle point

• Stationary point is when $f'(t^*) = 0$

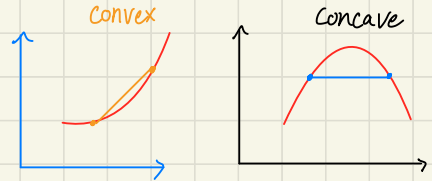
↳ function is flat

If t^* is local minimizer, it must be Stationary point

First order necessary condition?

Stationary \neq minimizer always

could be local min
local max
Saddle points



Convexity

• Convex means it lies below the chord at any 2 points

↳ Strictly convex:

↳ local min = global min
guarantees uniqueness

Every local minimum is a global minimum

Convex - function lies below or on chord

Strictly convex - strictly below chord between any 2 points

↳ guarantees unique minimizer

Gradient Inequality

• From Taylor Expansion with remainder
 $f(t) = f(t_0) + f'(t_0)(t - t_0) + \frac{1}{2} f''(\xi)(t - t_0)^2$

Meaning: When you approximate a function around point t_0 . There is a point ξ between t and t_0 , where second deriv is evaluated

If f is convex $f''(\xi) \geq 0$ then:

$$f(t) \geq f(t_0) + f'(t_0)(t - t_0)$$

This says the tangent line is below the curve

For convex functions, the tangent line is always global underestimator

• Convex helps find mins

4. Key Concepts to Distinguish

- Stationary point: derivative = 0 (necessary, not sufficient for min).
- Saddle point: stationary but neither max nor min.
- Convex function: guarantees well-behaved minima.
- Gradient inequality: shows why convexity makes optimization stable.

$$\underbrace{\frac{1}{2} f''(\xi)(t - t_0)^2}_{\text{Error Term}}$$

Class 04: Scalar Minimization - PDF Notes

- Vector Valued Problems are fine to have good approximations
- **Inexact methods**: aim to reduce computational effort while making progress

Line Search

• moving along given direction d from current estimate t_k

↳ Objective function: $f(t)$

Current guess t_k , function value f_k , derivative f'_k

direction $d = \pm 1$ (scalar)

Fixed Stepsize Search

Simplest approach: $t_{k+1} = t_k - s_0 f'(t_k)$

s_0 - constant stepsize

Too Big \rightarrow overshoot, divergence

Too Small \rightarrow slow convergence

Inexact Backtracking Search

• Start with user's stepsize s_0 , shrink if it is

too large. Keep multiplying by β .

• pick largest s improves $f(t)$

• want smaller stepsizes

Armijo Backtracking Condition

• Based on gradient inequality

Step s such that:

$$f(t_k + sd) \leq f(t_k) + \alpha_k sd$$

where $\alpha_k = f'(t_k)/2$

Reduce s by a factor β

New point must be strictly lower than predicted by tangent line.

✓ Key Takeaways

1. Fixed stepsize methods are simple but unstable.
2. Backtracking systematically reduces stepsize until descent is guaranteed.
3. Armijo condition formalizes the test: ensures sufficient decrease, not just any decrease.
4. This makes optimization efficient, robust, and less sensitive to poor stepsize choices.

Class 4 - Lecture Video Notes

Finding Stationary Point for Approximation

$$f'(t_{k+1}) = 0$$

Differentiate:

$$f'(t) \approx f'(t_k) + f''(t_k)(t - t_k)$$

$$0 = f'(t_k) + f''(t_k)(t_{k+1} - t_k)$$

$$t_{k+1} = t_k - \frac{f'(t_k)}{f''(t_k)}$$

Direction $d_k = -f'(t_k)$

Step size $s_k = \frac{1}{f''(t_k)}$

⇒
Iteration

$$t_{k+1} = t_k + s_k d_k$$

or

$$t_{k+1} = t_k - \alpha f'(t_k)$$

• Smaller Step Size is always better

$$f(t) = t^2$$
$$f'(t) = 2t$$
$$d_k = -2t_k$$

Backtracking/inexact line search is for back tracking

$$t_{k+1} = t_k - 2s t_k = (1-2s)t_k$$

Start with $t = 10, s = 0.25$

$$t_1 = (1-0.5) \cdot 10 = 5$$

$$t_2 = (1-0.5) \cdot 5 = 2.5$$

$$t_3 = 1.25$$

Converges to 0.

Backtracking

• Start with initial guess s

↳ factor of β : $s, \beta s, \beta^2 s, \beta^3 s$

Armijo Condition: new point should be

below a certain line to ensure progress

↳ Searching for a good step size

Exponential Backoff: $s, \beta s, \beta^2 s$

Class 05: Functions with a Vector Argument - PDF Notes

From Derivatives to Partial Derivatives

- multivariable calculus uses

Partial derivatives

$$\frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2}$$

How function changes as one coordinate changes

Directional Derivatives

- move in any direction \vec{v}

- Directional derivative: $D_{\vec{v}} f(\vec{w}) = \lim_{h \rightarrow 0} \frac{f(\vec{w} + h\vec{v}) - f(\vec{w})}{h}$

- Dot product between gradient and the direction vector.

Gradient (∇f)

- gradient is a row vector (1-form) for all partial derivatives

$$\nabla f(\vec{w}) = \left[\frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2}, \dots, \frac{\partial f}{\partial w_n} \right]$$

- gradient points in the direction of steepest ascent.

- negative direction is steepest descent.

Jacobian Matrix

- For vector valued function $\vec{f}(\vec{w}) = [f_1(\vec{w}), \dots, f_m(\vec{w})]^T$

- Jacobian is matrix of gradients

$$J_f(\vec{w}) = \begin{bmatrix} \nabla f_1 \\ \nabla f_2 \\ \nabla f_3 \end{bmatrix}$$

Level Curves

- level curve of $f(\vec{w})$ is the set of all points with same function value

$$SC(f, 1) = \{ \vec{w} : f(\vec{w}) = 1 \}$$

- gradient is always perpendicular to level curves.

Class 6 - Video Lecture Notes

Connecting multivariable Calc with Lin Alg

• Functions with multiple variables: $f(w_1, w_2, \dots, w_n)$

$$\hookrightarrow f(w_1, w_2) = w_1^2 + 3w_2$$

• Gradient is vector of partial derivatives: $\nabla f = \left[\frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2}, \dots \right]$

Example: $f(w_1, w_2) = w_1^2 + 3w_2$

$$\nabla f = [2w_1, 3], \text{ at point } 1, 2, \nabla f = [2, 3]$$

Directional Derivative:

rate of change of f in some direction v

$$D_v f = \nabla f \cdot v$$

At a point with certain gradient

$$D_v f = [2, 3] \cdot [1, 0] = 2$$

$\hookrightarrow f$ increases in w_1 direction by 2 steps per

1 Form

• linear functional that takes a vector and produces a number

vectors are directions

1-forms are row vectors that act on vectors

gradient is naturally a row vector (1 Form)

\hookrightarrow multiplies a column vector, gets directional derivatives

vector \rightarrow

one Form \rightarrow

tells change along vector

If a is vector $[\vec{a}]$, transpose is $[\vec{a}]^T = \frac{a}{[a]^T = \vec{a}}$

Jacobian Matrix

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

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

$$J = \begin{bmatrix} 2x & 1 \\ y & x \end{bmatrix}$$

Theorem: $a^T w$ (linear function)
gradient is a

Class 07: First order optimization - Steepest Descent - PDF Note

Optimization Setup

want: $w^* = \underset{w}{\text{arg min}} f(w)$

• Start at point w_0

At every step, we need a direction and a stepsize

• Updated rule

$$w_{k+1} = w_k + s d$$

next step starting step stepsize direction

Descent Directions

• d is descent direction if the directional derivative is negative

• Steepest descent is fastest decreasing f .

$$d = -\nabla f(w_0)$$

Stepsize Selection

• Fixed step size, too large \rightarrow overshoot, oscillate
 , too small \rightarrow crawl slowly

• Armijo Backtracking

• Start with large s_0 , shrink backwards

$$f(w_k + s d_k) \leq f(w_k) + \alpha s \nabla f(w_k) \cdot d_k \quad ?$$

Algorithms

• Input $w_0, s, K_{\max}, g_{\text{mag}}$

1) Find gradient: $g = \nabla f(w)$

2) Set $d = -g^T$

3) Update $w = w + s d$

4) Stop when $\|g\| < g_{\text{mag}}$ or max iterations

Examples

• Quadratic Function

$$f_1(w) = w^T K w, \quad K = \begin{bmatrix} 1/4 & 0 \\ 0 & 1 \end{bmatrix}$$

start at $w_0 = (-3.0, 3.2)$

Backtracking: ($s_0 = 1.5, \beta = 0.75$)

$$f_1(w) = \begin{bmatrix} -3.0 \\ 3.2 \end{bmatrix} \begin{bmatrix} 1/4 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -3.0 & 3.2 \end{bmatrix}$$

Class 09: Newton's Method - PDF Notes

Point of Newton's Method

- Steepest descent is cool, but often zigzags in deep valleys \rightarrow slow convergence
- Use curvature (Hessian, 2nd Derivative) to adjust step direction/length
- Use more info than just slope

Scaled Descent

- Some functions are sensitive to one variable more than others

$$f_1(w) = 100w_1^2 - 10w_1w_2 + w_2^2 + 100$$

\hookrightarrow changes in w_1 are 10x more impactful than w_2

Fix: Scale variables with matrix B

$$w = B M v, B \text{ reduces effect of } w_1$$

TIDR: Use inverse Hessian as B , for smoother and faster convergence

Newton's Method

- Use quadratic model

$$\text{Taylor Approx: } f(w) = f(w_0) + \nabla f(w_0)(w - w_0) + \frac{1}{2}(w - w_0)^T H(w - w_0)$$

$$d = -H^{-1} \nabla f(w_0)$$

$$w_{k+1} = w_k + d$$

Newton Direction:

Converges in 1 Step if model is exact?

Damped Newton's Method

- If steps are too big, Newton can overshoot
- Damping \rightarrow Scales back step size with armijo backtracking
??

Pitfalls

- Newton's method **only works if the Hessian is positive definite**
- Hessian is indefinite, method can diverge

Class 10: Nonlinear least squares and L-M Algorithm

Idea: Solve optimization problems, where the model being fitted are nonlinear. (similar to GPS)

Nonlinear Least Squares

• Trying to minimize: $f(\vec{w}) = \frac{1}{2} \|\vec{r}(\vec{w})\|^2$

\vec{w} = parameters you're optimizing

$r(\vec{w})$ = vector of residuals (error between predicted/observed)

Solving Gauss-Newton

• approximates as a locally linear problem

Iteration: $\vec{w}_{k+1} = \vec{w}_k + (J^T J)^{-1} J^T \vec{r}$

Jacobian Matrix (derivatives of residuals w.r.t parameters)

• scaled steepest descent method.

Levenberg-Marquardt Algorithm (LM)

• Smarter Gauss-Newton avoids problems when:

• your initial guess is bad

• Jacobian is near singular

• Add Damping term: $\vec{w}_{k+1} = \vec{w}_k + (J^T J + \pi I)^{-1} J^T \vec{r}$

• π is large, gradient descent

• π is small, acts like Gauss-Newton

Regularized the problem

hyperparam: $\pi \geq 0$

Non Linear Least Squares Solution

• Finding parameters x that fits observed data by minimizing sum of squared error

• Linear least squares: $f(x) = Ax - b$

↳ minimize $\|Ax - b\|^2$

• GPS has location (x, y, z) and clock bias t

receiver measures pseudorange: $d = \sqrt{(x-x_i)^2 + (y-y_i)^2 + (z-z_i)^2} + c \cdot t$

To solve: Gauss newton method

: Levenberg - Marquardt

Class 12: Convex Functions, Convex Sets and Level Sets

Convex vs Monotonic

- Monotonic: function decreases steadily
- Convex: function lies between chord connecting two points

Define

• Convex Function: $f((1-\theta)v + \theta v) \leq (1-\theta)f(v) + \theta f(v)$, $0 \leq \theta \leq 1$

- Meaning in words:
The value of the function at a point in between u and v is less than or equal to the average of the function values at u and v .
- Geometric meaning:
The graph of the function lies below or on the straight line (the "chord") joining the two points $(u, f(u))$ and $(v, f(v))$.
- Why important:
 - Ensures no local minima except the global one \rightarrow optimization is easier and more reliable.
 - Makes constrained optimization problems solvable with efficient algorithms (why convexity is a foundation for ML, optimization, and economics).

- Strictly convex: inequality is strict when $u \neq v$.
- Convex set: contains all linear interpolations between its points.

↓
Convex, if any two points inside it, entire line segment connecting is also in the set.

↪ Taking weighted average between two points

Important Convex Functions

- Affine functions ($f(w) = Mw + c$): always convex
- Quadratics form ($f(w) = w^T K w$): convex if K is positive definite

Operations Preserving Convexity

- Positive Scaling lf is convex if $l \geq 0$
 - sum of convex functions is convex
- Affine transformation: if f is convex, $g(v) = f(Mv + c)$ is convex

Tangent Plane Property

- Convex functions are always above tangent plane
- $f(w) \geq f(w_0) + \nabla f(w_0)(w-w_0)$
- Gradient inequality ensures convexity

Level Sets

- $SL(f, L) = \{w \mid f(w) \leq L\}$
- level sets of convex functions are convex

Class 12: In Person Notes

convexity, find point inside/outside set.
for any point, line is entirely in the set



$$L(\theta) = \theta \cdot 1 + (1-\theta) \cdot 1$$

$$f(\theta) =$$

monotonically increasing $\rightarrow f(\theta)$
on an interval
means θ

linear interpolant above line

Convex on $V \subseteq \mathbb{R}^1$

Recall scalar

Convex sets:

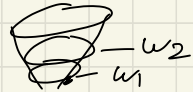
$V \subseteq \mathbb{R}^n$ is convex set
 $\forall \lambda \in [0, 1]$ or $w = v$

Convex function is bowl like
Find points and planes, line is
always underneath.

Example Plane:

$$f(w) = mw + b$$

\uparrow
offset
constant



$$= [w_1 \ w_2] \begin{bmatrix} \\ \\ \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$
$$= \vec{w}^T \kappa \vec{w}$$

Class 14: Neural Networks (single neuron)

Covering

- Activation function
- Optimization with steepest descent
- Vector representation of data
- Back propagation

How do we handle multiple observations?

In practice, $\vec{d} = \sum_{k=1}^m \vec{d}_k$

Structure of Single Neuron

- Artificial neuron: $u = \vec{x} \cdot \vec{w} + b$ and then $z = \phi(u)$
 - x : input vector
 - \vec{w} : weight vector
 - b : bias (Augment \vec{x} with 1)
 - $\phi(u)$: activation function

Common Activation Functions

- Logistic / sigmoid $\phi(u) = \frac{1}{1 + e^{-u}}$ Outputs between 0-1 Deriv: $\phi(u)(1 - \phi(u))$
- ReLU: $\phi(u) = \max(0, u)$ Outputs between 0 - ∞ Deriv 1, if $u > 0$
else 0
↳ Hidden layers for neural nets
- Heaviside: 0 if $u < 0$, 1 if $u \geq 0$ outputs between $\{0, 1\}$ Output Classification

Training via Steepest Descent

- minimized squared error loss: $g(r) = \frac{1}{2}r^2$, where $r = y - \phi(u)$
- Gradient Descent update: $w_{k+1} = w_k + \eta \cdot \vec{d}$, where $\vec{d} = -\nabla f(\vec{w})$
- Backpropagation Terms:
 - Residual error: $r = y - \phi(u)$
 - Descent Direction: $\vec{d} = -\vec{x} \cdot r$

After training, neuron learns a decision boundary

Output Layer Differentiation

- weight vector split into $\vec{w} = [w_1, w_2, \dots, w_n, b]$
- Linear term becomes: $u = \vec{w}_{1..n} \cdot \vec{x} + b$
- Backpropagation over layers

Class 15: Neural Networks: Multiple Layers

Goal: Extend neural net training with grad descent to multi-layer ANN's.
Calculating gradients (derivatives) using back propagation when multiple layers

Adding more layers and nonlinear activations allows neural nets to learn more complex decision boundaries.

Conceptually

- layer is a lin transformation followed by a nonlinear activation
- Linear: matrix vector Product: $wx + b$
- Activation: $\phi(u)$
- Forward pass: compute activations layer by layer
- Backward pass: (Back Propagation): compute gradients layer by layer to update weights

Structure of 3 layer Neural Network

- Layer 1: (Input) Raw input features
- Layer 2: (Hidden) Applies weights and activation functions
- Layer 3: (Output) Computes final prediction from Layer 2's output

Forward Passing

Layer 2: $v^2 = [x, 1] \cdot w^2$

Layer 2 Activation: $a^2 = \phi(v^2)$

Layer 3 Linear output: $v^3 = [\phi^2, 1] \cdot w^3$

Layer 3 Activation: $a^3 = \phi(v^3)$

$\psi = \frac{d\phi}{du}$: derivative of activation function

Gradients use Hadamard Product (element wise): $\psi \odot b$

Backpropagation

- gradients are computed layer by layer
- output layers use errors between prediction and true label (residual)
- hidden layers use chain rule to propagate error backward using jacobian

Each layers gradient is computed based on:

- activation derivative, error propagated from next layer, input data

Class 17: Back-Propagating Scale Factors of Gradient Components

- How propagating scale factors backward through neural networks instead of computing ^{Jacobian} matrices
- Scale factors tells how much error flows backward through each neuron

Forward and Backward Computation

- Forward Evaluation: Compute the output $z(w)$ of the network given an input x and then compute the objective function (squared error)
- Backward: Compute gradients layer by layer using chain rule.
 - Uses scale factors (partial derivatives) to propagate error backward

Optimization Formulation

- Neural networks are treated as an unconstrained optimization problem:

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

where η is the learning rate and f is the squared error function

- 2-layer network:

Layer 1 (input): 2 features

Layer 2 (Hidden): 2 neurons

Layer 3 (Output): 1 neuron

- Each layer has weights and activations, denoted $\phi(t)$ and its derivative $\psi(t)$

Backpropagation Steps

- residual: $r = y - z(\vec{w})$

- output layer scale factor: $b_3 = -r$

- Compute layer 3 gradient: $\vec{d}_3 = -\vec{x}_3 \cdot \psi_3 \cdot b_3$

- Back-propagate to L2: $\vec{b}_2 = \vec{w}_3 \cdot \psi_3 \cdot b_3$
 $\vec{d}_2 = -\vec{x}_2 \cdot \psi_2 \cdot b_2$

- Final descent vector: \vec{d}_3 and \vec{d}_2

- Equations: Activation: $\phi(t) = \frac{1}{1+e^{-t}}$, $\psi(t) = \phi(t)(1-\phi(t))$

Forward Evaluation Example

- Given: weights and observation vector
↳ Activation function: logistic

- Forward computations:

Layer 2 outputs: $\vec{v}_2, \vec{\phi}_2, \vec{p}_2$

Layer 3:

$$v_3 = \vec{x}_3 \cdot \vec{w}_3$$

$$\phi_3(v_3) = z(\vec{v}_3)$$

Observation:

- derivatives (ψ) are numerically small \rightarrow
training progresses at different rates
across layers

Class 18: Constrained Optimization Problems

Purpose: introduce constrained optimization. Point is to minimize objective function
subject to constraint be equality or inequality-based.

- Feasible sets: Allowed solutions
- Solve both Quadratics/Linear objectives
- 3 classic optimization types:

Quadratic Programming (QP), Convex Optimization, Linear Programming

Constrained Optimization

- Minimizing a function (squared distance) subject to constraints.
- 2 constraint types: Equality constraint $Mw = c$
: Inequality constraint $Mw \leq c$

Linear Equality Example

- $-w_1 + 2w_2 = 4$
(two variables)
- Closest point given constraints
↳ subject to $[-1 \ 2]w = 4$

Linear Inequality

- $w_1 - w_2 \leq 4$, feasible region in half space. Optimization form:
 $\min \|w - w_0\|^2$. Subject to $[1 \ -1]w + 4 \leq 0$.

Convex Optimization Basics

- Property function $p_i(w)$ is convex \rightarrow defines constraint
- Feasible point: satisfies all $p_i(w) \leq 0$
- Feasible set F : is all such w values that satisfy the property function/constraint
- If $Mw - c \leq 0$ defines constraint and M is full rank, feasible region is convex

Standard form for Convex Optimization

- $\min f(w)$, subject to $p_i(w) \leq 0$, $Mw = C$
 $f(w)$ is convex, $p_i(w)$ are convex
- Quadratic Programming: $\min w^T K w + a^T w$
 $Aw - b \leq 0$, $Mw = C$
where K is symmetric, positive semidef

Examples:

Class 20: Lagrange Multipliers for 2D Functional Convex Problems

Purpose: Lagrange multipliers help solve equality constraints. Shows how minimizer lies on level curve of objective function and constraint function.

The key takeaway is: At the optimal point, the gradients of the objective and constraint are parallel and scaled versions of each other, giving rise to the Lagrange multiplier condition.

Example problems

- Level curve and constraint curve.
- Looking for point where lowest possible level curve just touches the constraint. Geometrically this corresponds to gradients being parallel.

Ex 1

- Squared Length + Linear Equality
- Objective: Minimize $f_1(w) = \|w\|^2 = w^T w$
- Constraint: $P_1(w) = [-1 \ 1]w + 3 = 0$



Increasing values of l , (radius)

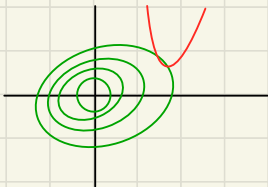
Optimal point $w_1^* = [1.5, -1.5]$

↓

$$\nabla f(w^*) = \mu \nabla P(w^*)$$

Ex 2

- Weighted Norm and Quadratic Constraint
- Objective: $f_2(w) = \frac{1}{2} w^T K w$
- Constraint: $P_2(w) = (w_1 - 3)^2 + 1 - w_2 = 0 \rightarrow$ parabola



Increasing l , will first touch

$w_2^* = [2.158, 1.714]$

Satisfies Lagrange condition

$$\nabla f_2(w^*) = \mu \nabla P_2(w^*)$$

- 1) Compute gradients
 $\nabla f(w) = [2w_1, 2w_2]$
 $\nabla P(w) = [-1, 1]$
 $[2w_1, 2w_2] = \mu [-1, 1]$
- 2) $2w_1 = -\mu$, $2w_2 = \mu$
 $w_1 = -\frac{\mu}{2}$, $w_2 = \frac{\mu}{2}$
- 3) Apply constraint
 $-w_1 + w_2 + 3 = 0$
 $\mu = -3$
 $w_1 = -\frac{-3}{2} = 1.5$, $w_2 = \frac{-3}{2}$
 $w^* = [1.5, -1.5]$, $\mu = -3$

Ex 3

Affine Objective + Quadratic Constraint

• Objective $f_3(w) = [1, 1] \vec{w} = w_1 + w_2$
 $p_3(\vec{w}) = (w_1 - 1)^2 + 1 - w_2 = 0$



Using Lagrange multipliers:

$$\nabla f_3(w^*) = \nu_3 \nabla p_3(w^*)$$

$$[1, 1] = \nu_3 [2w_1 - 2, -1]$$

$$1 = \nu_3 (2w_1 - 2)$$

$$1 = \nu_3 (-1)$$

$$1 = -\nu_3$$

$$1 = -1(2w_1 - 2)$$

$$1 = -2w_1 + 2$$

$$\frac{-1}{-2} = \frac{-2w_1}{-2}$$

$$w_1 = 0.5$$

Plug

$$(w_1 - 1)^2 + 1 - w_2 = 0$$

$$(0.5 - 1)^2 + 1 - 0.5 = 0$$

$$(0.5 - 1)^2 = -0.5$$

$$\nu_3 = 1.25$$

Insert
lecture photo
or meshes

Core Idea: $\nabla f(w^*) = -\nu \nabla p(w^*)$

ν is the Lagrange multiplier

gradients are parallel and opposite

- Linear Constraint, flat vertical plane
- Quadratic Constraint, curved surface

Class 21: Constrained Optimization using Lagrange Multipliers

- Introduces Lagrange function, Lagrange equation.
- Finding minimizer and multiplier using matrix algebra
- Solves using KKT system (Jacobian + gradient)

Lagrange function and Equation

$$L(w, \nu) = f(w) + \nu \cdot p(w)$$

Objective \rightarrow Lagrange multiplier
 $p(w) = 0$ is constraint

$$\frac{\partial L}{\partial w} = 0, \quad \frac{\partial L}{\partial \nu} = 0 \quad \text{Solve for both } w^* \text{ and } \nu^*$$

Ex 1: Squared Length + Linear Equality

- Objective $f_1(w) = w^T I w$
- Constraint $p_1(w) = [-1 \ 1] w + 3 = 0$

$$\text{Lagrangian: } L_1(w, \nu) = w^T w + \nu ([-1, 1] w + 3)$$

$$w^* = [1.5, -1.5], \quad \nu = 3$$

$$\text{Solve: } \begin{bmatrix} 2I & [-1; 1] \\ [-1, 1] & 0 \end{bmatrix} \begin{bmatrix} w^* \\ \nu^* \end{bmatrix} = \begin{bmatrix} 0 \\ -3 \end{bmatrix}$$

Lagrange function is $L(w, \nu) = f(w) + \nu p(w)$
Objective + Scalar multiple

- minimizer of $L(w, \nu)$ must be a stationary point.

EXAMPLE 2: Mechanical system with 2 springs

- Objective: minimize potential energy of 2 springs: $f_2(w) = \frac{1}{2} w^T K w$
- Constraint: $p_2(w) = m w - c_0 = 0$

$$\text{where } K = \text{diag}(k_1, k_2), \quad m = [-k_1, k_2]$$

$$\text{• Lagrangian system: } L_2(w, \nu) = \frac{1}{2} w^T K w + \nu (m w - c_0)$$

$$\text{• Solve: } \begin{bmatrix} K & m^T \\ m & 0 \end{bmatrix} \begin{bmatrix} w^* \\ \nu^* \end{bmatrix} = \begin{bmatrix} 0 \\ c_0 \end{bmatrix}$$

$$\text{Sample: } k_1 = 1, k_2 = 2, c_0 = 6$$

$$w^* = [-2; 2], \quad \nu = -2$$

General Case: Quadratic objective + Linear Equality Constraints

• Problem Form: $\min f(w) = \frac{1}{2} w^T K w + a^T w$ subject to $Mw = c$

• KKT system:
$$\begin{bmatrix} K & M^T \\ M & 0 \end{bmatrix} \begin{bmatrix} w \\ \lambda \end{bmatrix} = \begin{bmatrix} -a \\ c \end{bmatrix}$$

• used to solve minimization with constraints

• K must be symmetric positive definite, M full rank

• has both positive and negative eigenvalues due to form

Class 23: Dual Formulation of Lagrange Equations

Focus: Newton's method extension

- Taylor Approximation, gradient + Hessian

• Minimizing Functions:

1D Newton's Method: $x_{k+1} = x_k - \frac{f'(x_k)}{f''(x_k)}$

• From 1D to Multivariate:

gradient ∇f replaces slope,

Hessian matrix replaces second deriv

$$: x_{k+1} = x_k - H^{-1} \nabla f(x_k)$$

Feasible set:

$$F = \{ \vec{v} : m\vec{v} - c = 0 \}$$

Taylor Expansion view

$$f(x + \Delta x) \approx f(x) + \nabla f(x)^T \Delta x + \frac{1}{2} \Delta x^T H \Delta x$$

Computing Newton Step

• $H \Delta x = -\nabla f(x)$

• Linear System:

Ex $f(x) = \begin{bmatrix} 6 \\ 4 \end{bmatrix}$, $H = \begin{bmatrix} 20 & 0 \\ 0 & 4 \end{bmatrix}$

$$\Delta x = -H^{-1} \nabla f(x) = \begin{bmatrix} -3 \\ -1 \end{bmatrix}$$

Newton Decrement

• tells us how close we are to optimum

$$\pi(x) = \sqrt{\nabla f(x)^T H^{-1} \nabla f(x)}$$

Stop when $\frac{1}{2} \pi(x)^2 \leq \epsilon$

3 Phases of Convergence

- 1) Damped, take small steps unsure of curve
- 2) Quadratic, extremely fast convergence
- 3) Overshoot, may take large steps

Newton Direction is always a descent direction as long as Hessian H is positive definite.

Lagrangian Function: $b(w, \nu) = f(\vec{w}) + \nu p(\vec{w})$

At a stationary point w^* $\nabla f(\vec{w}^*) = -\nu \nabla p(\vec{w}^*)$

By Solving Dual, you get same as solving primal problem

Dual Function is always concave, even if original function isn't convex

Class 24: Inequality Constraints and KKT Conditions

- Linear inequality constraints define convex sets
- KKT conditions help determine optimal solutions
- Feasible Point: Solution that satisfies all inequality constraints

Problems have objective function and constraints:

$$f(w) = \frac{1}{2} w^T K w + a^T w \quad A\vec{w} \leq \vec{b}$$

$$L(w, \pi) = f(w) + \pi^T (A\vec{w} - b)$$

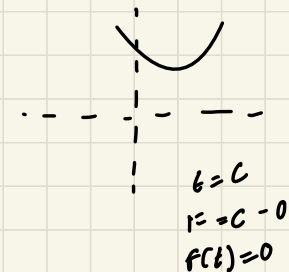
KKT conditions

- 1.) Primal Feasibility: $A\vec{w} \leq \vec{b}$
- 2.) Dual Feasibility: $\pi \geq 0$
- 3.) Stationarity: $K\vec{w} + \vec{a} + A^T \vec{\pi} = 0$
- 4.) Complementary slackness: $\pi_i (a_i \cdot \vec{w} - b_i) = 0$ for each constraint

If all 4 hold, the point is a KKT point which is optimal if the objective is convex.

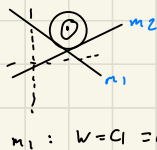
watch videos

Class #24: KKT Conditions for Constrained Optimization



Solution at $b=c$

$$F(w) = \|w - g\|$$



Minimize Objective Intersection

Dual formulation (Big \rightarrow Small)

Consider: 1D quadratic objective

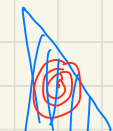
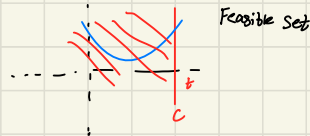
$$f(t) = a_1 t^2 + a_2 t + a_3$$

constrained by $t^* \leq c$

F : Feasible Set



Functions are locally convex



If the unconstrained minimizer is the only its in the feasible set



Intersection of 2 half Spaces



$$\begin{aligned}
 \text{Lagrange Eq} \\
 L(w, \lambda) &= f(w) \\
 &= \lambda_1 m(w_1) + \lambda_2 c_2(w_2)
 \end{aligned}$$

Property P_i is active

KKT Conditions:

For quadratic: $f(w) = \frac{1}{2} w^T K w + r^T w$

For linear inequalities: $A w = \vec{b} = \begin{pmatrix} \frac{\partial f}{\partial x} \\ \lambda \end{pmatrix} \vec{w}$

A point w^* is a constrained minimizer

Primal feasibility: $A w^* = \vec{b}$

Complementary Slackness

Class 25: Geometry AT KKT Points

Purpose: KKT conditions describe how constraints interact with objective functions:

- which constraints are active
- which are inactive

Interpreting Lagrange Multipliers

- $\lambda = 0$: constraint is inactive at optimum (Solution must be in feasibility region)
- $\lambda > 0$: constraint is active (solution lies on boundary defined by constraint)
- $\lambda < 0$: occurs only with equality constraints (Equivalent to 2 constraints)
↳ $\lambda_1 = 0, \lambda_2 > 0$

Main Geometric Takeaways

Quadratic objective + linear inequality constraints

- minimizing quadratic that produces level curves
- Constraints produce half-spaces
- Optimal point found when:
 - 1) objective is minimized
 - 2) constraints are satisfied

Inactive: $\lambda = 0$
active: $\lambda > 0$

Example 1: (3 constrained)

- only 3rd constrained is active $\rightarrow \lambda(0,0)$

KKT Conditions

1.) Primal Feasibility: $A\vec{w} \leq b$

2.) Dual Feasibility: $\lambda \geq 0$

Example 2: • When each constrained is used alone, λ values are larger

3.) Stationarity:

gradient of Lagrangian
equals zero: $K\vec{w} + \vec{q} + A^T\vec{\lambda} = 0$

Equality constraint interpretation

4.) Complementary slackness:

↳ $\lambda(a_i^T w - b_i) = 0$ for each constraint

- Equality $m\vec{w} - c = 0$, is treated as two inequalities
- $m\vec{w} - c \leq 0$
- $-m\vec{w} + c \leq 0$
- only one inequality is active \rightarrow explains negative λ in equality-case solutions

Class 26: Constrained Least Squares and Tikhonov Regularization

- Extension of ordinary least squares by adding constraints/penalties to stabilize solutions

Ordinary Least Squares

Goal: find weights w that make predictions $Xw \approx y$

OLS objective: $\min \|Xw - y\|^2$

Solution:

take gradient and set to 0, gives Normal Equation

$$X^T X w = X^T y$$

Assuming X has full rank: $w^* = (X^T X)^{-1} X^T y$

Issues: • Sensitive to outliers

• test error high (overfitting)

Constrained Least Squares

• Force solution vector to have limited magnitude

$$\bullet \|w\|^2 \leq \epsilon$$

CLS Problem: $\min \|Xw - y\|^2$, s.t. $\|w\|^2 \leq \epsilon$

Lagrange Multiplier: $L(w, \lambda) = \|Xw - y\|^2 + \lambda (\|w\|^2 - \epsilon)$

2 Cases:

1) OLS solution already satisfies constraint

2) OLS violates constraint (need $\lambda > 0$)

$$\hookrightarrow (X^T X + \lambda I) w = X^T y$$

CLS pulls solution closer to smaller norm.

Tikhonov Regularization

• Instead of constraint, add penalty to loss

General form: $T(w, \lambda) = \|Xw - y\|^2 + \lambda \|Rw\|^2$

Special case:

$$\bullet R = I$$

$$\bullet \lambda = \alpha^2$$

$$\min \|Xw - y\|^2 + \lambda \|w\|^2$$

Ridge Regression

Tikhonov, λ is chosen by user

Class 27: Ridge Regression and The Lasso

Goal: Improving Linear regression with messy, high dimensional data

Standardizing Data (Z-Score)

- mean = 0

- variance = 1

$$z = \frac{x - \mu}{\sigma}$$

- Standardized data has lower RMSE

Ridge Regression (L2 Regularization)

- Ordinary Least Square creates nonsense coefficients

When : • variables are correlated

• data is noisy

• many features

Ridge adds penalty to large weight:

$$\min_w \|Xw - y\|^2 + \lambda \|w\|_2^2$$

Shrinks weights towards 0

Lasso (L1 Regularization)

• imposes constraint on sum of absolute values of weights

$$\|w\|_1 \leq 0$$

• Key SuperPower:

sets some coefficients exactly to 0

This means automatic feature selection

$$w_{\text{ridge}} \approx (0.254, 0.967)$$

$$w_{\text{lasso}} \approx (0, 1)$$

Elastic Net

• Combines ridge + lasso and gives best of both worlds

Class 29: Support Vector Machines

Core Idea: SVM finds best separating hyperplanes between 2 classes (class +1, class -1)

Basic Intuition

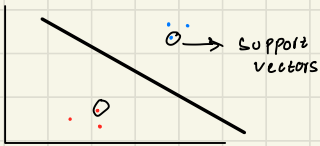
- Hyperplane: $w^T x + b = 0$
 - weight vector
 - bias
- Perceptron can find separating lines, but not the best one.
SVM finds line with largest margin between 2 classes
- Support vectors are data points closest to hyperplane that define edges of margin, hyperplane is midway between closest pos/neg points

Marginal Max

$$y_i(w^T x_i + b) \geq 1$$

Smaller weight vector means flatter slope \rightarrow wider margin
margin = $1/\|w\|$, maximizing margin

minimize $w^T w$



$$\vec{n} = \frac{\vec{w}}{\|w\|}, \quad a = \frac{b}{\|w\|}$$

Margin (r) is distance between H and closest data points.

$$r = \vec{n} \cdot \vec{x}_s + a$$

$$r = x_s^T \vec{n} + a$$

Primal Form:

- minimize $\frac{1}{2} \|w\|^2$
- subject to $y_i(w^T x_i + b) \geq 1$
- impossible to solve normally, use dual formula

$$L(w, b, \alpha) = \frac{1}{2} w^T w - \sum_i \alpha_i [y_i (w^T x_i + b) - 1]$$

1) derive w.r.t w

2) derive w.r.t b

dual:
$$\max_{\alpha} \left(\sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j (x_i^T x_j) \right)$$

only dot products left

Class 30: Dual Formulation of SVM

Goal: Find best separating hyperplane: $\min \frac{1}{2} w^T w$ subject to $y_i (w^T x_i + b) \geq 1$

Lagrange Multipliers

- introduce 1 multiplier for every constraint
- $$L(w, b, \alpha) = \frac{1}{2} w^T w + \sum_i \alpha_i [1 - y_i (w^T x_i + b)]$$

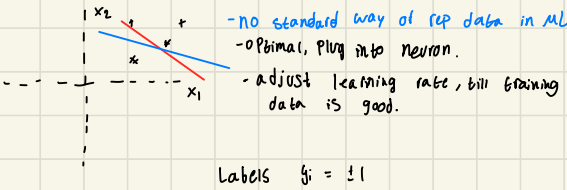
Derivatives

- set gradients to 0 to remove w and b
- w is just weighted combination of data points

Plug back into Lagrangian: $\alpha_i \geq 0, \sum_i \alpha_i y_i = 0$

Once α are found:

- Compute $w = X^T Y \alpha$
- Compute b using support vector $y_i (w^T x_i + b) = 1$
- Classify new data with $\text{sigm}(w^T x + b)$

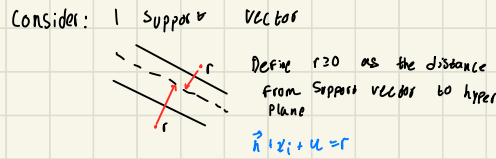


Concept: Hyperplane Separation

H is a unit normal, $\|H\|=1$, bias scalar = a

Idea: maximize margin from hyperplane to data

Define: Support vector is point closest to hyperplane



Sneaky:
use $y_i = \pm 1$
For $y = +1$ have $w^T (x_i + b)$
For $y = -1$ have $y_i (x_i + a) = r$

What if its not a support vector,

- has to be atleast r away
- For $y_i = +1$ have $n \cdot x_i + a < r$
($n \cdot x_i + a > 1$)

For any vector x_i have $|n \cdot x_i + a| \geq r$

Consider an unrestricted vector $\vec{w} \approx 0$

Define $\frac{1}{r} = \|w\|$ s.t. $\vec{n} = \frac{w}{\|w\|}$
or $w = \frac{\vec{n}}{r}$

$w \cdot x_i + b = 1$

Rewrite: $n \cdot x_i + a = r$
 $= \frac{1}{r} (n \cdot x_i + a)$

Maximize: r
Minimize: $\frac{1}{r}$
or minimize: $\frac{1}{\|w\|}$
or minimize: $w^T w$

Class 32: SVM: Soft Margins

Motivation: Data are not always linearly separable

→ when data are not linearly separable, misclassified points create false positives and false negatives

Slack Variables (ξ) Allowing Constraint Violations

• a slack variable turns optimization problem to turn an inequality constraint → equality constraint

→ hard margin constraint: $y_i (w^T x_j + b) \geq 1$

→ soft margin constraint: $y_i (w^T x_j + b) + \xi_j \geq 1, \xi_j \geq 0$

- allow some points inside the margin
- allow some points to be misclassified.

Incorporating Slack Variables

• Instead of ξ as free variables, we regularize them.

Objective: $\min \frac{1}{2} \|w\|^2 + C \sum \xi_j$

• L1 penalty: ξ_j , just sum of violations

• C controls the trade off

- Large C → fewer violations (harder margin)
- Small C → more flexibility (softer margin)

Constraints in matrix Form

1. Modified classification constraint:

$$1 - YXw - by - \xi \leq 0$$

2. Non negativity of slacks

$$-\xi \leq 0$$

Lagrangian:

$$L(w, b, \xi, \alpha, \beta) = \frac{1}{2} w^T w + \alpha^T (1 - YXw - by - \xi) + C \sum \xi_j - \beta^T \xi$$

Dual Formulation:

- Differentiating w.r.t the primal variables using KKT

Conditions shows: $\alpha + \beta = C$

- Final dual constraint $0 \leq \alpha_j \leq C$

Case 0: $\xi = 0$, correctly classified outside margin

Case 1: $0 < \xi < 1$, inside margin correct side

Case 2: Misclassified $\xi_j > 1$ (full mistake)

Class 33: The Kernel Trick

Why kernels: because datasets are not linearly separable

We embed data into higher dimensional space where you can separate it.

$$(u_1, u_2) \rightarrow (u_1, u_2, u_1^2 + u_2^2)$$

Issue: In explicit embedding computing $\psi(x)$ in high dimension is expensive

↳ Solution: SVM only requires dot products between embedded points: $\phi(x_i)^T \phi(x_j)$

Kernel Trick

- Kernel K computes dot product without explicitly forming $\psi(x)$
 $k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$
- replace $x_i^T x_j$ in the dual SVM with $K_{ij} = K(x_i, x_j)$
- This bypasses the embedding and allows nonlinear SVM's

Gram Matrix

$$K_{ij} = K(x_i, x_j)$$

- symmetric
- PSD
- plugged directly into Dual

$$\frac{1}{2} \alpha^T Y K Y \alpha - 1^T \alpha$$

Common kernels

- Linear $k(u, v) = u^T v$
- Polynomial $(u^T v + c)^p$
- Gaussian $\exp\left(-\frac{\|u - v\|^2}{2\sigma^2}\right)$

Class 34: SVM Kernel Classification

Goal: How SVM's classify new points when using kernels

- After training SVM with kernel, must decide how to score a new vector and assign a class label.

SVM: Classification Score

• The score of the data vector is $z_j = x_j^T w + b$

↳ classify using $a_j = \text{sign}(z_j)$

Using Lagrange Multipliers

• Primal Lagrange equation, with equality

$$\vec{w} = x^T \gamma \vec{\alpha}$$

↳ sub into score:

$$z(x_j) = \sum_{i=1}^m \alpha_i y_i (x_i \cdot x_j) + b$$

Kernel Version (Key Formula)

• replace dot product with kernel: $(x_i \cdot x_j) \rightarrow K(x_i, x_j)$

↳ Final classification score: $z(x) = \sum_{i=1}^m \alpha_i y_i K(x_i, x) + b$

• only support vectors ($\alpha > 0$)

• classification $a(x) = \text{sign}(z(x))$

Gaussian (RBF) Kernel Behaviour

• Gaussian kernel: $K(u, v) = e^{-\frac{\|u-v\|^2}{2\sigma^2}}$

• Close points \rightarrow kernel ≈ 1

• Far points \rightarrow kernel ≈ 0

• Classification becomes non linear curves

• SVM separating ± 1 uses curved boundaries

• Adding one new point changes data boundary significantly

Class 33: Kernel Trick - In person NOTES.

Concept: instead of doing embedding use kernel function in a lower dimension

How can we use a kernel in SVM:

Ex Kernel function of $u \cdot v^T + (u \cdot v)^2$
Embedding costs you, twice (space, time)

Recall Lagrange dual Function

$$\mathcal{L}(\vec{\alpha}, \beta) = \frac{1}{2} \alpha^T X^T X X^T Y Z + \beta^T \alpha$$

Rewrite individual terms

$$Y X X^T \begin{bmatrix} y_1 & & & \\ & y_2 & & 0 \\ & & & y_n \\ 0 & & & y_n \end{bmatrix} \begin{bmatrix} \\ \\ \\ \end{bmatrix}$$

In vector terms, v_x

$$\mathcal{L}(\vec{v}, b) = \frac{1}{2} Y^T K + \beta^T Z$$

replacing LSVM

Common Functions

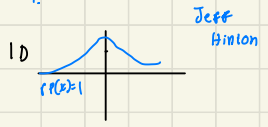
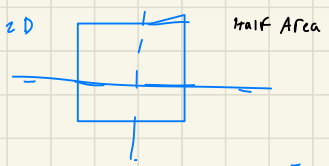
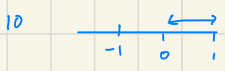
Linear: $k(v, v) = v^T v$

Quadratic: $\psi(v, v) = (v^T v + c)^2$

where $c \geq 0$ (works well up to 6 dimensions)

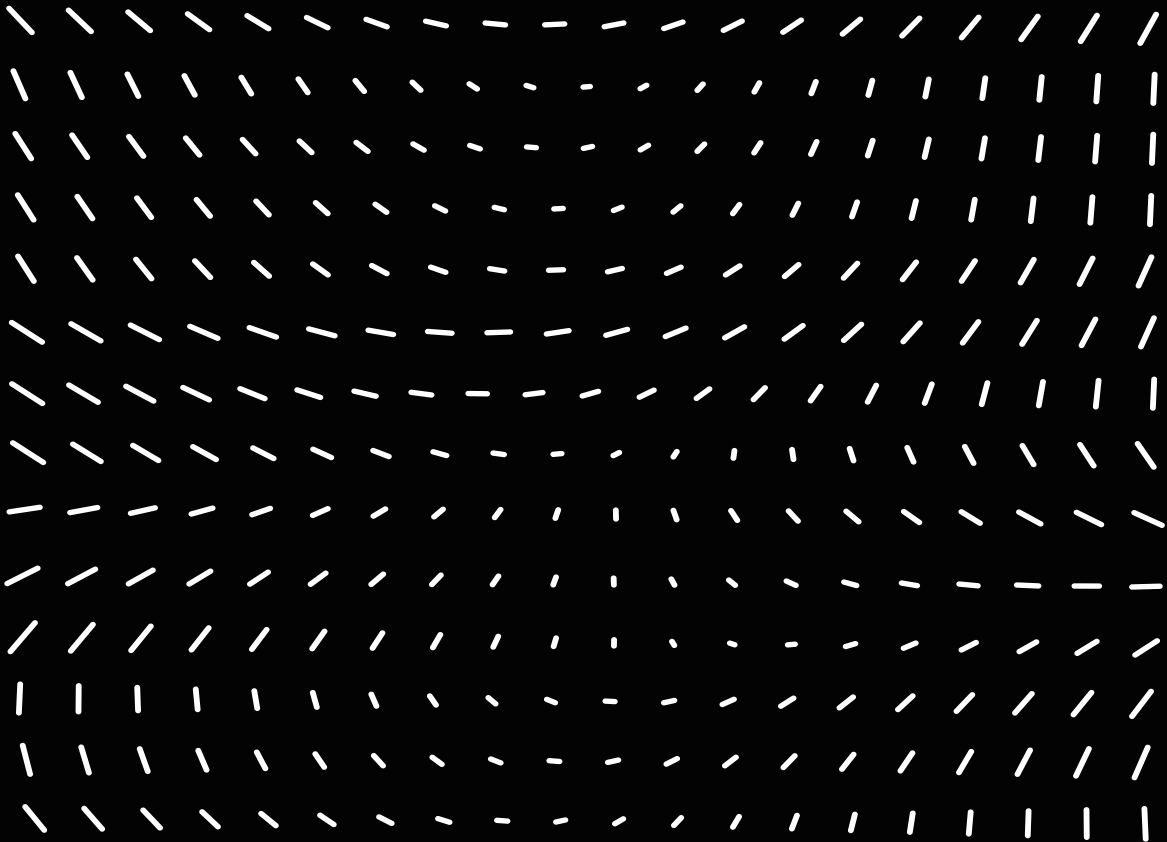
Gaussians: $(v, v) = \frac{1}{c} \exp(-\frac{1}{2c} \|v\|^2)$

Confirm data



Quiz 1

Prep

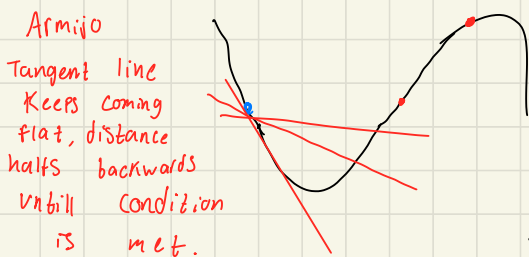


Quiz 1 Prep: Main concepts

- Scalar Calculus
- Derivatives, first and second
- Stationary points, convexity, second derivative

Scalar Optimization

- Taylor series approximations
- Iterative methods:
 - ↳ Step size iteration $t_{k+1} = t_k - s f'(t_k)$
- Backtracking Armijo method
 - ↳ α, β parameters



Fixed

- avoid oscillation in 3D
- Backtracking is a better way to find fixed step size

Taylor Series / Quadratic

input Point, first deriv, second deriv

↓ gives function that's easier to work with
Set that to 0, find min.

Hopefully that min is close to og function min

Partial Derivatives

rate of change of f , with respect to variables
for direction \vec{v}

key relation to partials

$$D_v f = \nabla f \cdot v$$

∇f always points in direction of steepest ascent

Directional derivative in which direction you ask for

Practise 1: Prep for Online Quiz 1

Content: classes 1-5
(no hessian)

Main Themes

- Find first or second derivatives in matlab
- Compare 3 1-D optimization methods
- Determine the kind of Stationary point in 2D (using Hessian)

Question 1: Scalar optimization

- given 3 functions, start at right endpoint
- 3 methods to implement:
 - ↳ Quadratic Taylor
 - ↳ Fixed Step Size GD:
 - ↳ Armijo Backtracking

Stopping rule: $|f'(t_k)| \leq 10^{-6}$

Deliver: Table with estimated minimizer and iteration count for each method x function combo.

1) plug in FD, SD, Function itself

f = function

g = first derivative

h = second derivative

2) unify functions allow us to handle

f , given scalar t , returns two outputs

Fixed stepsize, Armijo functions expect obj + gradient

3) compute steepest descent / backtracking estimates

$w = b - a$, a lot of setup

4) steep fixed function

- moves downhill along negative gradient with constant step

5) Steepline function

- direction of $+1/-1$, move opposite to slope
- user armijo condition (be under a certain line, otherwise keep on shrinking it)

Question 2: Stationary Points

- Finding and classifying stationary points using gradient and Hessian

↳ Define 3 functions

compute gradient and its square

solves $\nabla f(w) = 0$ for real stationary points

evaluate Hessian at each point, take eigen values

↳ print table with point, eigen values and classification

Week 1: Homework 1

- Differentiate simple scalar functions
- Find stationary points (when first derivative = 0)
- Use second deriv test to classify points

5 functions

$$f_1(t) = 2t^3 - 9t^2 + 5$$

$$f_2(t) = t^3 - 12t + 8$$

$$f_3(t) = 3t^4 + 4t^3 - 72t^2 + 1$$

$$f_4(t) = t^4$$

$$f_5(t) = \frac{-t}{t^2+1}$$

Bounds $t \in [-5, 5]$

$$f_5(t) = \frac{-t}{t^2+1}$$

$$f_5'(t) = \frac{(-1)(t^2+1) - (t)(-2t)}{(t^2+1)^2}$$

$$0 = \frac{t^2-1}{(t^2+1)^2}$$

$$t = \pm 1$$

$$f_5''(t) = \frac{(2t)(t^2+1)^2 - (t^2-1)(2)(t^2+1)(2t)}{((t^2+1)^2)^2}$$

$$= \frac{2t(-t^2+3)}{(t^2+1)^3}$$

$$f_5''(-1) = -0.5 \text{ local max}$$

$$f_5''(1) = 0.5 \text{ local min}$$

1) $f_1(t) = 2t^3 - 9t^2 + 5$

$$f_1'(t) = 6t^2 - 18t$$

$$0 = 6t^2 - 18t$$

$$0 = 6t(t-3)$$

$$t=0 \quad t=3$$

Stationary points @ $(0, 0)$ and $(3, 0)$

$$f_1''(t) = 12t - 18$$

$$f_1''(0) = -18 \therefore \text{local max}$$

$$f_1''(3) = 18 \therefore \text{local min}$$

2) $f_2(t) = t^3 - 12t + 8$

$$f_2'(t) = 3t^2 - 12$$

$$0 = 3(t^2 - 4)$$

$$t = \pm 2$$

$$(2, 0), (-2, 0)$$

$$f_2''(t) = 6t$$

$$f_2''(2) = 12 \therefore \text{local min}$$

$$f_2''(-2) = -12 \therefore \text{local max}$$

$$f_3(t) = 3t^4 + 4t^3 - 72t^2 + 1$$

$$f_3'(t) = 12t^3 + 12t^2 - 144t$$

$$0 = 12t(t^2 + t - 12)$$

$$t=0 \quad t=4 \quad t=-3$$

$$f_3''(t) = 36t^2 + 24t - 144$$

$$f_3''(0) = -144 \therefore \text{local max}$$

$$f_3''(-3) = 108 \therefore \text{local min}$$

$$f_3''(4) = 528 \therefore \text{local min}$$

Week 2: Homework 2

Workflow:

1) gradient ∇f

2) Directional Derivative at Point w_0 in Vector v

3) One Step Steepest descent

↳ one step from w_0 in the direction $-\nabla f(w_0)$ with step size $s_0 = 1$

4) Solve $\nabla f(w) = 0$ for (w_1, w_2)

1) $f_1(w) = w_1^2 + \cosh(w_2)$

gradient is partial derivatives

$$\nabla f(w) = (2w_1, \sinh(w_2))$$

2) Directional Derivative in direction v

$$D_v f(w) = \nabla f(w) \cdot v$$

compute gradient at $\vec{w}_0 = \begin{bmatrix} -1 \\ 0 \end{bmatrix}$

$$\hookrightarrow (-2, 0)$$

$$D_v = (-2, 0) \cdot (1, 1) = -2$$

Slope is in direction $v = (1, 1)$ is -2

3) one Steepest descent update

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

$$w^1 = (-1, 0) - 0.1(-2, 0) = (-0.8, 0)$$

4) Stationary Point

$$\nabla f(w) = (2w_1, \sinh(w_2)) = (0, 0)$$

$$2w_1 = 0$$

$$\sinh(w_2) = 0$$

$$w^* = (0, 0)$$

Function $f_3(w) = (w_2 - w_1)^2 + w_1^3$

$$w_0 = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \quad v = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

1) gradient

$$\nabla f(w) = (2w_1 - 2w_2 + 3w_1^2, 2(w_2 - w_1))$$

2) DD at w_0 $\nabla f(-1, 0) = (1, 2)$

along v

$$\begin{aligned} \nabla f(-1, 0) \cdot v &= (1, 2) \cdot (1, 1) \\ &= 3 \end{aligned}$$

3) Steepest desc Step

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

$$\begin{aligned} w^{(1)} &= (-1, 0) - 0.1(1, 2) \\ &= (-1.1, -0.2) \end{aligned}$$

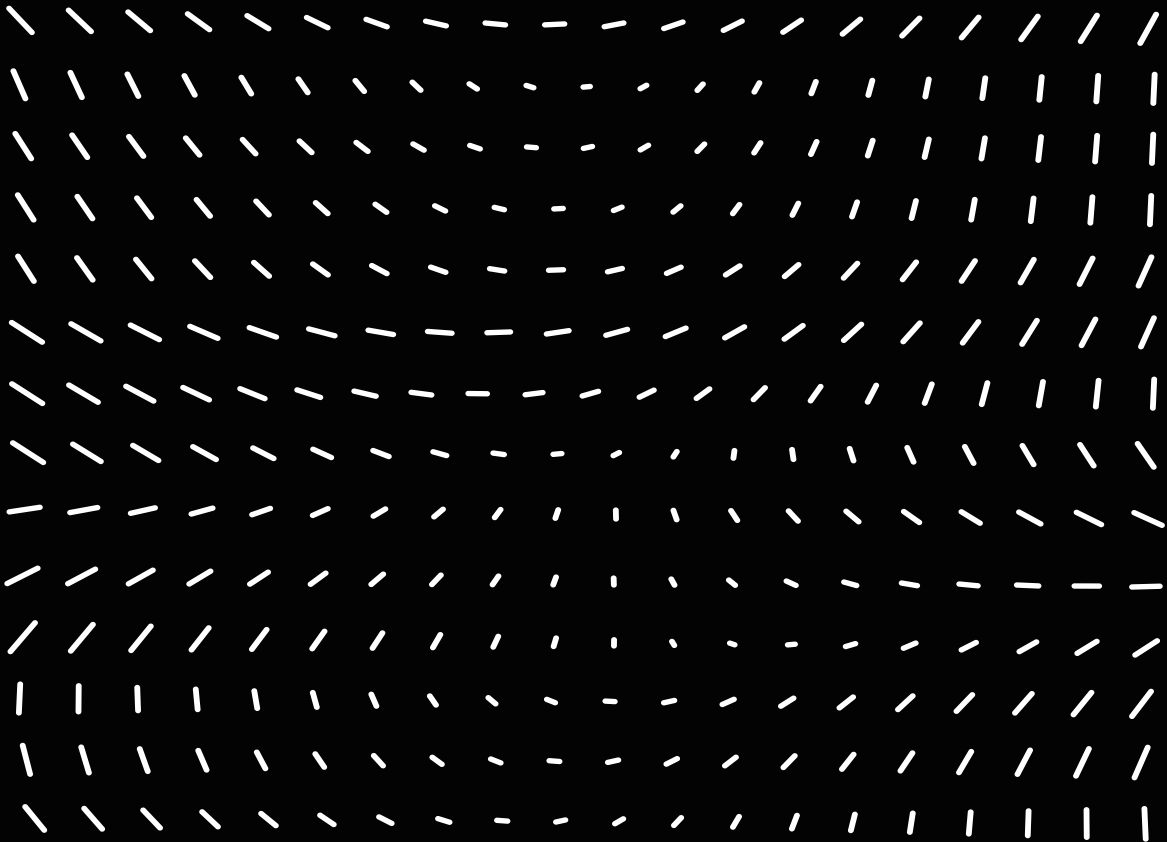
4) Stationary points

$$w^* = (0, 0)$$

Leave Hessian Stuff

Test 1

Prep



Class II: Test 1 Prep

Basic Scalar Calculus

- Derivatives for basic functions
- Chain rule, product rule, quotient rule
- Find Stationary Points, Classify Point
- Convexity (second deriv test)

Scalar Optimization

- Iterations with fixed stepsize (potential problems): $t_{k+1} = t_k - \beta f'(t_k)$
 - Initial guess
 - fixed stepsize
 - deriv of function @ initial guess
- moves opposite to gradient
 - Too Big \rightarrow Divergence
 - Too small \rightarrow slow convergence
- Stepsize selection, backtracking, armijo method

\hookrightarrow Dynamically choose good stepsize, armijo condition

\hookrightarrow Armijo Condition:

$$f(t_k + \beta \gamma s_0 d_k) \leq f(t_k) + \alpha_k \beta \gamma s_0 d_k$$

new function value after step

initial guess

shrink factor

current shrink factor

initial step size

$-f'(t_k)$ descent direction

predicted using slope margin

- Taylor series approx, with/without derivatives \rightarrow Explain grad desc
 - \hookrightarrow This helps predict function behaviour near a point

First order: $f(t) \approx f(t_0) + f'(t_0)(t - t_0)$

Second order $f(t) \approx f(t_0) + f'(t_0)(t - t_0) + \frac{1}{2} f''(t_0)(t - t_0)^2$

$f''(t_0) > 0$: curves upward U (min) (convex)

$f''(t_0) < 0$: curves down \cap (max) (concave)

Vector Calc

• Partial derivatives, forms gradient

• Directional deriv $D_v f(\vec{w}_0) = \nabla f(\vec{w}_0) \cdot \vec{v}$

gradient

direction vector

• Gradient Points in direction of steepest increase

Fall 2024 - Practise Test 1

1) $f(t) = \sin(\cos(t))$

$$f'(t) = \cos(\cos(t))(-\sin(t))$$

$$= -\cos(\cos(t)) \sin(t)$$

$$\begin{aligned} \cos(0) &= 1 \\ -\sin(0) &= 0 \\ -\cos(0) &= -1 \\ &= 0 - \frac{0}{-1} \\ &= 0 \end{aligned}$$

2) $f_2(t) = \cos(t)$, $t_0 = 0$

Find Quad approx?

$$f(0) = 1, f'(0) = 0, f''(0) = -1$$

$$f(t) = \cos(t)$$

$$f(0) = 1$$

$$f'(t) = -\sin(t), f'(0) = 0$$

$$f''(t) = -\cos(t), f''(0) = -1$$

Plug in

$$h_2(t) = 1 + 0 + \frac{1}{2}(-1)t^2$$

$$= 1 - \frac{t^2}{2}$$

3) For $f(t) = t^3 + 3t^2$ describe stationary point $t^* = 0$

$$f'(t) = 3t^2 + 6t$$

$$= 3t(t + 2)$$

$$t = 0 \quad t = -2$$

$$f''(t) = 6t + 6$$

$$f''(0) = 0 + 6 \therefore \text{Pos} \therefore \text{local min}$$

4) $f(t) = t^3 - 12t + 11$, $t_1^* = -2$, $t_2^* = 2$

$$f'(t) = 3t^2 - 12$$

$$f''(t) = 6t$$

↓
local max
↓
concave

↓
local min
↓
convex

5) $f(t) = t^2 - t$, $t_1 = 1$, $S = 1$

find new estimate of minimizer

$$t_{k+1} = t_k - S f'(t_k) \quad f'(t) = 2t - 1$$

$$t_{k+1} = 1 - 1(1)$$

$$t_2 = 0$$

Point
Stepsize
of Steps

* Review

6) $f(t) = t^2 - t$, $\beta = \frac{1}{2}$, $t_1 = 1$, $S = 1$

find stepsize $S = \beta \gamma_{s_0}$ after one

step of armijo backtracking

$$f(t) = t^2 - t, f'(t) = 2t - 1$$

$$f'(1) = 1$$

$$d_k = -f'(t) = -1$$

$$\alpha_k = \frac{1}{2}$$

$$\text{RHS: } \alpha_k \cdot S \cdot d_k = -0.5$$

$$t_{\text{trial}} = t_1 + S \cdot d_k = 0$$

$$0 \leq -0.5$$

False

Shrink to 0.5

$$-0.25 \leq -0.25$$

$$\text{RHS: } t_{\text{trial}} = 1 + 0.5(-1) = 0.5$$

$$\therefore s_0 = 0.5$$

$$\text{LHS: } 0 + \frac{1}{2} \cdot 0.5 \cdot (-1) = -0.25$$

Fall 2024 - Test 2A Practise

1) $f_1(\vec{w}) = \sin(w_1) + \cos(w_2)$, find $\nabla f(\vec{w})$

$$\nabla f(\vec{w}) = \left[\frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2} \right]$$

$$\frac{\partial}{\partial w_1} \sin(w_1) = \cos(w_1) \quad \nabla f(\vec{w}) = ?$$

$$\frac{\partial}{\partial w_2} \cos(w_2) = -\sin(w_2)$$

2) $f(\vec{w}) = w_1^2 + (w_2 + 1)^2$ $\vec{w}_0 = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$, $\vec{v} = \begin{bmatrix} -2 \\ -1 \end{bmatrix}$

Find directional derivative $\text{Dr} F(\vec{w}_0)$

$$\nabla f \cdot \vec{v}$$

$$\nabla f = 2w_1 + 2(w_2 + 1) @ \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

$$\nabla f = [4, 4]$$

$$\begin{aligned} \text{Dr} f(\vec{w}_0) &= [4] \begin{bmatrix} -2 \\ -1 \end{bmatrix} \\ &= -8 - 4 \\ &= -12 \end{aligned}$$

3) $f_3(\vec{w}) = (w_1 + 1)^2 + w_2^2$ has a stationary point at $w^* = \begin{bmatrix} -1 \\ 0 \end{bmatrix}$

Hessian is matrix of second order derivatives

$$\begin{aligned} \text{at } w^* &= 2(w_1 + 1) = 2w_2 \quad ? \text{ Check } \\ &= 2w_1 + 2 = 0 \quad \text{Second Deriv} \\ &= 2(-1) + 2 = 2 \\ &= 2 \end{aligned}$$

$$\nabla^2 f = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

Eigenvalues: 2, 2 \rightarrow imply locally convex, min

All pos	min	convex
All neg	max	concave
mixed	Saddle	

\therefore Positive Definite,

4) $f(\vec{w}) = w_1^2 + 2w_2^2$, $w_0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$, $s_0 = 1$

find \vec{w}_1 , that is in direction of steepest descent with stepsize $s_0 = 1$

$$\vec{w}_1 = w_0 - s_0 \cdot \nabla f(w_0) \quad \vec{w}_1 = \begin{bmatrix} -1 \\ -3 \end{bmatrix}$$

$$\vec{w}_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} - 1 \cdot \begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

$$\vec{w}_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

5) $f(\vec{w}) = w_1^2 - w_2^2 + w_2^4$, $\vec{w}_0 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$

find \vec{d} descent vector using newtons method

$$\vec{d} = -H^{-1} \nabla f(w_0)$$

$$\vec{d} = -[\nabla^2 f(w_0)]^{-1} \nabla f(w_0)$$

$$\nabla f = [2w_1, -2w_2, +4w_2^3]$$

$$\nabla f = [-2, -2(-1), +4(-1)^3]$$

$$\nabla f = [-2, -2, -2]$$

$$\nabla^2 f = \begin{bmatrix} 2 & 0 \\ 0 & -2 + 12w_2^2 \end{bmatrix}$$

$$\text{at } w_2 = -1 = \begin{bmatrix} 2 & 0 \\ 0 & 10 \end{bmatrix}$$

$$\text{Inverse } [\nabla^2 f]^{-1} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/10 \end{bmatrix}$$

$$= - \begin{bmatrix} 1/2 & 0 \\ 0 & 1/10 \end{bmatrix} \begin{bmatrix} -2 \\ -2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0.2 \end{bmatrix}$$

6.) $f(w) = (w_1 + 1)^2 + w_2^2$

Both terms are quadratic

Both are convex \therefore this general function is convex

Use Hessian

Fall 2024 - Test B - Practise

1) $f(w) = \cos(w_1^2) + \sin(w_2^2)$
 $\nabla f(w) = -2 \sin(w_1^2) + \cos(w_2^2)$

2) $F(\vec{w}) = (w_2 + 1)^3 + w_1^3$, $w_0 = [1]$, $\vec{v} = [-1]$
 $\nabla f \cdot v$

$$3(w_2 + 1)^2 + 3w_1^2$$

$$3(2)^2 + 3$$

$$12, 3 \quad \begin{bmatrix} 12 \\ 3 \end{bmatrix} \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

$$= -12 - 3$$

$$= -15$$

3) $f(w) = w_1^2 - (w_2 + 1)^2$, stationary point at $w^* = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$

Find eigenvalues, describe convexity at point w^* .

$f'(w) = 2w_1 - 2(w_2 + 1)$
 $f''(w) = 2 \quad -2$

If variables are left in Hessian plug in your point.

$H = \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix}$ \therefore Saddle Point
 Curves upward then downward
 convex, then concave

4) $f(w) = 2w_1^2 + w_2^2$, $w_0 = [1]$, $s_0 = 1$

$t_{k+1} = t_k - \text{dKSK} \rightarrow \text{gradient}$
 $t_{k+1} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 4 \\ 2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$
 $d = \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 4 \\ 2 \end{bmatrix}$
 $= 4w_1, 2w_2 = \begin{bmatrix} -3 \\ -1 \end{bmatrix}$
 $= 4, 2$

6) $f(w) = w_1^2 - (w_2 + 1)^2$
 $= 2w_1 - 2(w_2 + 1)$
 $= 2 \quad -2$

$\begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix}$ \therefore Saddle Point
 The function is convex then concave

5) $f(w) = w_1^2 + 3w_2^2 - w_2^4$, $w_0 = [-1]$

Find desc vector using N.M using Hessian, which requires invertible / positive definite

$\nabla f = 2w_1 + 6w_2 - 4w_2^3$
 $\nabla^2 f = 2 + 6 - 12w_2^2$

$H = \begin{bmatrix} 2 & 0 \\ 0 & 6 - 12w_2^2 \end{bmatrix}$ plug in
 $= \begin{bmatrix} 2 & 0 \\ 0 & -6 \end{bmatrix}$ Newtons Method
 Assume Positive Definite

2019 Test 1A Practise

1) $f(t) = te^t$

Stationary point at $t = -1$

$$f'(t) = e^t \cdot (t + 1)$$

$$f''(t) = e^t + (t + 1)e^t$$

$$f''(-1) = 0.367 \therefore \text{Pos}$$

$\therefore \text{min, Convex}$

Answer: A

2) $f(t) = \cosh(t) = \frac{e^t + e^{-t}}{2}$ $t_1 = 1$

Using Quadratic approx to find a t_2

$$t_2 = t_1 - \frac{f'(t_1)}{f''(t_1)}$$

$$f'(t) = \sinh(t)$$

$$f''(t) = \cosh(t)$$

$$t_2 = 1 - \frac{1.175}{1.543}$$

$$t_2 = 0.2395$$

Answer A

$$t_{k+1} = t_k - \frac{f'(t_k)}{f''(t_k)}$$

3) $f(t_1 + \beta^r s) \leq f(t_1) + \frac{\alpha}{2} \beta^r s |f'(t_1)|^2 d$

To find number of r backtracking steps:

1) Find direction and slope

$$f(t) = t^2 - 1$$

$$f'(t) = 2t$$

$$f'(1) = 2, d = -2$$

Trial Point: $t = 1 + 1 \cdot 2 \cdot -2$
 $t = -3$
 $f(-3) = 8$

$$8 \leq -4$$

Third Trial	RHS
$f(1 + 0.5^2 \cdot 2 \cdot -2)$	$0 + \frac{0.5}{2} (1/2)^2 \cdot 2 \cdot (4) \cdot -2$
$f(0) = -1$	$= -1$

$\therefore r = 2$

4) $f(\vec{w}) = \frac{1}{3} w_1^3 + w_1 w_2^2 + w_1 w_2 + 5$

$$H = \begin{bmatrix} 2w_1 & 2w_2 + 1 \\ 2w_2 + 1 & 2w_1 \end{bmatrix}$$

$$H = \begin{bmatrix} 2(0.5) & 2(-0.5) + 1 \\ 2(-0.5) + 1 & 2(0.5) \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \therefore \text{local minimum}$$

2nd Trial

$$f(1 + 0.5 \cdot 2 \cdot -2)$$

$$f(-1) = 0$$

RHS

$$0 + \frac{0.5}{2} \cdot 0.5^2 \cdot 2 \cdot 4 \cdot -2$$

$$= -2$$

$$0 \leq -2$$

5) gradient of $\tanh(\vec{w} \cdot \vec{x})$

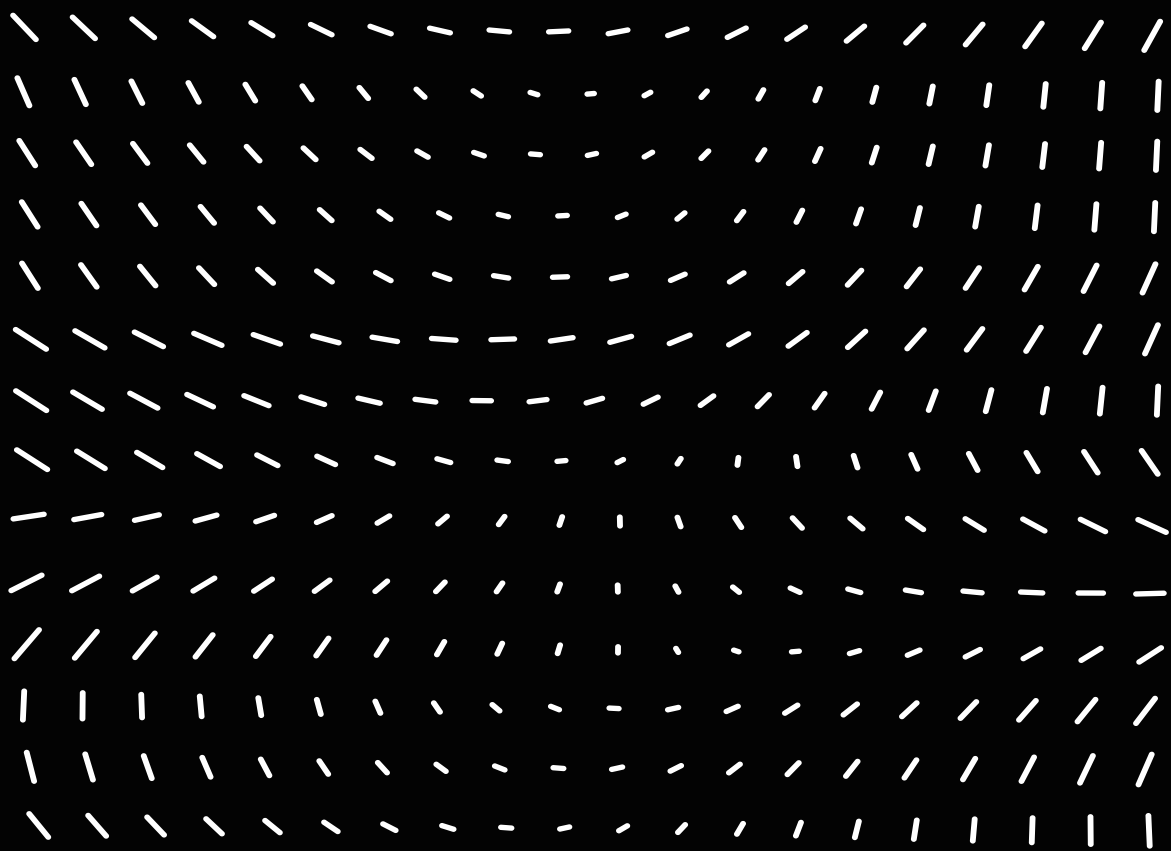
Chain rule: $\frac{1}{\cosh^2(w \cdot x)} \cdot \vec{x}$

6) $f(\vec{w}) = \frac{1}{3} w_1^3 + w_1 w_2^2 + w_1 w_2 + 5$

$$\vec{w}_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Quiz 2

Prep



Key concepts - Practise 2

- Find Stationary points for 2 variables
- Hessian to classify
- Iterative solvers; second order
- Compare convergence behaviour
- Formulate and solve nonlinear squares problem

Part 1: Stationary points

- 3 Functions
 - ↳ Compute Hessians, gradients
 - Solve $\nabla f = 0$ to find stationary points
 - Classify, visualize w mesh / contour

1) Set gradient to 0 to find stationary points
 $\nabla f = 0$

EX

$$f(w) = 2w_2^3 - 6w_2^2 + 3w_1^2 w_2$$

$$\frac{\partial f}{\partial w_1} = 6w_1 w_2, \quad \frac{\partial f}{\partial w_2} = 6w_2^2 - 12w_2 + 3w_1^2$$

$$0 = 6w_1 w_2$$

$$0 = 6w_2^2 - 12w_2 + 3w_1^2$$

2) Hessian

$$H = \begin{bmatrix} 6w_2 & 6w_1 \\ 6w_1 & 12w_2 - 12 \end{bmatrix}$$

$$\text{At } (0, 0) = \begin{bmatrix} 0 & 0 \\ 0 & -12 \end{bmatrix}$$

$$\text{At } (0, 2) = \begin{bmatrix} 12 & 0 \\ 0 & 12 \end{bmatrix}$$

Both eigenvalues $> 0 \therefore \text{min}$

Both eigenvalues $< 0 \therefore \text{max}$

one pos, one neg $\therefore \text{saddle}$

any 0 \rightarrow inconclusive

Part 2: Vector Optimization

- numerically minimize each function

Fixed stepsize steepest Descent

- update rule: $w_{k+1} = w_k - s \nabla f(w_k)$

Backtracking line search steepest descent

update rule: $w_{k+1} = w_k - s_k \nabla f(w_k)$

↳ start with $s_0 = 0.1$, shrink by $\beta = 0.5$

↳ converges faster

Damped Newtons Method

• updated rule: $w_{k+1} = w_k - s_k H^{-1}(w_k) \nabla f(w_k)$

- uses Hessian curvature info

For each function; must report the minimizer w^* , number of iterations

Part 3: Non Linear Squares

- estimation location of GPS receiver

↳ given: 3D coordinates of 6 satellites

• measured distance between each satellite and receiver

$$d_{i,j} = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2}$$

Content Overview

Vector Calc

- PartialS, gradient, Directional deriv
 $\frac{\partial f}{\partial w}$ ∇f $D_v f = \nabla f \cdot \vec{v}$

- Jacobian matrix, vector valued function

$$J = m \times n$$

$$J_{ij} = \frac{\partial f_i}{\partial w_j}$$

Jacobian gives matrix, best linear approximation to a vector valued function.

$$\hookrightarrow f(w) = \begin{bmatrix} w_1^2 + w_2 \\ \sin(w_1 w_2) \end{bmatrix} \quad \text{Jacobian} = \begin{bmatrix} 2w_1 & 1 \\ w_2 \cos(w_1 w_2) & w_1 \cos(w_1 w_2) \end{bmatrix}$$

→ can do mappings and plug ins.

- Hessian $\nabla^2 f$, Second partials, eigenvalues = convexity

Steepest Descent

- Desc direction, steepest descent $= - [\nabla f(w)]^T$
- Fixed stepsize $w_{k+1} = \vec{w}_k + \alpha \vec{d}$
- Back tracking

Newtons Method

- HELPS Find minimum by using second order information (∇ grad, $\nabla^2 f$ Hessian)
- must be pos definite
- Damped newtons method, adds scaling factor, (newton steps too aggressive)
- Iterate until convergence

Non linear Least Squares

- useful when model is nonlinear in params
- minimizes sum of squared residuals
- Reduce $\frac{1}{2} \|\vec{r}(\vec{w})\|^2$

1) residual vector

$$\vec{r}(\vec{w}) = \begin{bmatrix} r_1(\vec{w}) \\ r_2(\vec{w}) \\ \vdots \\ r_m(\vec{w}) \end{bmatrix} \quad \text{model output} - \text{actual value}$$

2) linearized residual, helps linearize

$$a_i(\vec{w}) = r_i(\vec{w}_0) + \nabla r_i(\vec{w}_0)^T (\vec{w} - \vec{w}_0)$$

3) Linearized objective function:

$$f(\vec{w}) = \frac{1}{2} \|\vec{a}(\vec{w})\|^2$$

iteratively update.

2 Main NLS Algorithms



Gauss Newton

$$w_{k+1} = w_k + [J_k^T J_k]^{-1} J_k^T \vec{r}(w_k)$$

J_k is the jacobian matrix at point w_k

$$\text{Hessian} = J_k^T J_k$$

Levenberg Marquadt method

$$w_{k+1} = w_k + [J_k^T J_k + \tau I]^{-1} J_k^T \vec{r}(w_k)$$

Adds damping factor $\tau \geq 0$
to stabilize updates
to prevent 0 eigenval problem

Example

$$r_1(w) = w^2 - 2, \quad r_2(w) = \sin(w) - 0.5$$

$$\text{objective function, } f(w) = \frac{1}{2} (r_1(w)^2 + r_2(w)^2)$$

$$\text{residual vector } \vec{r}(w) = \begin{bmatrix} w^2 - 2 \\ \sin(w) - 0.5 \end{bmatrix}$$

$$w_1 = w_0 + (J^T J)^{-1} J^T \vec{r}$$

$$w_1 = w_0 + \begin{pmatrix} -1.8154 \\ 4.2929 \end{pmatrix}$$

$$w_1 = 0.5771$$

$$\text{Jacobian} = \begin{bmatrix} 2w \\ \cos w \end{bmatrix}$$

$$J(w_0) = \begin{bmatrix} 2.0 \\ 0.5403 \end{bmatrix}$$

$$r(1.0) = \begin{bmatrix} 1.0^2 - 2 \\ \sin(1.0) - 0.5 \end{bmatrix}$$

$$= \begin{bmatrix} -1.0 \\ 0.3415 \end{bmatrix}$$

Level curve:

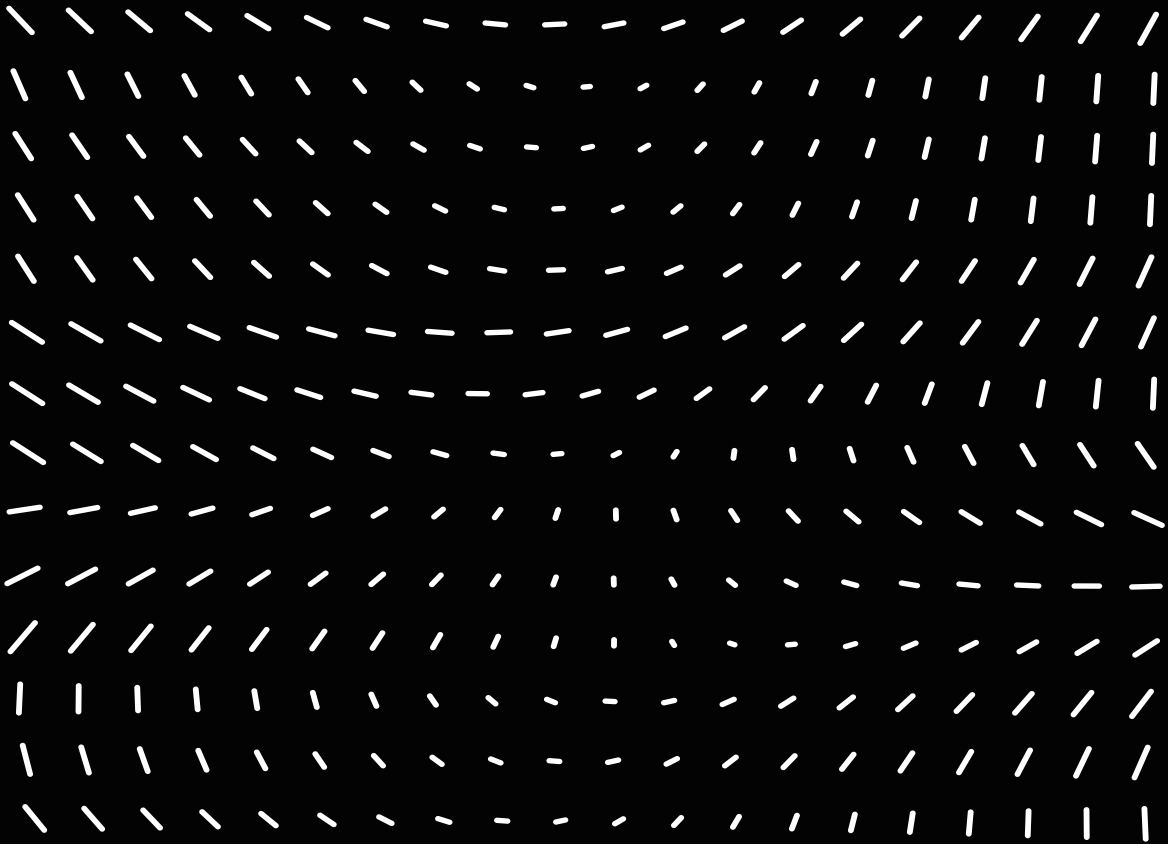
- Multiple curves, contour plot

$$z = f(x, y)$$

- Concentric circles

Test 2

Prep



Test 2 Prep

Vector Calc

$$f(w_1, w_2) = 3w_1^2 + 2w_1w_2 + w_2^3$$

$$\frac{\partial f}{\partial w_1} = 6w_1 + 2w_2$$

$$\frac{\partial f}{\partial w_2} = 2w_1 + 3w_2^2$$

→ Plug in (1, 2)

$$= 10, 14$$

Directional Deriv

$$DD = \nabla f \cdot v$$

$$f_2(w) = 3(w_2 + 1)^2 + 2w_1$$

Plug in (1, 1)

$$= (12, 2) \cdot (-1, -1) = -14$$

Gradient 1 form

$$f(w_1, w_2) = w_1^2 + 5w_1w_2 + 4w_2^2$$

$$\frac{\partial f}{\partial w_1} = 2w_1 + 5w_2$$

$$\frac{\partial f}{\partial w_2} = 5w_1 + 8w_2$$

$$-\nabla f \rightarrow \text{steepest Descent}$$

Jacobian

$$\begin{bmatrix} \frac{\partial f_1}{\partial w_1} & \frac{\partial f_1}{\partial w_2} \\ \frac{\partial f_2}{\partial w_1} & \frac{\partial f_2}{\partial w_2} \end{bmatrix}$$

Jacobian contains derivatives of vector valued function

Each row is gradient of residual

$$f_1(w) = w_1^2 + w_2$$

$$f_2(w) = \sin(w_1 w_2)$$

$$J_i(w) = \|a_i - w\|$$

For each row $J_i = \frac{a_i - w^T}{\|a_i - w^T\|}$

$$\frac{\partial f_1}{\partial w_1} = 2w_1$$

$$\frac{\partial f_1}{\partial w_2} = 1$$

$$\frac{\partial f_2}{\partial w_1} = \cos(w_1 w_2) \cdot w_2$$

$$\frac{\partial f_2}{\partial w_2} = \cos(w_1 w_2) \cdot w_1$$

$$\text{At } w = (1, \pi) = \begin{bmatrix} 2 & 1 \\ -\pi & -1 \end{bmatrix}$$

Each residual: $\|a_i - w\|$

$$\frac{\text{direction vector}}{\text{length of vector}}$$

$$A = [3, 12, 1]$$

$$b_0 = 2$$

$$J(b) = \begin{bmatrix} \frac{1}{1} & \frac{10}{10} & \frac{-1}{1} \\ 1 & 1 & -1 \end{bmatrix}$$

$$J_i = \frac{a_i - w}{\|a_i - w\|}$$

2D Example

$$= a_i - w$$

$$= [4 - 1, 3 - 1] = [3, 2]$$

$$J_i = \frac{[3, 2]}{\sqrt{13}}$$

$$J_i = \frac{a_i - b_0}{\|a_i - b_0\|}$$

One D

Gauss-Newton

residual $r(w) = [w - 2, w - 3]$

$$J(w) = [1, 1]$$

From $w_0 = 0$

$$w_{k+1} = w_k + [J^T J]^{-1} J^T r(w_k)$$

$$r = [-2, -3]$$

$$0 + [2]^{-1} (-2) + (-3)$$

$$= 0 + \frac{-2}{2} = -2.5$$

3 by 2 · 2 by 5
valid
Resultant

Hessian

$$f(w) = w_1^2 + w_2^4$$

$$e [1, 0]$$

$$H = 2w_1 + 4w_2^3$$

$$= 2, 12w_2^2$$

$$\begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix}$$

Hessian is not positive definite
∴ Newton's method can't be applied

Hessian

$$f(w_1, w_2) = w_1^2 + 4w_1w_2 + 4w_2^2$$

$$H = \nabla^2 f$$

$$\frac{\partial^2 f}{\partial w_1^2} = 2w_1 + 4w_2$$

$$\frac{\partial^2 f}{\partial w_2^2} = 4w_1 + 8w_2$$

$$\begin{bmatrix} 2 & 4 \\ 4 & 8 \end{bmatrix}$$

min stability
∴ convex

$$\nabla f = 0$$

$$0 = 6w_1 + 2w_2$$

$$2w_2 + 2w_2$$

$$(0, 0)$$

$$f(w_1, w_2) = \log(1 + w_1^2) + e^{w_2}$$

Convexity

$$\log(1 + w_1^2)$$

is convex from second deriv test

Optimization For Steepest Descent

- minimizing function by moving in direction of negative gradient. $\vec{J} = -\nabla F(w)$

1) $\vec{w}, x, b =$

Steepest desc vector: $\vec{J} = -b \cdot \vec{x} \cdot \vec{w}$

$$\vec{J} = -b \cdot \text{avg } x$$

$$A = \begin{bmatrix} 2 \\ 7 \\ 10 \end{bmatrix}, \quad t_0 = 5, \quad \tau = 3$$

$$r = \begin{bmatrix} 2 & -5 \\ 7 & -5 \\ 10 & -5 \end{bmatrix} = \begin{bmatrix} -3 \\ 2 \\ 5 \end{bmatrix}, \quad \tau = 3$$

$$J = \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix} = 5 + [3 + 3]^{-1} \\ 5 + \frac{1}{6} \cdot 10 \\ = 6.6$$

Levenberg Example

5) $A = \begin{bmatrix} 4 \\ 8 \end{bmatrix}, \quad t_0 = 2, \quad \tau = 1$

$$b_1 = b_0 + (J^T J + \tau I)^{-1} J^T r$$

1) Compute r : $\begin{bmatrix} 1 & -2 \\ 4 & -2 \\ 8 & -2 \end{bmatrix} = \begin{bmatrix} -1 \\ 2 \\ 6 \end{bmatrix}$

2) Jacobian: $\frac{a_i - t}{|a_i - t|} = \left[\frac{-1}{1}, \frac{1}{1}, \frac{1}{1} \right]$

$$J = \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix} \\ J^T J = 3$$

$$J^T r = \begin{bmatrix} -1 \\ 2 \\ 6 \end{bmatrix} \begin{bmatrix} -1 & 1 & 1 \end{bmatrix} \\ = 9$$

3)

$$= t_0 + \frac{J^T r}{J^T J + \tau}$$

$$= 2 + \frac{9}{3+1}$$

$$= 3.75$$

$$A = \begin{bmatrix} 0 \\ 5 \\ 6 \end{bmatrix}, \quad t_0 = 3, \quad \tau = 2$$

$$w_{k+1} = 3 + [J^T J + \tau]^{-1} J^T r$$

$$r = \begin{bmatrix} 0 & -3 \\ 5 & -3 \\ 6 & -3 \end{bmatrix} \quad r = \begin{bmatrix} -3 \\ 2 \\ 3 \end{bmatrix}$$

$$J = \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix} \quad \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} -1 & 1 & 1 \end{bmatrix}$$

$$= 3 + [3+2]^{-1} \begin{bmatrix} -1 & 1 & 1 \end{bmatrix} \begin{bmatrix} -3 \\ 2 \\ 3 \end{bmatrix} \\ = 3 + \frac{1}{5} \cdot 8 = 4.6$$

NLS

$$r(w) = [r_1(w), r_2(w), \dots, r_m(w)]$$

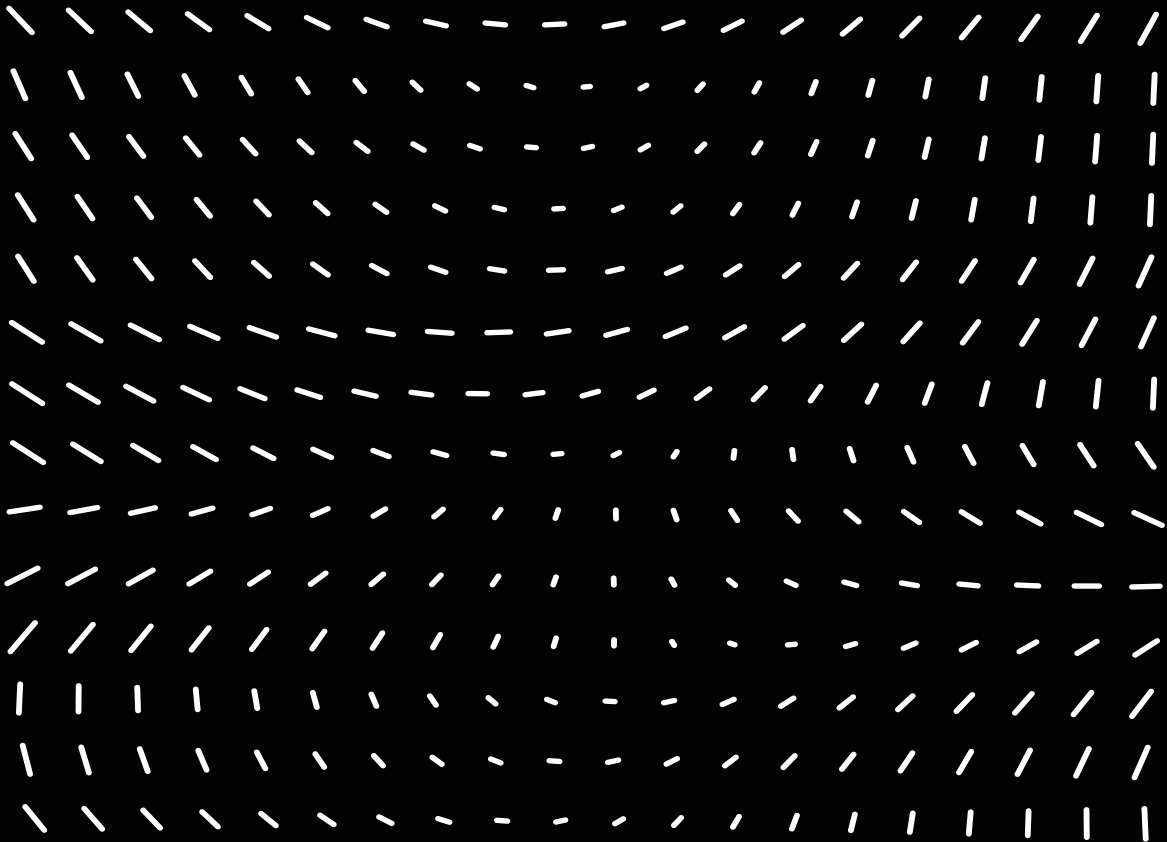
$$F(w) = \frac{1}{2} \|r(w)\|^2$$

Descent Direction: $d = -\nabla f(w_0)$

$$w_{k+1} = w_k + \alpha_k d_k$$

Newton's Method: $d_k = -H^{-1} \nabla f(w_k)$

Quiz 3/ Test Prep



Unit 3 - Important Concepts

Single Artificial Neuron

• Smallest building block of Neural net

• Takes input vector \rightarrow Applies linear function \rightarrow Passes through activation function

\rightarrow Produces output score

• Data vector is transpose of observation

• Neurons include bias term. $[x \ 1]$

• For Classification, labels are either 0 or 1, because of logistic activation

• Linear response is $u(\vec{w}) = [x \ 1] w$

Ex

$$u = x \cdot w \text{ given } \vec{w} = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \quad A = \begin{bmatrix} 2 & 3 \\ 0 & 2 \end{bmatrix}$$

$$\begin{bmatrix} 2 & 3 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \rightarrow \vec{v} = \begin{bmatrix} 0.00 \\ -1.00 \end{bmatrix}$$

• Activation function is ϕ

$$\phi(u) = \frac{1}{1+e^{-u}} \rightarrow \text{Logistic}$$

Output is between 0 and 1

• ReLU activation

$$\phi_R(u) = \max(0, u)$$

for hidden layers

Outputs 0 if input is negative

Outputs u if input is positive

• derivative of activation function

$$\psi(u) = \phi'(u)$$

$$\phi'(u) = \phi(u)(1 - \phi(u))$$

$$\text{Ex } \phi(u) = 0.4, \quad \psi = 0.4(1 - 0.4) = 0.240$$

• residual error: $r(\vec{w}) = y - \phi(\vec{w})$

(actual - predicted)

• Objective $f(w) = \frac{1}{2} r(w)^2$

Always reduce residuals

• Learning rate η is like weight updated:

$$\vec{w}_{\text{new}} = \vec{w}_{\text{old}} + \eta \vec{d}$$

Neural Networks with hidden layers

Layer 1 - Input

$$x = [x_1, x_2]$$

$$[x] = [x_1, x_2]$$

Layer 2: Hidden Layer (new part)

multiple neurons \rightarrow weight matrix

2 hidden neurons \rightarrow 2 weight columns

Layer 3 - Output neuron

receives input from hidden layer, weight vector. $w^{(3)}$

$$z^{(3)} = [z_1^{(2)} \quad z_2^{(2)} \quad 1]$$

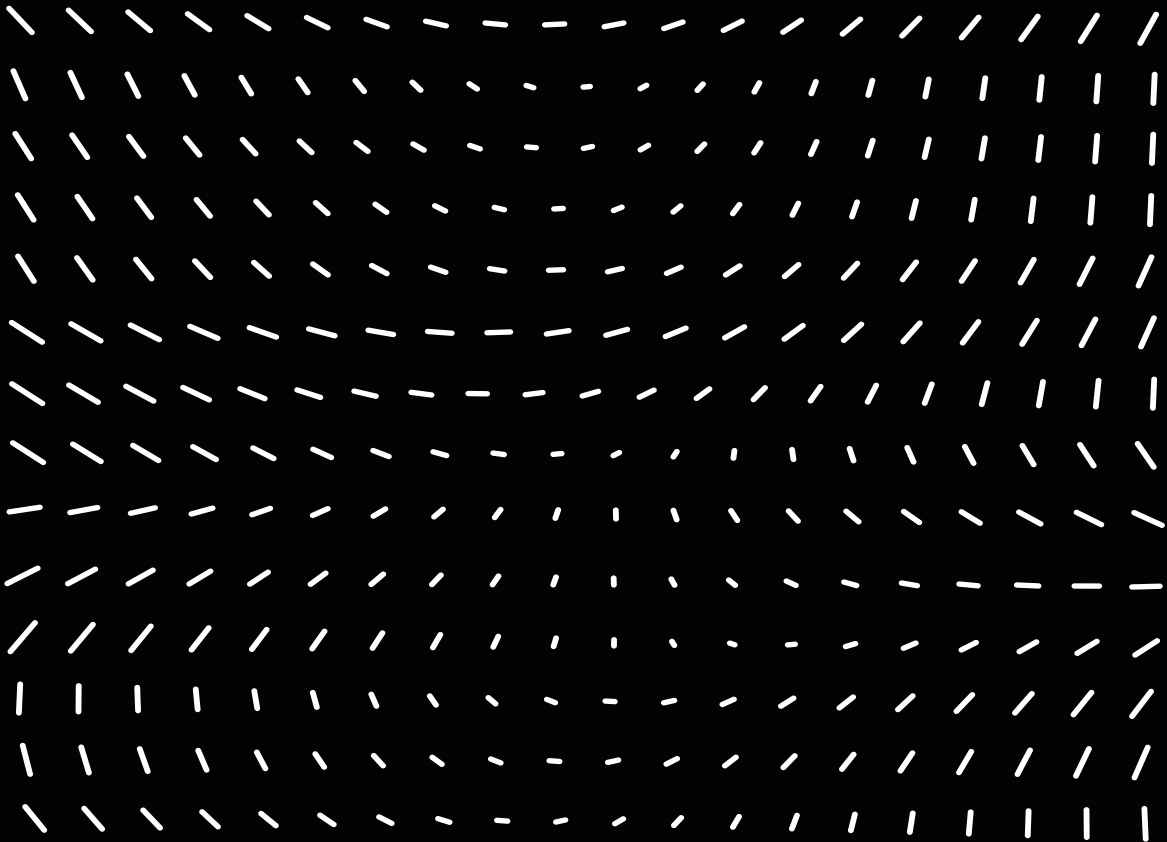
• back Prop is the same

$$b = -f'$$

• just extends single neurons to vectors/matrices

Quiz 4

Prep



Key Topics for Quiz 4

Lagrange Multipliers and Constrained Optimization

- goal: solve constrained optimization problem to minimize function subject to constraints
- if you have constraints like $g(x) = 0$ or $h(x) \leq 0$ can't take deriv, set deriv to 0 like unconstrained problems
- combine objective and constraints into 1 expression
This is the Lagrangian and uses Lagrange multipliers

Equality Constraint

Ex Lagrangian: $L(x, y, \nu) = f(x, y) + \nu \cdot g(x, y)$

Minimize $f(x, y)$

subject to $g(x, y) = 0$

Solve $\nabla_{x,y} L = 0$ and $g(x, y) = 0$

↓
Lagrange multiplier \rightarrow how sensitive the optimal solution is to changes in constraint

ν : large means constraint heavily affects solution

Ex Inequality Constraint

minimize $f(x)$

subject to $h(x) \leq 0$

Lagrangian: $L(x, \lambda) = f(x) + \lambda \cdot h(x)$

But: we require $\lambda \geq 0$, $\lambda \cdot h(x) = 0$

Not all inequality constraints are active at solution

$\lambda \cdot h(x) = 0$: condition forces 2 cases

Case 1: constraint is active

its "tight" $x - 5 \leq 0$, $x = 5$

λ can be non zero

Case 2: constraint is inactive

strictly within bounds

$x - 5 \leq 0$, solution is $x = 3$

$h(x) = -2 < 0$

Then $\lambda = 0$, because constraint does not affect solution

Example

$$\min f(x) = x^2, \text{ subject to } x \geq 2$$

$$h(x) = 2 - x \leq 0$$

1) Lagrangian:

$$L(x, \pi) = x^2 + \pi(2-x)$$

complementary Slackness

Primal Feasibility

KKT conditions

1.) Stationarity: $\frac{dL}{dx} = 2x - \pi = 0$

From ① $\pi = 2x$

2.) Primal Feasibility: $2 - x \leq 0$

Sub into ①

$$2x(2-x) = 0$$

3.) Dual Feasibility: $\pi \geq 0$

4.) Complementary Slackness $\pi(2-x) = 0$

So, $x = 2$ or $x = 0$

↓
Active Constraint

↓
Constraint Violation

Check $\pi = 2x = 4 \geq 0 \checkmark$

Optimal $x = 2$

Stationarity: critical point w.r.t x is zero

Primal Feasibility: original constraint checker

Dual Feasibility: $\pi \geq 0$

Complementary Slackness: can be active/inactive ($\pi_i \cdot h_i(x) = 0$)

If objective is quadratic: $f(x) = \frac{1}{2} x^T Q x + c^T x$

Constraints are linear equalities: $Ax = b$

KKT Matrix System: $\begin{bmatrix} Q & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} x \\ \mu \end{bmatrix} = \begin{bmatrix} -c \\ b \end{bmatrix}$

- gives optimal x and multipliers

Ex

$$Q = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

$$c = \begin{bmatrix} -2 \\ -5 \end{bmatrix}$$

$$A = \begin{bmatrix} 1 & 1 \end{bmatrix}$$

$$b = 3$$

$$\begin{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} & \begin{bmatrix} 1 \\ 1 \end{bmatrix} \\ \begin{bmatrix} 1 & 1 \end{bmatrix} & 0 \end{bmatrix} \begin{bmatrix} x \\ \mu \end{bmatrix} = \begin{bmatrix} \begin{bmatrix} -2 \\ -5 \end{bmatrix} \\ 3 \end{bmatrix}$$

1) Minimizing $\frac{1}{2} (2x_1^2 + 2x_2^2) - 2x_1 - 5x_2$

2) Subject to $x_1 + x_2 = 3$

$$x^T Q x = \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 4 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 4x_1^2 + 2x_1x_2 + 3x_2^2$$

How to make sense of this

KKT Matrix System

→ Used when Quadratic + Linear before and linear equality constraints

Ex
$$\begin{bmatrix} Q & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} x \\ \nu \end{bmatrix} = \begin{bmatrix} -c \\ b \end{bmatrix}$$

$a = \begin{bmatrix} 2 & 0 & 1 \\ 0 & 2 & 1 \\ 1 & 1 & 0 \end{bmatrix}$
 $A^T = \begin{bmatrix} 1 \\ 1 \\ -2 \end{bmatrix}$
 $A = \begin{bmatrix} 1 & 1 \\ -2 & -5 \end{bmatrix}$
 $c = \begin{bmatrix} -2 \\ -5 \end{bmatrix}$
 $b = 3$

$$\begin{bmatrix} 2 & 0 & 1 \\ 0 & 2 & 1 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \nu \end{bmatrix} = \begin{bmatrix} 2 \\ 5 \\ 3 \end{bmatrix}$$

Lin alg this
↙ 0 v b

$\nu = 2 - 2x_1$

$2x_1 + \nu = 2$

$\nu = 2 - 2x_1$

$2(3-x_1) - 2x_1 = 3$

$2x_2 + (2 - 2x_1) = 5$

$2x_2 + \nu = 5$

$2x_2 + (2 - 2x_1) = 5$

$x_1 = \frac{3}{4}$

$2x_2 - 2x_1 = 3$

$x_1 + x_2 = 3$

$x_2 = 3 - x_1$

$x_2 = \frac{9}{4}$

$x_1 = 3 - x_2$

$2x_2 - 2(3 - x_2) = 3$

$2x_2 - 6 + 2x_2 = 3$

$\frac{4x_2}{4} = \frac{9}{4}$
 $x_2 = \frac{9}{4}$

$\nu = \frac{1}{2}$

Primal Problem - the one you start with

Dual Formulation with Linear Constraints

• write in terms of Lagrange multipliers ν

1.) $f(w), \text{ s.t. } A\vec{w} = \vec{b}$

$L(w, \nu) = f(w) + \nu^T(Aw - b)$

↓
vector of Lagrange multipliers

2. take w.r.t \vec{w} , express as function of ν

$g(\vec{\nu}) = \inf L(w, \nu)$

$g(\nu)$ dual function → always concave
maximizing

Ex

$\min \frac{1}{2} \|w - c\|^2 \text{ s.t. } a^T w = b$

$c = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$
 $a = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$
 $b = 10$

① $L(w, \nu) = \frac{1}{2} \|w - c\|^2 + \nu(a^T w - b)$
 $= \frac{1}{2} (w - c)^T (w - c) + \nu(a^T w - b)$

② w.r.t $w = w - c + \nu a = 0$
 $w = c - \nu a$

w.r.t gradient

$w^T w = 2w$

$a^T w = a$

$Aw = A^T$

$\frac{1}{2} \|w - c\|^2 = w - c$

3) Plug $w = c - \nu a$ into Lagrangian:

Objective: $\frac{1}{2} \nu^2 (25)$

Constraint: $\nu(1 - 25\nu)$

Then find ν
Then get w

4) Simplify →

$g(\nu) = -\frac{25}{2} \nu^2 + \nu$

2 ways to solve quadratic optimization problem with linear equality

constraints:

- 1) Solve with KKT (Primal)
- 2) Derive, solve dual problem

B-Dual matrix. Problem becomes how
to maximize: $-\frac{1}{2} \mu^T B \mu - \mu^T d$

• Always concave

• KKT Conditions must be satisfied to be optimal solution

Basically, given any $f(\vec{w})$ if it has equality or inequality constraints for any optimal points there is $\lambda \geq 0$, μ to satisfy:

- 1) Stationarity \rightarrow finding gradient = 0
- 2) Primal feasibility \rightarrow must satisfy original constraints $g_i(\vec{w}^*) \leq 0$, $h_j(\vec{w}^*) = 0$
- 3) Dual feasibility $\rightarrow \lambda_i$ for all i , Lagrange multipliers for inequality are non negative.
- 4) Complementary slackness \rightarrow if constraint active $\lambda_i \geq 0$
if constraint inactive $\lambda = 0$

KKT Conditions together determine a solution to a specific problem

Constrained Least Squares (CLS) and Tikhonov Regularization

• Minimize $\|A\vec{w} - \vec{b}\|^2$

Control size or penalize large weights, you add constraints or regularization.

Both CLS and Tikhonov, you control $\|\vec{w}\| = w_1^2 + w_2^2 + \dots + w_n^2$, helps avoid overfitting

• Fits training data, sacrifices some accuracy for better smooth curve

• Tikhonovs Regularization: $\min_{\vec{w}} \|A\vec{w} - \vec{b}\|^2 + \lambda \|\vec{w}\|^2$

without regularization: $w = \begin{bmatrix} 1 \\ 1000 \end{bmatrix}$, $w = \begin{bmatrix} 0.99 \\ 10 \end{bmatrix}$
 \downarrow Penalty on large weights

• CLS Problem is like $\min \|A\vec{w} - \vec{b}\|^2$ subject to $\|\vec{w}\|^2 \leq \theta$

• θ max allowed size of w

• $\min \|A\vec{w} - \vec{b}\|^2 + \lambda \|\vec{w}\|^2$

\hookrightarrow Tunable constant
penalizing size of w

Tikhonov Type Questions

1) Given: $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$, $b = \begin{bmatrix} 5 \\ 6 \end{bmatrix}$, $\lambda = 1$

Find minimizer w^* for $\|Aw - b\|^2 + \lambda \|w\|^2$

$$w^* = (A^T A + \lambda I)^{-1} A^T b$$

$$\text{Tikhonov: } \|Aw - b\|^2 + \lambda \|w\|^2$$

$$(A^T A + \lambda I) = \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix} = \begin{bmatrix} 10 & 14 \\ 14 & 20 \end{bmatrix} = \begin{bmatrix} 10 & 14 \\ 14 & 20 \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 11 & 14 \\ 14 & 21 \end{bmatrix}^{-1} \rightarrow \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix} = \frac{1}{25} \begin{bmatrix} 21 & -14 \\ -14 & 11 \end{bmatrix}$$

Larger λ
Shrinks w more

$$A^T b = \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix} \begin{bmatrix} 5 \\ 6 \end{bmatrix} = \begin{bmatrix} 23 \\ 34 \end{bmatrix} = \begin{bmatrix} \frac{1}{5} \\ \frac{34}{5} \end{bmatrix}$$

2.) CLS Behaviour

$$\min \|Aw - b\|^2 \text{ subject to } \|w\|^2 \leq \theta$$

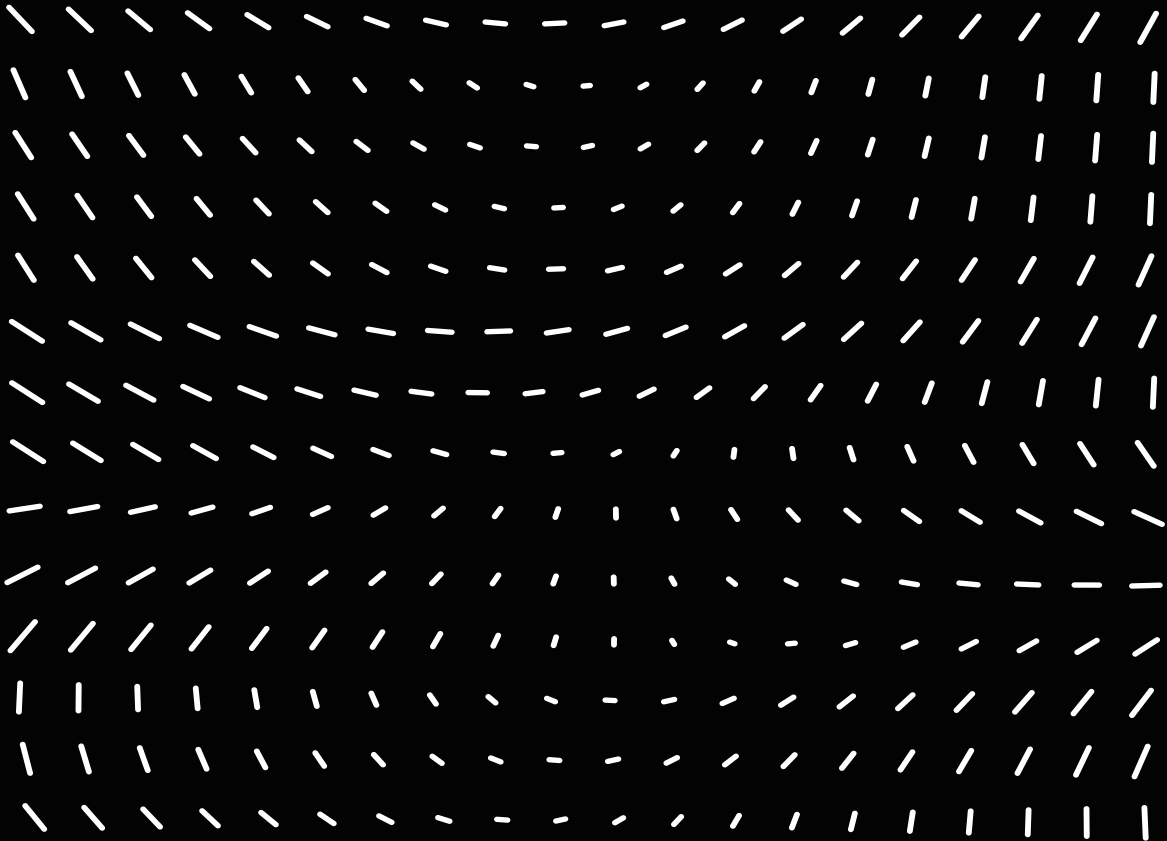
$$\|w\|^2 = 25$$

$$\theta = 10$$

$$\|w^*\|^2 = 25 > \theta = 10 \text{ violates constraint}$$

Test 4

Prep



2024 - Test 4

2) Objective function and non linear equality constraint:

$$f(\vec{w}) = (w_1 - 1)^2 + (w_2 - 1)^2$$

$$p(\vec{w}) = w_1^2 - w_2^2 = 0$$

Find Lagrange function

$$L(x, y, \nu) = f(x, y) + \nu \cdot g(x, y)$$

$$= (w_1 - 1)^2 - (w_2 - 1)^2 + \nu (w_1^2 - w_2^2)$$

3) Objective function: $f(w) = \frac{1}{2} \vec{w}^T K \vec{w}$ and linear equality constraint $m \vec{w} - c = 0$, where

$$K = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \quad m = [1 \quad 1] \quad c = 2, \text{ find minimizer that satisfies constraint}$$

Note: no linear term in $\frac{1}{2} \vec{w}^T K \vec{w}$
so, means $c = 0$ in general form

Note: RHS of KKT
Linear $\rightarrow \frac{1}{2} x^T Q x \rightarrow \begin{bmatrix} b \\ 0 \end{bmatrix}$
Quadratic $\rightarrow \frac{1}{2} x^T Q x + c^T x \rightarrow \begin{bmatrix} b \\ -c \end{bmatrix}$

1) Setup KKT matrix :

$$\begin{bmatrix} K & m^T \\ m & 0 \end{bmatrix} \begin{bmatrix} w \\ \nu \end{bmatrix} = \begin{bmatrix} 0 \\ c \end{bmatrix}$$

$$\begin{bmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \nu \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 2 \end{bmatrix}$$

2) Setup linear systems:

$$\textcircled{1} \quad 2w_1 + w_2 + \nu = 0$$

$$\textcircled{2} \quad w_1 + 2w_2 + \nu = 0$$

$$\textcircled{3} \quad w_1 + w_2 = 2$$

$$\therefore w^* = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

3) solve:

$$\textcircled{3} \quad w_2 = 2 - w_1$$

sub in $\textcircled{1}$

$$2w_1 + (2 - w_1) + \nu = 0$$

$$w_1 + 2 + \nu = 0$$

$$\nu = -w_1 - 2$$

$$w_1 + 2(2 - w_1) + (-w_1 - 2) = 0$$

$$w_1 + 4 - 2w_1 - w_1 - 2 = 0$$

$$-2w_1 + 2 = 0$$

$$w_1 = 1$$

$$w_2 = 2 - 1 \\ w_2 = 1$$

$$\begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\nu = -1 - 2$$

$$\nu = -3$$

Q4) Objective $f(\vec{w}) = \frac{1}{2} w^T K \vec{w} + q^T w$ and linear equality constraint $m \vec{w} - c = 0$

$$K = \begin{bmatrix} 8 & 0 \\ 0 & 2 \end{bmatrix} \quad m = [1 \ 1] \quad q = \begin{bmatrix} -4 \\ -2 \end{bmatrix} \quad c = 2$$

Find minimizer w^* and dual matrix β

1) Setup KKT matrix: $\begin{bmatrix} K & m^T \\ m & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \nu \end{bmatrix} = \begin{bmatrix} -q \\ c \end{bmatrix}$

$$\begin{bmatrix} 8 & 0 & 1 \\ 0 & 2 & 1 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \nu \end{bmatrix} = \begin{bmatrix} 4 \\ 2 \\ 2 \end{bmatrix}$$

WATCH OUT
FOR $-q$

$$A^{-1} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

2) Setup and solve for ν

$$\left. \begin{aligned} 8w_1 + \nu &= 4 \\ 2w_2 + \nu &= 2 \\ w_1 + w_2 &= 2 \end{aligned} \right\}$$

$$w_1 = 2 - w_2$$

$$\nu = 2 - 2w_2$$

$$\frac{1}{16-0} \begin{bmatrix} 1 & 1 \\ -1 & 0 \end{bmatrix}$$

$$8(2 - w_2) + \nu = 4$$

$$8(2 - w_2) + (2 - 2w_2) = 4$$

$$16 - 8w_2 + 2 - 2w_2 = 4$$

$$-10w_2 = 4 - 18$$

$$-10w_2 = -14$$

$$w_2 = \frac{14}{10}$$

$$w_1 = 2 - \frac{14}{10}$$

$$w_1 = \frac{3}{5}$$

$$\nu = 2 - 2w_2$$

$$\nu = 2 - 2\left(\frac{14}{10}\right)$$

$$\nu = -0.8$$

$$\nu = 2 - 2\left(\frac{7}{5}\right)$$

$$\nu =$$

Finding Dual Matrix

Use $B = mK^{-1}m^T$

$$B = (mK^{-1}m^T) \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

$$B = \left[[1 \ 1] \begin{bmatrix} \frac{1}{8} & 0 \\ 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right]$$

$$B = \left[\begin{bmatrix} \frac{1}{8} & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right] = \begin{bmatrix} \frac{5}{8} \end{bmatrix} \rightarrow \frac{5}{8} ?$$

$$B = 0.625$$

5.) Convex Objective: $f(t) = t^2 + 4t$, nonlinear inequality constraint $p_1(t) = t^2 - 6t + 8$ with $p_1(t) \leq 0$ find minimizer t^*

1) Lagrangian: $f(t) + \nu \cdot p_1(t)$
 $(t^2 + 4t) + \nu(t^2 - 6t + 8)$

1) minimize $f(t)$:

$$f(t) = t^2 + 4t$$

$$f'(t) = 2t + 4$$

$$= 2(t + 2)$$

$$t = -2$$

$$(-2)^2 - 6(-2) + 8$$

$$4 + 12 + 8 \quad \times$$

Constraint

violation

how what?

For a convex function on closed interval the minimum should occur at

• end points

• or stationary point (all tried here)

2) Optimize over feasible region,

$$p_1(t) = t^2 - 6t + 8$$

$$\text{Factor: } (t - 4)(t - 2)$$

$$t = 4 \quad t = 2$$

3) Minimize on this interval $t \in [2, 4]$

$$f(2) = 12$$

$$f(4) = 32$$

$$\therefore \min_{t \in [2, 4]} f(t) = 12$$

This is inspection method: only use Lagrange if you know constraints active

Q6) Convex objective function $f(\vec{w}) = (w_1-3)^2 + (w_2-4)^2$, nonlinear inequality constraint
 $p_1(\vec{w}) = w_1^2 + w_2^2 - 4$ with $p_1(\vec{w}) \leq 0$, find minimizer \vec{w}^* .

1) Intuition: The constraint is circle of radius 2, centered at the origin.

Objective: convex quadratic point (3,4)

Trying to minimize distance from (w_1, w_2) to (3,4)

2) Try unconstrained: $f(w_1, w_2) = (w_1-3)^2 + (w_2-4)^2$
 $w_1=3 \quad w_2=4$

Plug into p_1 : $p_1(3,4) = 21 \leq 0$ X

3) Does not work, must be Lagrange:

$$f(w_1, w_2) + \mu(p_1(w_1, w_2))$$

$$(w_1-3)^2 + (w_2-4)^2 + \mu(w_1^2 + w_2^2 - 4)$$

4) KKT $w_1: 2(w_1-3) + \mu(2w_1)$

$w_2: 2(w_2-4) + \mu(2w_2)$

Stationary:

w.r.t w_1, w_2, μ

$\lambda: w_1^2 + w_2^2 - 4$

$$2w_1 - 6 + 2\mu w_1 = 0$$

$$2w_1 + 2\mu w_1 = 6$$

$$(w_1-3) + \mu w_1 = 0$$

$$\frac{w_1-3}{w_1} = (-\mu) -$$

$$\frac{3-w_1}{w_1} = \mu$$

5) Set all equal to 0:

$$\begin{cases} 2(w_1-3) + 2\mu w_1 = 0 \rightarrow \mu = \frac{3-w_1}{w_1} \\ 2(w_2-4) + 2\mu w_2 = 0 \rightarrow \mu = \frac{4-w_2}{w_2} \\ w_1^2 + w_2^2 = 4 \end{cases}$$

6) Solve :

$$\frac{3-w_1}{w_1} = \frac{4-w_2}{w_2}$$

$$w_2(3-w_1) = w_1(4-w_2)$$

$$3w_2 - w_1w_2 = 4w_1 - w_1w_2$$

$$3w_2 = 4w_1$$

$$w_2 = \frac{4}{3}w_1$$

$$w_1^2 + \left(\frac{4}{3}w_1\right)^2 = 4$$

$$w_1^2 + \frac{16}{9}w_1^2 = 4$$

$$w_1^2 \left(1 + \frac{16}{9}\right) = 4$$

$$w_1^2 \left(\frac{25}{9}\right) = 4$$

$$w_1^2 = \frac{4}{\frac{25}{9}}$$

$$w_1^2 = \frac{36}{25}$$

$$w_1 = \pm \sqrt{\frac{36}{25}}$$

$$w_1 = \pm 1.2$$

$$w_2 = \frac{4}{3}w_1, \pm \frac{4}{3} \cdot 1.2$$

$$w_2 = \pm 1.6$$

7) Possible points:

(1.2, 1.6)

(-1.2, -1.6)

8) Evaluate

$$f(1.2, 1.6) = 9$$

$$f(-1.2, -1.6) = 49$$

$$\therefore \vec{w}^* = \begin{bmatrix} 1.2 \\ 1.6 \end{bmatrix}$$

$\nabla f = -\mu$ in vector form

Test 5 - Breakdown

1) Convex objective function $f(w) = [w-g]^T [\vec{w}-\vec{g}]$, where $\vec{g} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$

linear inequality constraints

$A = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 1 \end{bmatrix}$, $b = \begin{bmatrix} 0 \\ 0 \\ 4 \end{bmatrix}$ find number of ineq constraints

$A\vec{w} \leq \vec{b}$

1) $f(w) = \|w-g\|^2$

minimized at $w=g = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$, if no constraints

Sub g as \vec{w}

2) check if $\vec{g} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$ satisfies constraints

$A\vec{w} = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \end{bmatrix} = \begin{bmatrix} -1 \\ -3 \\ 4 \end{bmatrix}$

$Ag = b?$

Compare to
B

$\begin{bmatrix} 0 \\ 0 \\ 4 \end{bmatrix}$

$-1 \leq 0$

$-3 \leq 0$

$4 \leq 4$

\therefore One constraint is active

1) $f(w) = (w_1-1)^2 + (w_2-1)^2$

$p(w) = w_1^2 - w_2^2$

$L(w, \nu) = (w_1-1)^2 + (w_2-1)^2 + \nu(w_1^2 - w_2^2)$

4) $f(t) = t^2 + 4t$, $p(t) = t^2 - 6t + 9$
with $p(t) \leq 0$

1) check if uncon works

$f'(t) = 2t + 4 = 2(t+2)$
 $t = -2$

2) check condition: $p(-2) = -2^2 - 6(-2) + 9 = 4 + 12 + 9 \leq 0$ X

3) must lie on bounds, find crits on P

4) $P(t) = t^2 - 6t + 8$ check both endpoints
 $0 = (t-4)(t-2)$
 $t = 4$ $t = 2$
 $f(4) =$
 $p(2) =$

check if this works

\therefore min @ 2

2) $f(w) = \frac{1}{2} w^T K w$, lin equality constraint $m\vec{w} - c = 0$

$\begin{bmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \nu \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 2 \end{bmatrix}$

$w_1 = 1$ \therefore minimizer $w^* = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$w_2 = 1$
 $\nu = -3$

3) $\begin{bmatrix} 8 & 0 & 1 \\ 0 & 2 & 1 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \nu \end{bmatrix} = \begin{bmatrix} 4 \\ 2 \\ 2 \end{bmatrix}$ $\frac{1}{16} \begin{bmatrix} \frac{1}{8} & 0 \\ 0 & \frac{1}{2} \end{bmatrix}$

$w_1 = \frac{3}{8}$
 $w_2 = \frac{7}{8}$
 $\nu = -\frac{4}{5}$

$B = M K^{-1} M^T$
 $B = \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{8} & 0 \\ 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$\nu = -0.8$ $B = \begin{bmatrix} \frac{1}{8} & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$B = \frac{5}{8}$

5) $f(w) = (w_1-3)^2 + (w_2-4)^2$
 $p(w) = w_1^2 + w_2^2 - 4$, $p(w) \leq 0$, find min

1) try $w_1 = 3, w_2 = 4$

$p(3,4) = 21 \leq 0$ X

2) setup lagrange:

Homework 9 - Practise

Linear equality constraints:

$$M\vec{w} = c \quad \text{where}$$

$$M = \begin{bmatrix} 1 & 1 \end{bmatrix} \quad \text{and} \quad c = 4$$

Quadratic Objective with Linear Constraints

1. Find primal Lagrange equation
2. Solve KKT matrix to find w^* and ν^*
3. formulate DUAL
4. dual formulation to find ν^* then w^*

Task #1

$$f(w) = \frac{1}{2} w^T K_1 \vec{w} + a_1^T \vec{w}, \quad \text{where } K_1 = \begin{bmatrix} 18 & 8 \\ 8 & 10 \end{bmatrix}, \quad g_1 = \begin{bmatrix} 5 \\ 2 \end{bmatrix}, \quad \vec{c}_1 = \begin{bmatrix} -106 \\ -60 \end{bmatrix}$$

$$\begin{aligned} \text{i) } L(\vec{w}, \nu) &= \frac{1}{2} w^T K_1 \vec{w} + a_1^T \vec{w} + \nu (M\vec{w} - c) \\ &= \frac{1}{2} [w_1, w_2] \begin{bmatrix} 18 & 8 \\ 8 & 10 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + \begin{bmatrix} -106 \\ -60 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + \nu \left(\begin{bmatrix} 1 \\ 1 \end{bmatrix} [w_1, w_2] - 4 \right) \\ &= \frac{1}{2} \begin{bmatrix} 18w_1 + 8w_2 \\ 8w_1 + 10w_2 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \end{aligned}$$

$$= \frac{1}{2} \begin{bmatrix} 18w_1^2 + 8w_1w_2 \\ 8w_1w_2 + 10w_2^2 \end{bmatrix}$$

$$= \frac{1}{2} (18w_1^2 + 16w_1w_2 + 10w_2^2) - 106w_1 - 60w_2 + \nu(w_1 + w_2 - 4)$$

Do this way, matrix form, omits constants

$$\text{ii) } \begin{bmatrix} K & M^T \\ m & 0 \end{bmatrix} \begin{bmatrix} w \\ \nu \end{bmatrix} = \begin{bmatrix} -a_1 \\ c \end{bmatrix}$$

$$\begin{bmatrix} 18 & 8 & 1 \\ 8 & 10 & 1 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \nu \end{bmatrix} = \begin{bmatrix} 106 \\ 60 \\ 4 \end{bmatrix}$$

$$18w_1 + 8w_2 + \nu = 106$$

$$8w_1 + 10w_2 + \nu = 60$$

$$w_1 + w_2 = 4$$

$$w_1 = 4 - w_2$$

$$\nu = 60 - 8w_1 - 10w_2$$

$$\nu = 60 - 8(4 - w_2) - 10w_2$$

$$18(4 - w_2) + 8w_2 + 60 - 8(4 - w_2) - 10w_2 = 106$$

$$72 - 18w_2 + 8w_2 + 60 - 32 + 8w_2 - 10w_2 = 106$$

$$-12w_2 = 6$$

$$w_2 = -\frac{1}{2}$$

$$w_1 + (-\frac{1}{2}) = 4$$

$$w_1 = 4 + \frac{1}{2}$$

$$w_1 = \frac{9}{2}$$

$$8(\frac{9}{2}) + 10(-\frac{1}{2}) + \nu = 60$$

$$31 + \nu = 60$$

$$\nu = 29$$

$$w^* = \begin{bmatrix} 4.5 \\ -0.5 \end{bmatrix} \quad \nu^* = 29$$

$$\text{iii) } \text{dual: } g(\nu) = \min L(w, \nu)$$

$$g(\nu) = f(w(\nu)) + \nu(Mw(\nu) - c)$$

$$L(w, \nu) = \frac{1}{2} (w - g)^T K (w - g) + \nu (Mw - c)$$

$$h(w - g) = -\nu M^T$$

$$w(\nu) = g - \nu K^{-1} M^T$$

$$w(\nu) = \begin{bmatrix} 5 \\ 2 \end{bmatrix} - \nu \begin{bmatrix} \frac{1}{58} \\ \frac{5}{58} \end{bmatrix}$$

$$\begin{aligned} w_1 + w_2 &= 7 - \frac{6\nu}{58} - 4 \\ &= 3 - \frac{3\nu}{29} \end{aligned}$$

$$\text{Dual: } \frac{3\nu(58 - \nu)}{58}$$

↳ Derive, set to 0, find ν then w_1, w_2 .

HW-9 - Task 2

$$m = [1 \ 1] \quad \text{and} \quad c = 4$$

$$\begin{aligned} 2) \quad i) \quad f(w) &= \frac{1}{2} w^T K w + c^T w \\ &= \frac{1}{2} [w_1 \ w_2] \begin{bmatrix} 8 & 2 \\ 2 & 12 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + [-34 \ 26] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \\ L &= \frac{1}{2} (8w_1^2 + 14w_1w_2 + 12w_2^2 - 34w_1 + 26w_2) + \mu (w_1 + w_2 - 4) \end{aligned}$$

$$ii) \quad \begin{bmatrix} K & m^T \\ m & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \mu \end{bmatrix} = \begin{bmatrix} 34 \\ -26 \\ 4 \end{bmatrix}$$

$$\begin{bmatrix} 8 & 2 & 1 \\ 2 & 12 & 1 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \mu \end{bmatrix} = \begin{bmatrix} 34 \\ -26 \\ 4 \end{bmatrix}$$

$$8w_1 + 2w_2 + \mu = 34$$

$$2w_1 + 12w_2 + \mu = -26$$

$$w_1 + w_2 = 4$$

$$w_1 = 4 - w_2$$

$$\mu = 34 - 8(4 - w_2) - 2w_2$$

$$\mu = 34 - 32 + 8w_2 - 2w_2$$

$$\mu = 2 + 6w_2$$

$$2(4 - w_2) + 12w_2 + 2 + 6w_2 = -26$$

$$8 - 2w_2 + 12w_2 + 2 + 6w_2 = -26$$

$$10 + 16w_2 = -26$$

$$\frac{16w_2}{16} = \frac{-36}{16}$$

$$w_2 = \frac{-9}{4}$$

$$w_1 = 4 + \frac{9}{4}$$

$$w_1 = \frac{25}{4}$$

$$2\left(\frac{25}{4}\right) + 12\left(\frac{-9}{4}\right) + \mu = -26$$

$$\mu = -11.5$$

could use Calc: Ea → 3 variable

$$3) \quad B = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 8 & 2 \\ 2 & 12 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\frac{1}{ad-bc} \begin{bmatrix} d & -c \\ -b & a \end{bmatrix}$$

$$B = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} \frac{3}{23} & -\frac{1}{46} \\ -\frac{1}{46} & \frac{2}{23} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\frac{1}{92} \begin{bmatrix} 12 & -2 \\ -2 & 8 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} \frac{5}{46} & \frac{3}{46} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$= \frac{4}{23} \rightarrow \mu^* = \frac{23}{4} (-4) \quad ? \text{ OFF}$$

$$= -23$$

Homework 10:

Q1: Vector Argument, Linear Inequality Constraints

1) $f(w) = 4w_1^2 + w_2^2$, with $Aw \leq \vec{b}$
where $A = [-1 \ 0]$, $\vec{b} = [-2]$ A $Aw - b = [-1 \ 0] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} - (-2) = -w_1 + 2$

a) Yes it is convex

b) $L(w, \mu) = 4w_1^2 + w_2^2 + \mu(-w_1 + 2)$

c) Because: $4w_1^2 = \frac{1}{2}(8w_1^2)$
 $w_2^2 = \frac{1}{2}(2w_2^2)$
 $K = \begin{bmatrix} 8 & 0 \\ 0 & 2 \end{bmatrix}$, $q = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ (Because no constant term)

d) Unconstrained minimizer:

$$\Delta f = 8w_1 = 0 \quad w_1 \geq 2$$
$$2w_2 = 0 \quad 0 \geq 2?$$

Does $(0, 0)$ work? X Wrong

The unconstrained minimizer is not feasible.

e) only one active constraint: $w_1 \geq 2$
 $A = [-1 \ 0]$

f) KKT Matrix for active constraint:

$$\begin{bmatrix} K & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \mu \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ -2 \end{bmatrix}$$

$$\begin{bmatrix} 8 & 0 & -1 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \mu \end{bmatrix}$$

$$w_1 = 2$$

$$w_2 = 0$$

$$\mu = 6$$

2) $f(w) = 5w_1^2 - 8w_2 - 10w_1 + 2w_2^2 + 13$
with $A \leq b$, where $A = \begin{bmatrix} 10 & 0 \\ 0 & 4 \end{bmatrix}$ $b = \begin{bmatrix} -2 \\ -2 \end{bmatrix}$

a) Yes convex $\begin{bmatrix} 10 & 0 \\ 0 & 4 \end{bmatrix}$

b) $Aw - b \leq 0$

$\begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$

c) K is the hessian

$K = \begin{bmatrix} 10 & 0 \\ 0 & 4 \end{bmatrix}$

$q = \begin{bmatrix} -10 \\ -8 \end{bmatrix}$

$-w_1 \leq -2$
 $w_2 \leq -2$

d) w.r.t $w_1 = 10w_1 - 10 \rightarrow w_1 = 1$

w.r.t $w_2 = -8 + 4w_2 \rightarrow w_2 = 2$

=

\therefore not feasible

Q2: $f(t) = t^4 + t^2$, $p_4(t) \leq 0$, $p_4(t) = t^2 - 3t + 2$

a) $f(t) = 4t^3 + 2t$
 $f'(t) = 12t^2 + 2 > 0$

for all $t \therefore$ convex

b) $L(w, \rho) = t^4 + t^2 + \rho(t^2 - 3t + 2)$

c) $0 = 4t^3 + 2t$

$0 = 2t(2t^2 + 1)$ violation $p(0) \leq 0 \times$

$t=0$ $t=-\frac{1}{2}$

e) t feasible

region is between 2 and 3

d) $p_4(t) = t^2 - 3t + 2 = 0$

$(t-1)(t-2) = 0$

$t=1$ or $t=2$

f) $t=1$ is the minimizer

P5) $f(t) = (t+1)^2$, $p(t) \leq 0$

$p(t) = t^2 + 3t + 2$

a) $f'(t) = 2(t+1)$

$= 2t + 2$

$= 2$ $t > 0$ for all t

\therefore convex

b) $= (t+1)^2 + \rho(t^2 + 3t + 2)$

$t=-1$

c) $p(-1) = -1^2 + 3(-1) + 2$

$= 1 - 3 + 2$

$= 0$ $t=1$, feasible

7) $f(w; g) = \frac{1}{2}(w-g)^T(w-g)$

squared Euclidean distance from w and g .

unconstrained minimizer is always: $w = \vec{g}$

1) $f(w) = \frac{1}{2}(w - \begin{bmatrix} 4 \\ 4 \end{bmatrix})^T(w - \begin{bmatrix} 4 \\ 4 \end{bmatrix})$

$p(w) = w_1^2 \leq 4$

$|w_1| \leq 2$

2) unconstrained optimizer was $w_1 = 4$

which violates constraint

3) That means constraint is active

$w_1^2 = 4 \Rightarrow w_1 = \pm 2$

$L = \frac{1}{2}(w_1 - 4)^2 + (w_2 - 4)^2 + \lambda(w_1^2 - 4)$

7) $f(w; g) = \frac{1}{2}(w-g)^T(w-g)$

hyperparameter: $g = \begin{bmatrix} 4 \\ 4 \end{bmatrix}$

$p(w) = w_1^2$

Find global minimizer of f subject to $p \leq 4$.

$w_1^2 \leq 4$ or $|w_1| \leq 2$

$f(w) = \frac{1}{2}\|w-g\|^2$, squared distance

1) gradient: $\frac{1}{2}((w_1-4)^2 + (w_2-4)^2)$

w.r.t w_1 : $(w_1-4) + 2\lambda w_1 = 0$

w.r.t w_2 : $(w_2-4) = 0 \rightarrow w_2 = 4$

if $w_1^2 = 4 \Rightarrow w_1 = \pm 2$

2) candidates: $w = \begin{bmatrix} 2 \\ 4 \end{bmatrix}$, $w = \begin{bmatrix} -2 \\ 4 \end{bmatrix}$

$\lambda = 0.5$

w.r.t w_1 : $(w_1-4) + 2\lambda w_1 = 0$

w_2 : $w_2^* = 4$

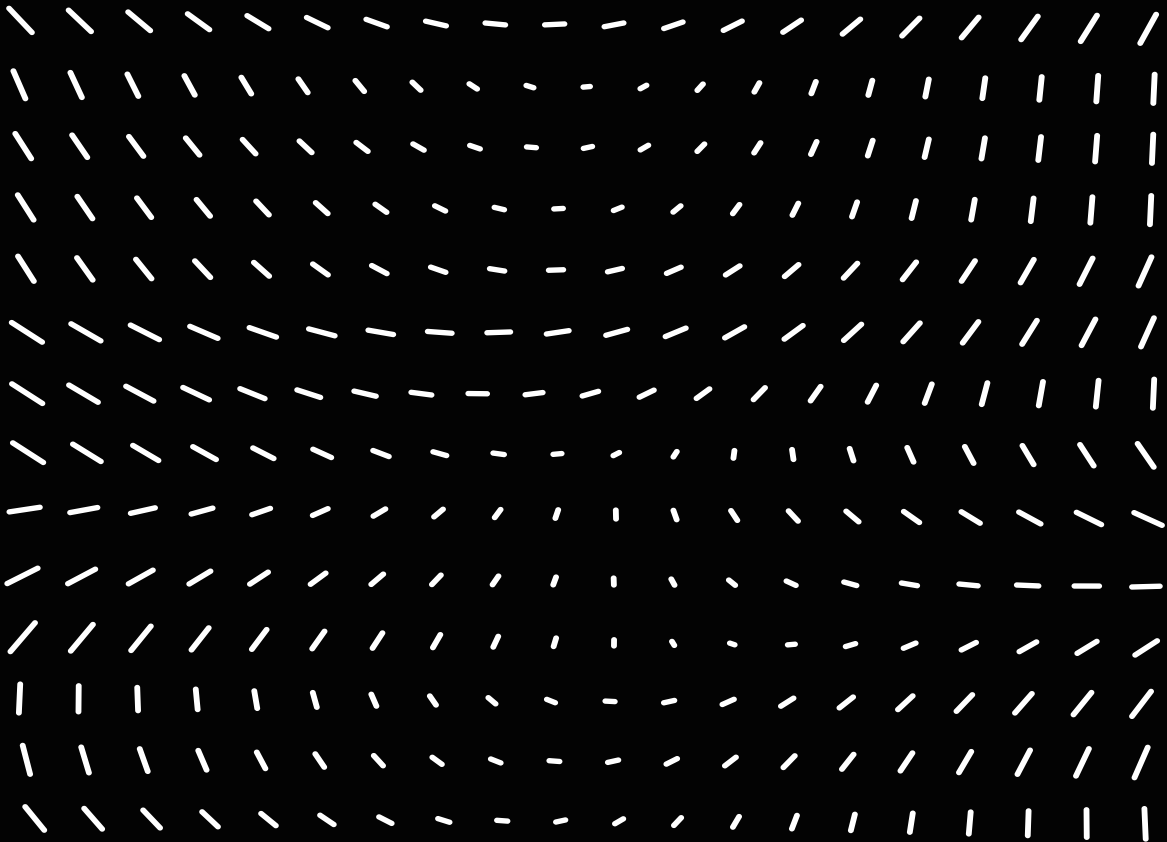
we know $w_1 = \pm 2$

Try $(-2, 4)$

$(2, 4)$

Test 5

Prep



Fall 2024 Test 5:

2) For data vectors and label values

$$x_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, x_2 = \begin{bmatrix} 3 \\ 0 \end{bmatrix}, x_3 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}, x_4 = \begin{bmatrix} -3 \\ 0 \end{bmatrix}, y^T = [1 \ 1 \ -1 \ -1]$$

SVM computation separates data with weight vector $\vec{w} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$ and $b = 0$

Find N_S of indexes of support vectors.

1) Find raw score for each point: $z = w^T x_i + b$

$$z_1 = 0.5(1) + 0.5(1) = 1$$

$$z_2 = 0.5(3) + 0.5(0) = 1.5$$

$$z_3 = 0.5(-1) + 0.5(-1) = -1$$

$$z_4 = 0.5(-3) + 0.5(0) = -1.5$$

2) Check if $y_i \cdot z_i = 1$

$$y_1 = +1, y_1 z_1 = (1)(1) = 1 \therefore \text{support vector}$$

$$y_2 = +1, y_2 z_2 = (1)(1.5) = 1.5 \therefore \text{not S.V.}$$

$$y_3 = -1, y_3 z_3 = (-1)(-1) = 1 \therefore \text{yes S.V.}$$

$$y_4 = -1, y_4 z_4 = (-1)(-1.5) = 1.5 \therefore \text{not S.V.}$$

Support vectors are $N_S = \{1, 3\}$

3) Consider the information $x_1 = \begin{bmatrix} -1 \\ 3 \end{bmatrix}, x_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, y = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \alpha = \begin{bmatrix} 0.25 \\ 0.25 \end{bmatrix}, b = -1$
find score $Z(g)$ for $g = \begin{bmatrix} 0.5 \\ -0.75 \end{bmatrix}$?

1) Compute dot products $x_i \cdot g$ 2) Multiply each term by α_i and y_i 3) $Z(g) = (-0.6875) + (0.0625) + b$

$$x_1 = \begin{bmatrix} -1 \\ 3 \end{bmatrix}, g = \begin{bmatrix} 0.5 \\ -0.75 \end{bmatrix}$$

$$= -2.75$$

$$x_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, g = \begin{bmatrix} 0.5 \\ -0.75 \end{bmatrix}$$

$$= -0.25$$

$$= -2.75(1)(0.25)$$

$$= -0.6875$$

$$= 0.25(-1)(-0.25)$$

$$= 0.0625$$

$$Z(g) = -1.625$$

4) For the embedding $\begin{bmatrix} a_1 \\ a_2 \end{bmatrix} \mapsto \begin{bmatrix} a_1^2 \\ a_2^2 \\ -\sqrt{2} a_1 a_2 \end{bmatrix}$, find the kernel function $K(\vec{u}, \vec{v})$.

$$u = (u_1, u_2)$$

$$v = (v_1, v_2)$$

$$\phi(u) = \begin{bmatrix} u_1^2 \\ u_2^2 \\ -\sqrt{2} u_1 u_2 \end{bmatrix}$$

$$\phi(v) = \begin{bmatrix} v_1^2 \\ v_2^2 \\ -\sqrt{2} v_1 v_2 \end{bmatrix}$$

Write it out

$$= (u_1^2)(v_1^2) + (u_2^2)(v_2^2) + (-\sqrt{2} u_1 u_2)(-\sqrt{2} v_1 v_2)$$

$$= (u \cdot v)^2$$

Shortcut: look at kernel is degree 2 \rightarrow quadratic, deg=3 \rightarrow cubic

: mixed term $a_1 a_2, a_1^2 a_2$ \rightarrow indicates polynomial

: $\|x\|^2 \rightarrow$ RBF

5) For Kernel $k(u,v) = (u \cdot v + 1)^2$, S.V., label values y_i , Lagrange multipliers, b scalar
Find score for $z(g) = \begin{bmatrix} 0.5 \\ -0.75 \end{bmatrix}$

1) Write scoring equation $z = \alpha_1 y_1 K + \alpha_2 y_2 K + b$

2) Compute kernel function

$$(x_1, g) = [-1, 3] \begin{bmatrix} 0.5 \\ -0.75 \end{bmatrix} = (-2.75 + 1)^2 = 3.0625$$

$$(x_2, g) = [0, -1] \begin{bmatrix} 0.5 \\ -0.75 \end{bmatrix} = (0.75 + 1)^2 = 3.0625$$

3) Sub in all values, compute: $z = (0.171)(1)(3.0625) + (0.171)(-1)(3.0625) - 1$
 $= -1$

To find S.V

$$y_i(x_i^T w) + b$$

$$x_1: 1(1) + 0 = 1 \quad \therefore$$

$$x_2: 1(2) + 0 = 2 \quad \therefore \text{N.S.}$$

$$x_3: -1(-1) + 0 = 1 \quad \{1, 3\}$$

$$x_4: -1(-2) + 0 = 2$$

$$z(g) = \left(\sum \alpha_i y_i (x_i \cdot g) + b \right)$$

$$= (0.25)(1)(-2.75) + (0.25)(-1)(-0.25) - 1$$

$$= -1.625$$

$$z(g) = (0.171)(1)(3.0625) + (0.171)(-1)(3.0625) - 1$$
$$= -1$$

Preparation - Test 5 Notes

CLS (Constrained Least Squares)

• minimizing squared error with constraints

• Normal Least Squares: $\min \|Xw - y\|^2$

↳ with constraint $\|w\|^2 \leq \theta$

θ is a hyperparameter that limits how large w can be

θ helps prevent overfitting

Tikhonov Regularization

• adds penalty to size of weight vector

• minimizes both error and model complexity

$$\min \|Xw - y\|^2 + \lambda \|w\|^2$$

λ is regularization parameter that controls trade off

• Might see regularization matrix R , instead of $\|w\|^2$

• Using $\lambda = 0$, removes the regularization term making it original squared error.

• Both CLS and Tikhonov use controlling w ,

$$\|w\|^2 = \sum_i w_i^2$$

• In CLS: it's a constraint

• In Tikhonov: it's a penalty

AS λ increases, weights decrease

Tikhonov Example

$$x = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, y = \begin{bmatrix} 3 \\ 4 \end{bmatrix}, \lambda = 1$$

Find w ?

$$v = (x^T x + \lambda I)^{-1} x^T y$$

$$w = (I + I)^{-1} \cdot I \cdot y = \begin{bmatrix} 1.5 \\ 2 \end{bmatrix}$$

- Primal formulation is the original problem that tries to maximize the margin between classes. Assuming the data is linearly separable.

(Hard Margin SVM)

Core Idea: Finding best separating hyperplane with the widest margin possible

Hard Margin

- given an objective $\frac{1}{2} \|w\|^2$ try to minimize weight vector which maximizes margin

- Constraint ensures no misclassification:

$$y_i(w^T x_i + b) \geq 1 \text{ for all } i$$

- To solve primal SVM problem, use Lagrangian method

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y_i(w^T x_i + b) - 1]$$

- convert it to solve with derivatives and KKT

- Stationary point, optimal w, b is found by taking partial derivatives of L and setting them to 0:

$$\frac{1}{2} \|w\|^2 - \alpha_i [y_i(w^T x_i + b) - 1]$$

• Taking derivative w.r.t to w_1/w_2 results in

$$\begin{aligned} w_2 - \alpha & & w_2 = \alpha \\ w_1 - \alpha & \text{ or } & w_1 = \alpha \end{aligned}$$

- α_i is a Lagrange multiplier

if $\alpha_i = 0$, point is not on margin / not support vector

if $\alpha_i > 0$, point is on margin, point is support vector

Dual Formulation
with Hard SVM:

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j$$

Subject to $\alpha_i > 0$

$$\sum \alpha_i y_i = 0$$

Quadratic Problem

Soft Margins as Slack Variables

- In SVM's, margin is gap between classes
- Slack variables allow some points to be inside margin or even be misclassified.

Usual constraint: $y_i(w^T x_i + b) \geq 1 - \xi_i$
 $\xi_i \geq 0$

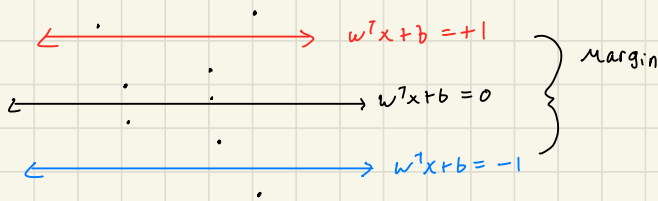
$\xi = 0$, correctly classified, outside margin \rightarrow Ideal

$0 < \xi < 1$, inside margin but correctly classified \rightarrow less than ideal

$\xi > 1$, misclassified point \rightarrow worst

C Controls how harshly margin violations are penalized

The decision boundary is $w^T x + b = 0$



Kernel Function to find support vectors and bias value.

• Dual SVM: $f(x) = \sum_i \alpha_i y_i K(x_i, x) + b$

Points with $\alpha_i > 0$ are support vectors

To find b : any support vector with $0 < \alpha_i < C$

$$b = y_i - \sum_j \alpha_j y_j K(x_j, x_i)$$

Kernel replaces dot products: $K(x, z) = x^T z$
Poly nomial: $(x^T z + 1)^d$

Example

i	α_i	y_i	x_i	
1	0.4	+1	$[1, 1]$	→ S.V
2	0.0	+1	$[2, 2]$	
3	0.6	-1	$[0, 1]$	→ S.V

Find bias.

$$b = 1 - (\alpha_i y_i K(x_1, x_3) + \alpha_i y_i K(x_1, x_1))$$

$$K(x_1, x_1) = 2$$

$$K(x_1, x_3) = 1$$

$$= 1 - (0.4(1)(2) + (0.6)(-1)(1))$$

$$b = 0.8$$

S.V are points closest to decision boundary

Support boundary and they lie on or inside margin

Predicting 2 questions

- Tikhonov / CLS
- Slack variables
- Dual constraints for soft margins

Scoring a data vector (kernel function)

$$f(x) = \sum_i \alpha_i y_i K(x_i, x) + b$$

Then if $f(x) \geq 0$ class 1, $f(x) < 0$, class 2

Ex

given: $x_1 = [1, 1]$, $\alpha = 0.4$, $y = +1$

$x_3 = [0, 1]$, $\alpha = 0.6$, $y = -1$

$$b = 0.8$$

Score point $x = [2, 1]$

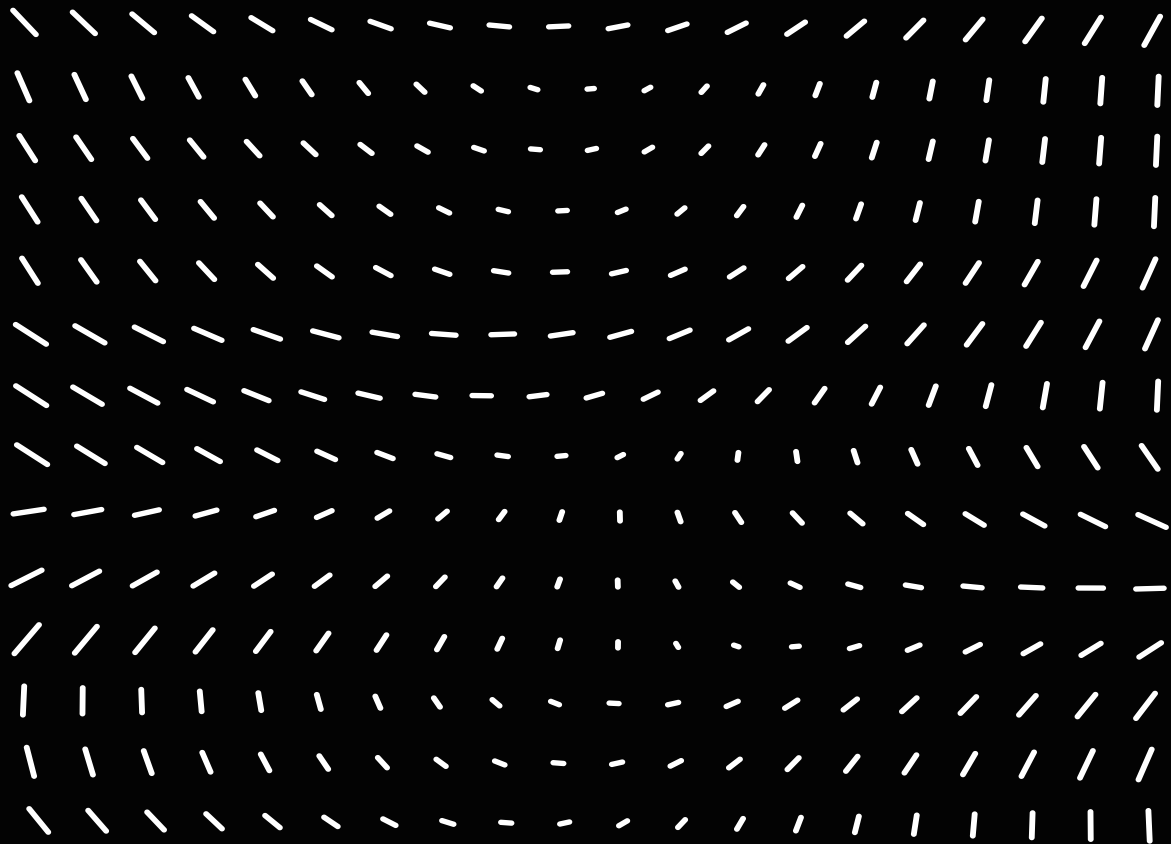
$$K(x_1, x) = 3$$

$$K(x_3, x) = 1$$

$$f(x) = (0.4)(1)(3) + (0.6)(-1)(1) + 0.8$$

$$= 1.4, \text{ positive class } +1$$

Exam Review



Test 1 Takeup

- 1) For function and point, $f(t) = \ln(1+t)$ $t_0 = 0$
find quadratic approx using truncated Taylor series?

$$f(t) = f(t_0) + f'(t_0)(t-t_0) + \frac{1}{2} f''(t_0)(t-t_0)^2$$
$$\left. \begin{array}{l} f(0) = 0 \\ f'(t) = \frac{1}{t+1} \quad f'(0) = 1 \\ f''(t) = -\frac{1}{(t+1)^2} \quad f''(0) = -1 \end{array} \right\} = 1(t) + \frac{1}{2}(-1)(t)^2 = \frac{-t^2}{2} + t$$

- 2) For a function, describe stationary point at $t^* = 0$? $f(t) = t^3 - 3t^2$

$$f(t) = t^3 - 3t^2 \quad \text{Use Second Derivative test}$$
$$f'(t) = 3t^2 - 6t \quad \text{This is concave, local maximizer}$$
$$f''(t) = 6t - 6$$
$$f''(0) = -6$$

- 3) For function $f(t) = t^4 + 4t^3$ has stationary points $t_1^* = -3$, $t_2^* = 0$. Describe the convexity.

$$f(t) = t^4 + 4t^3$$
$$f'(t) = 4t^3 + 12t^2$$
$$f''(t) = 12t^2 + 24t$$

Second Deriv Test:
 $f(-3) > 0 \therefore \text{min, convex}$
 $f(0) = 0 \therefore \text{indeterminate}$

- 4) Backtracking: given $f(t) = t^3 - t^2$, $\beta = \frac{1}{2}$, $t_1 = 1$, $s_0 = 1$ *

Find stepsize $s = B^r s_0$ for $f(t)$ after single step of backtracking.

Armijo condition: $f(t_k + B^r s_0 \alpha_k) \leq f(t_k) + \alpha_k B^r s_0 \alpha_k$
 $s = B^r s_0$

1) Plug in numbers

$$f(1) = 0$$

$$f'(t) = 3t^2 - 2t \quad f'(1) = 1, \quad d = -1$$

$$f''(t) = 6t - 2$$

For $r=0$, $s=1$

$$f(0) = 0 \leq -\frac{1}{2}$$

For $r=1$, $s=0.5$

$$f(1-\frac{1}{2}) = -\frac{1}{8} \leq -\frac{1}{4}$$

For $r=2$

$$f(1-\frac{1}{4}) = -\frac{9}{64} \leq -\frac{1}{8} \therefore s = \frac{1}{4} \quad s = \frac{1}{2} \cdot 1 = 0.25$$

5) For function f and point w_0 and direction \vec{v}

$$f(w) = w_1^2 + (w_2 + 1)^3 \quad w_0 = \begin{bmatrix} 4 \\ 1 \end{bmatrix}, \quad \vec{v} = \begin{bmatrix} -2 \\ -1 \end{bmatrix}$$

Find directional deriv?

1) Partial, $f'(w) = 2w_1 + 3(w_2 + 1)^2$

2) Plug in w_0 $[4, 12]$

$$\nabla f \cdot v = [4, 12] \begin{bmatrix} -2 \\ -1 \end{bmatrix} = -20$$

6) Function $f(w) = (w_1 - 1)^2 + w_2^3$ has stationary point at $w^* = [0]$

From eigenvalues of Hessian, describe convexity?

$$f(w) = 2(w_1 - 1) + 3w_2^2$$

$$f''(w) = 2, \quad 6w_2$$

$$\begin{bmatrix} 2 & 0 \\ 0 & 6w_2 \end{bmatrix} \rightarrow \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix}$$

\therefore semi definite

because one eigenvalue = 0

$$\vec{w}_1 = w_0 - s_0 \nabla f$$

$$f(w) = 6w_1 + 4w_2 + 1$$

$$\|Rw\|^2$$

$$f'(w) = -[6, 4] + 1$$

$$= \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 6 \\ 4 \end{bmatrix}$$

$$= \begin{bmatrix} -5 \\ -3 \end{bmatrix}$$

$w_0 - s_0 \nabla f \rightarrow$ Steepest descent

$$w_0 - s_0 \nabla f$$

$$w_0 - s_0 \nabla f$$

$$w_0 - s_0 \nabla f$$

$$d = -\nabla f(w_0)$$

Test 2. Takeup

1) For a function $g_1(\vec{w}) = w_1^4 + w_2^4 + 2w_1^2 w_2^2 - 1$, with 3 new functions

$f_{1,1}(\vec{w}) = g_1(\vec{w}) + g_1(\vec{w})$, $f_{1,2}(\vec{w}) = (g_1(\vec{w}))^2$ and $f_{1,3}(\vec{w}) = \sqrt{g_1(\vec{w})} + 1$, State how many are convex?

1) first step is recognize its factorable: $g_1(\vec{w}) = w_1^4 + w_2^4 + 2w_1^2 w_2^2 - 1$

$$g_1(\vec{w}) = (w_1^2 + w_2^2)^2 - 1$$

Convexity Closure Rules:

- Positive scaling and sums of convex functions are convex
- Add/subtracting constants does not change convexity
- norms and squared norms are convex

1. $f_{1,1}(\vec{w}) = g_1(\vec{w}) + g_1(\vec{w}) = 2g_1(\vec{w})$

Positive scaling convex functions \rightarrow convex

\therefore 2 functions are convex

2. $f_{1,3}(\vec{w}) = \frac{2\sqrt{g_1(\vec{w})} + 1}{2\sqrt{(w_1^2 + w_2^2)^2 - 1} + 1}$

Squared norm is convex. Quadratic function is convex

3. $f_{1,2}(\vec{w}) = (g_1(\vec{w}))^2 \rightarrow f_{1,2}(\vec{w}) = \|w\|^8 - 2\|w\|^4 + 1$

- takes negative values

- w^2 is decreasing on negative inputs

2.) For a function $f_2(w)$, starting point w_0 , stepsize s_0

$$f_2(w) = 3w_1^2 + 2w_2^2, w_0 = [1], s_0 = 1$$

find point w_1 single step in direction of steepest descent?

$$w_1 = -\nabla f + s_0$$

$$= -[6, 4]$$

$$= 6w_1 + 4w_2$$

$$= [6, -4] + [1, 1] \quad w_1 = [-5, -3]$$

3.) For a function $f_3(w) = w_1^4 - w_2^2 + w_2^4$ and starting point $w_0 = [-1]$

Find descent vector \vec{d}

$$\vec{d} = -|\nabla f(w_0)|^T$$

$$\begin{aligned} \nabla f' &= 4w_1^3 - 2w_2 + 4w_2^3 \\ &= 4(-1)^3 - 2(-1) + 4(-1)^3 \\ &= [-4, -2] \end{aligned}$$

$$\begin{aligned} \nabla f'' &= 12w_1^2 - 2 + 12w_2^2 \\ &= \begin{bmatrix} 12w_1^2 & 0 \\ 0 & 12w_2^2 - 2 \end{bmatrix} = \begin{bmatrix} 12 & 0 \\ 0 & 10 \end{bmatrix} \end{aligned}$$

Evaluate at w_0

$$\begin{aligned} d &= - \begin{bmatrix} 1/12 & 0 \\ 0 & 1/10 \end{bmatrix} \begin{bmatrix} -4 \\ -2 \end{bmatrix} \\ &= \begin{bmatrix} 1/3 \\ 1/5 \end{bmatrix} \\ \vec{d} &= \begin{bmatrix} 0.33 \\ 0.2 \end{bmatrix} \end{aligned}$$

4) For an NLS problem where each entry of a function $\vec{r}(w)$ is

$$r_i(\vec{w}) = \|x_i - w^T\|_2 = ([x_i - w^T] [x_i - w^T]^T)^{1/2}$$

Describe row J_i of Jacobian matrix $J_r(w)$?

$$r_i(w) = \|x_i - w^T\|$$

$$J_i = \frac{\partial r_i}{\partial w}$$

$$r_i'(w) = -([x_i - w^T] [x_i - w^T]^T)^{-1/2}$$

$$\nabla r_i(w) = -\frac{x_i - w^T}{\|x_i - w^T\|}$$

NLS: residuals are non linear functions of \vec{w}

$$r_i(\vec{w}) = \|x_i - w^T\|$$

need gradient of loss function

Affects distance to each point

5) 1D Fermat Problem with X , damping coefficient and initial estimate t_0 ?

$$X = \left[\frac{1}{9} \right], \quad t_0 = 2, \quad \tau = 1$$

$$b_1 = b_0 + (J^T J + \tau I)^{-1} J^T r$$

$$t_1 = 2 + 3 + 1 \quad \text{3) } J^T r = [-1 \ 1] \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

$$= 7$$

$$b_1 = 2 - 4^{-1}(7) = 0.25 \text{ or } 3.75$$

\pm

1) Residuals:

$$r = [x - t_0]$$

$$r = \begin{bmatrix} 1-2 \\ 3-2 \\ 9-2 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \\ 7 \end{bmatrix}$$

2) Jacobian:

$$1/x_i - t_i$$

$$\frac{a_i - t}{|a_i - t|}$$

$$= J \begin{bmatrix} \frac{1-2}{|1-2|} \\ + \\ \frac{9-2}{|1-2|} \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix}$$

6) Match the contour Plot.

$$f_1 = (w_1^2 + w_2^2)$$

$$f_2 = (w_1 + w_2)^2$$

$$f_3 = (w_1^2 - w_2^2)$$

$$f_4 = (w_1^3 - 3w_1 w_2^2)$$

$$\text{Level Curve: } \{w \in \mathbb{R}^2 \mid f(w) = c\}$$

$$\text{Circles: } x^2 + y^2$$

$$\text{Ellipses: } ax^2 + by^2$$

$$\text{hyperbolas: } x^2 - y^2$$

horizontal opening hyperbola

Can't be parabola; d

no constant lines

one constant line

Test 3: Takeup

- 1.) For artificial neuron with logistic function, weight vector \vec{w} , and observations A , find linear response \vec{v} .

$$A = \begin{bmatrix} 2 & 3 \\ 0 & 2 \end{bmatrix}, \quad \vec{w} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \quad \begin{array}{l} 2 \text{ by } 3 \quad 3 \text{ by } 1 \end{array}$$

Design matrix is A
with ones vector

$$X = \begin{bmatrix} 2 & 3 & 1 \\ 0 & 2 & 1 \end{bmatrix} \quad \vec{w} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$\vec{v} = X\vec{w} \rightarrow [0, -1]$$

- 2) 3 Layer ANN, logistic function everywhere, 2 neuron hidden layer with linear response $\vec{v} = [2, 5]$. Find derivative as Jacobian?

1) Logistic function $\sigma(u) = \frac{1}{1+e^{-u}}$

$J_i = \text{diag}(\sigma'(z_1), \sigma'(z_2))$ 2) Activation $\phi = [0.8808 \quad 0.9933]$

The chance it's in class 1

3) $\psi = [0.8808(1-0.8808), 0.9933(1-0.9933)]$

$$\psi = [0.1050, 0.0066]$$

$$J\vec{\psi} = \begin{bmatrix} \psi_1 & 0 \\ 0 & \psi_2 \end{bmatrix} = \begin{bmatrix} 0.1050 & 0 \\ 0 & 0.0066 \end{bmatrix}$$

Jacobians always diagonal

- 3) Given weight vector \vec{w} , input vector \vec{x} , back-prop factor b . Find steepest descent vector \vec{d} ?

$$\vec{w} = \begin{bmatrix} -1 \\ 1 \end{bmatrix} \quad x = [3 \quad 2] \quad y = 1 \quad b = -0.5$$

1) Calculate linear response: $U = XW \rightarrow \begin{bmatrix} 2 \\ 1 \end{bmatrix} \quad [-1 \quad 1]$
 $v = 0$

2) Activation: $\phi(0) = \frac{1}{1+e^0} = 0.5, \quad \psi = 0.5(1-0.5) = 0.25$

3) Descent Vector: $\vec{d} = -[x \quad 1] \psi b$

$$\vec{d} = -[3 \quad 2 \quad 1] (0.25)(-0.5)$$

$$\vec{d} = 0.125[3 \quad 2 \quad 1] \longrightarrow \vec{d} = [0.375, 0.250, 0.125]$$

4) 3 layer ANN with ReLU activation function, with weight matrix

zW for hidden layer 2, weight vector $z\vec{w}$ and activation deriv $z\psi$ and input x .

$$zW = \begin{bmatrix} 1 & 0 \\ 1 & 2 \\ 1 & 1 \end{bmatrix} \quad z\vec{w} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad x = [3 \ -2]$$

1) Linear response: $zU = [x \ 1]$
 $U = xW$ $U = [3 \ -2 \ 1] \begin{bmatrix} 1 & 0 \\ 1 & 2 \\ 1 & 1 \end{bmatrix}$
 $zW = [2 \ -3]$

2) Hidden Activation $z\phi = [2 \ 0]$

3) Output Response: $[z\phi \ 1] z\vec{w} = 3$
 $[2 \ 0 \ 1] \begin{bmatrix} 1 \\ 1 \end{bmatrix} z\vec{w} = 3$

4) Activation: $z\phi = 3$

5) For a 3 layer ANN, with ReLU, weight vector $z\vec{w}$ and back-prop value $z\vec{b}$ from objective computation, $z\vec{w} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$, find back propagation values $z\vec{b}$ provided to layer 2? $z\vec{b} = 8$

error signal from loss function

hidden backprop = (output-weight) (output slope) (deriv)

1) Find weights: $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $z\vec{w} = [1 \cdot 8, 1 \cdot 8]$
 $= [8, 8]$

2) Find $m = z\vec{b} \cdot z\psi$
 $m = 1 \cdot 8$
 $m = 8$

$z\psi$ is usually always 1

If given logistic function

find $\frac{1}{1+e^u}$, (would be given u)

Then find $z\psi$, then multiply by b .

Test 4: Take up

1) For objective function, nonlinear equality constraint. Find Lagrange function;

$$f(\vec{w}) = (w_1 - 1)^2 + (w_2 - 1)^2, \quad p(\vec{w}) = w_1^4 - w_2^2$$
$$L(w, \nu) = (w_1 - 1)^2 + (w_2 - 1)^2 - \nu(w_2^2 - w_1^4)$$

Just watch Lagrange is usually for sign switch. $f(w) + \nu(P(w))$

2.) For Obj function $F(w) = \frac{1}{2} w^T K w$ and lin equality constraint $mw - c = 0$, where

$$K = \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix}, \quad m = [1 \ 1], \quad c = 2. \quad \text{Find minimizer}$$

1.) $\begin{bmatrix} K & m^T \\ m & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \nu \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ c \end{bmatrix}$ 3) Use Lin alg or calculator to solve for w_1, w_2 .

2.) $\begin{bmatrix} 2 & 1 & 1 \\ 1 & 4 & 1 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \nu \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 2 \end{bmatrix}$ 4) $w^* = \begin{bmatrix} 1.5 \\ 0.5 \end{bmatrix}$

3.) For Obj function $f(\vec{w}) = \frac{1}{2} w^T K w + q^T w$ and lin equality constraint $mw - c = 0$

where $K = \begin{bmatrix} 2 & 0 \\ 0 & 8 \end{bmatrix}, \quad m = [1 \ 1], \quad q = \begin{bmatrix} -4 \\ -8 \end{bmatrix} \quad c = 1$

Find minimizer: ν^* and Dual matrix B.

1) $\begin{bmatrix} 2 & 0 & 1 \\ 0 & 8 & 1 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \nu \end{bmatrix} = \begin{bmatrix} 4 \\ 8 \\ 1 \end{bmatrix} \rightarrow$ plug into Calc and find minimizer

3) $K^{-1} = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{8} \end{bmatrix}$

$$B = [1 \ 1] \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{8} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$B = [1/2 \ 1/8] \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$B = 0.625$$

4) For Convex function $f(t) = t^2 + 4t$ and nonlinear inequality constraint $p_1(t) = t^2 - 5t + 6$ with $p_1(t) \leq 0$ find minimizer t^*

1) Check unconstrained

$$f'(t) = 2t + 4 = 0$$
$$t = -2$$

2) Check feasible

$$p(-2) = 20 > 0 \quad \therefore \text{not possible}$$

3) Find feasible region

$$t^2 - 5t + 6 = 0$$
$$(t-3)(t-2)$$
$$t=3 \quad t=2$$

4) Check on boundaries

$$f(2) = 12$$

$$f(3) = 21$$

5) \therefore minimizer is $t^* = 2$

5) For objective function $f(w) = (w_1 - 3)^2 + (w_2 - 4)^2$ and nonlinear inequality constraint $P(w) = w_1^2 + w_2^2 - 9$,
 Find minimized w^* $P(w \leq 0)$

1) Try unconstrained, given circle eq with roots (3,4), Try $w = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$.

2) Check feasible $3^2 + 4^2 = 25 > 9$

Shortcut: $x_{proj} = R \frac{x}{\|x\|}$

$$\begin{aligned}
 R = \text{radius } (3) &\rightarrow = 3 \cdot \frac{1}{3} \begin{bmatrix} 3 \\ 4 \end{bmatrix} \\
 x = (3, 4) & \\
 \|x\| = 5 & \\
 &= \begin{bmatrix} 1.8 \\ 2.4 \end{bmatrix}
 \end{aligned}$$

6) For convex objective $[w-g]^T [w-g]$, where $g = \begin{bmatrix} 0 \\ 4 \end{bmatrix}$ and linear inequality constraints $Aw \leq b$ with $A = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 1 \end{bmatrix}$ $b = \begin{bmatrix} 0 \\ 0 \\ 4 \end{bmatrix}$
 find lin. ineq. constraints, exactly satisfied

$$Ag = b?$$

$$\begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 4 \end{bmatrix} = \begin{bmatrix} 0 \\ -4 \\ 4 \end{bmatrix} \rightarrow \begin{matrix} 0 \\ 4 \end{matrix}$$

\therefore 2 active constraints

Test 5 - Takeup

1) For CONVEX Objective: $f(t) = (t-1)^2$, with hyper parameter $\lambda=2$ and $R=1$ find minimizer of Tikhonov $T(t, \lambda)$.

$$\begin{aligned} \text{Tikhonov Objective: } T(t) &= (t-1)^2 + \lambda R t^2 \\ &= t^2 - 2t + 1 + 2t^2 \\ &= 3t^2 - 2t + 1 \\ &= z(3t-1) \\ T(t) &= (t-1)^2 + 2t^2 & 0 &= 6t - 2 & t &= \frac{1}{3} \end{aligned}$$

Just plug into Tikhonov and expand, derivative, set to 0.

2) For data vectors and labels

$$x_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, x_2 = \begin{bmatrix} 3 \\ 0 \end{bmatrix}, x_3 = \begin{bmatrix} -1 \\ -1 \\ 0 \end{bmatrix}, x_4 = \begin{bmatrix} -3 \\ 0 \\ 0 \end{bmatrix}, y = \begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \end{bmatrix}$$

$$\text{for } \vec{w} = \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix} \text{ and bias } b = 1/3$$

$$x_i: y_i(xw - b) = 1?$$

$$\begin{aligned} x_1: 1 \left(\frac{4}{3} - 1 \right) &= 1 \\ x_2: 1 \left(\frac{5}{3} - 1 \right) &= \frac{2}{3} \end{aligned}$$

$$\begin{aligned} x_3: -1 \left(-\frac{2}{3} - \frac{1}{3} \right) &= 1 \\ x_4: -1 \left(-1 - \frac{1}{3} \right) &= \frac{4}{3} \end{aligned}$$

\therefore Support Vectors $N: \{1, 3\}$

$$\begin{aligned} \|Rt\|^2 \\ 2t^2 &= 2t \cdot 2t \\ &= 4t^2 \end{aligned}$$

3) For the following, $x_1 = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$, $x_2 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$, $y = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$, $\alpha = \begin{bmatrix} 0.1667 \\ 0.1667 \end{bmatrix}$, $b = -1$
For $z(g)$, $\vec{g} = \begin{bmatrix} 0.5 \\ -0.75 \\ 0.0 \end{bmatrix}$

$$z(g) = (0.1667)(1)(-2.75) + (0.1667)(-1)(-0.25) - 1 = -1.4167$$

$$= (u \cdot v)^2$$

4) For embedding $\begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix} \leftrightarrow \begin{bmatrix} a_{11}^2 \\ a_{12}^2 \\ a_{11}a_{12} \\ -\sqrt{2}a_{11}a_{12} \\ -\sqrt{2}a_{11}a_{13} \\ -\sqrt{2}a_{12}a_{13} \end{bmatrix}$

For $u = [u_1, u_2, u_3]$

$v = [v_1, v_2, v_3]$

$$(-\sqrt{2}u_1v_2)(-\sqrt{2}u_1v_2) = 2(u_1u_2v_2v_2)$$

$$(-\sqrt{2}u_1v_3)(-\sqrt{2}u_1v_3) = 2(u_1u_3v_1v_3)$$

$$(\sqrt{2}u_2v_3)(\sqrt{2}u_2v_3) = 2(u_2u_3v_2v_3)$$

5) For kernel $k(u, v) = (u \cdot v + 1)$, S.V., labels, Lagrange

$$x_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, x_2 = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \vec{y} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \alpha = \begin{bmatrix} 0.0143 \\ 0.0143 \end{bmatrix}, b = -1$$

$$\begin{aligned} q_1 &= u_1v_1 \\ q_2 &= u_2v_2 \end{aligned}$$

$$z(g) \text{ for } g = \begin{bmatrix} 0.5 \\ -0.75 \\ 0 \end{bmatrix}$$

$$\begin{aligned} k_1 &= 3.0625 \\ k_2 &= 3.0625 \end{aligned}$$

$$z(g) = \alpha_1 y_1 k_1 + \alpha_2 y_2 k_2 + b$$

$$z(g) = (0.0143)(1)(3.0625) + (0.0143)(-1)(3.0625) - 1$$

$$= -1$$