16B  Prof. Anant Sahai

Today: Finish up solving systems of equations (nonline-)

Revisiting Classification.

Think about "losses"

How can we find a $\vec{u}_0$ s.t. $\vec{f}(\vec{x}_0, \vec{u}_0) = \vec{0}$?

i.e. How can I solve a system of nonlinear equations?

Key Idea: **Iteration**:                     Newton's Method.

To simplify notation:  Let $\vec{g}(\vec{u}) = \vec{f}(\vec{x}_0, \vec{u})$

↑ fixed      ↰ what I want to fix

Want to solve  $\vec{g}(\vec{u}) = \vec{0}$

Use linearization.  $\vec{g}(\vec{u}) \approx \vec{g}(\vec{u}_0) + \frac{\partial \vec{g}}{\partial \vec{u}}\Big|_{\vec{u}_0} (\vec{u} - \vec{u}_0)$

Can I solve this??  Let $A_0 = \frac{\partial \vec{g}}{\partial \vec{u}}\Big|_{\vec{u}_0}$

Want: $A_0 (\vec{u} - \vec{u}_0) \approx -\vec{g}(\vec{u}_0)$

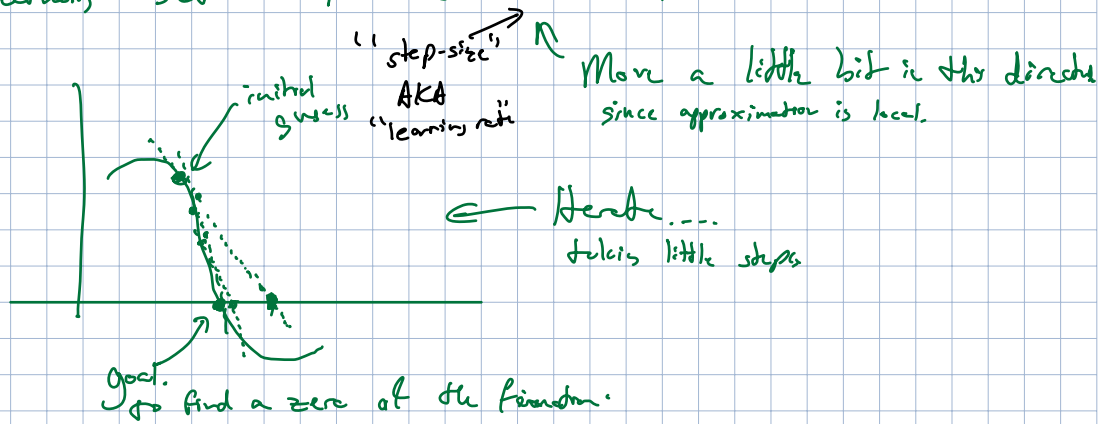Call  $\vec{\tilde{u}} = \vec{u} - \vec{u}_0$        $A_0 \vec{\tilde{u}} \approx -\vec{g}(\vec{u}_0)$

↖ If $A_0$ invertible (square-case),
can invert it.

Can solve using least-squares or min-norm approaches

depending on whether $A_0$ is tall or wide.

Suppose it were like least-squares. So $A_0$ is tall or square.

After I solve this, I'll get some $\vec{\tilde{u}}_q$.

I'll actually set $\vec{u}_1 = \vec{u}_0 + \alpha \cdot \tilde{u}_1$.

"step-size" AKA "learning rate"

Move a little bit in this direction since approximation is local.

initial guess

← Iterate.... takes little steps

goal. to find a zero of the function.

Can combine with the "transpose heuristic" ← From the Homework to iteratively solve least-squares.

$$\vec{u}_{i+1} = \vec{u}_i - \alpha\, A_i^\top\, g(\vec{u}_i)$$

where $A_i = \dfrac{\partial \vec{g}}{\partial \vec{u}}\Big|_{\vec{u}_i}$.

"Treat $A^\top$ as a cheap proxy for $A^{-1}$"

Often works to get you to a solution in the neighborhood of initial guess.

If we have a way of solving systems of nonlinear equations, then we can find local minima & maxima.

Recall 16A: Module 1: Solve Linear Equations. Module 3: Do least-squares by Solving Linear Equations.

Find places where no matter which direction you move in, nothing gets better.

Gives rise to a system of equations.

So we can bring the same approach in 16B...

Aside :

What do we understand?

1) Solving systems of linear equations
2) Linear Systems
3) Least-squares & Min-norm problems

What are we faced with in real-life?

A) Solving systems of nonlinear equations
B) Nonlinear Systems
C) Non-quadratic cost functions to minimize

How to deal with this?

We do an approximate reduction.

Use "linearisation" ideas (Taylor Expansion) to locally convert problems into things we know how to solve.

We take the approximation error $\longrightarrow$ call it a disturbance

Make sure we're "stable"

If needed, we iterate: Turns the algorithm into a dynamic system.

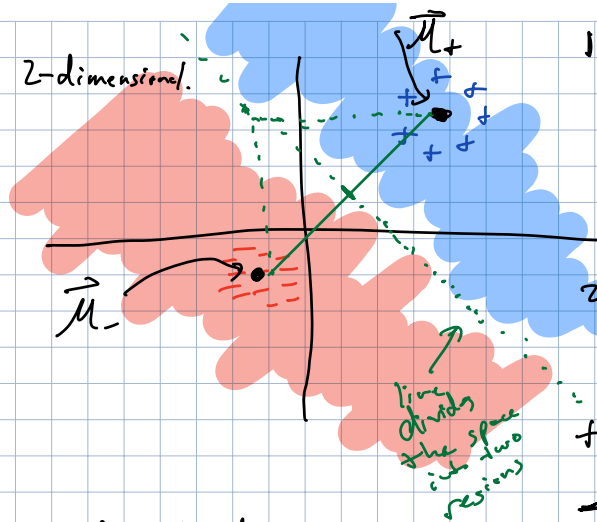Moral: When you don't know what to do, just try to make a little progress. And repeat.

Revisit Classification: Focus on binary classification:

We have data: $(\vec{x}_i, l_i)$ where $l_i = +1$ or $-1$

one kind of thing

A different thing

Most Intuitive Approach: (Lab approach)

Suppose $\vec{x}_i$ are 2-dimensional.



1. Compute means
$\vec{\mu}_+$ for + points
$\vec{\mu}_-$ for − points

2. Given a new point $\vec{x}$, we can classify it as

+ : if $\|\vec{x} - \vec{\mu}_+\| < \|\vec{x} - \vec{\mu}_-\|$

− : otherwise.

(line divides the space into two regions)

Question 1: What is the boundary between where we deem a point $\vec{x}$ as being + vs −??

Question 2: How do we express this decision rule so it is more clearly something related to a line?

$$\|\vec{x} - \vec{\mu}_+\|^2 < \|\vec{x} - \vec{\mu}_-\|^2$$

$$(\vec{x} - \vec{\mu}_+)^T (\vec{x} - \vec{\mu}_+) < (\vec{x} - \vec{\mu}_-)^T (\vec{x} - \vec{\mu}_-)$$

$$\vec{x}^T\vec{x} - 2\vec{\mu}_+^T\vec{x} + \vec{\mu}_+^T\vec{\mu}_+ < \vec{x}^T\vec{x} - 2\vec{\mu}_-^T\vec{x} + \vec{\mu}_-^T\vec{\mu}_-$$
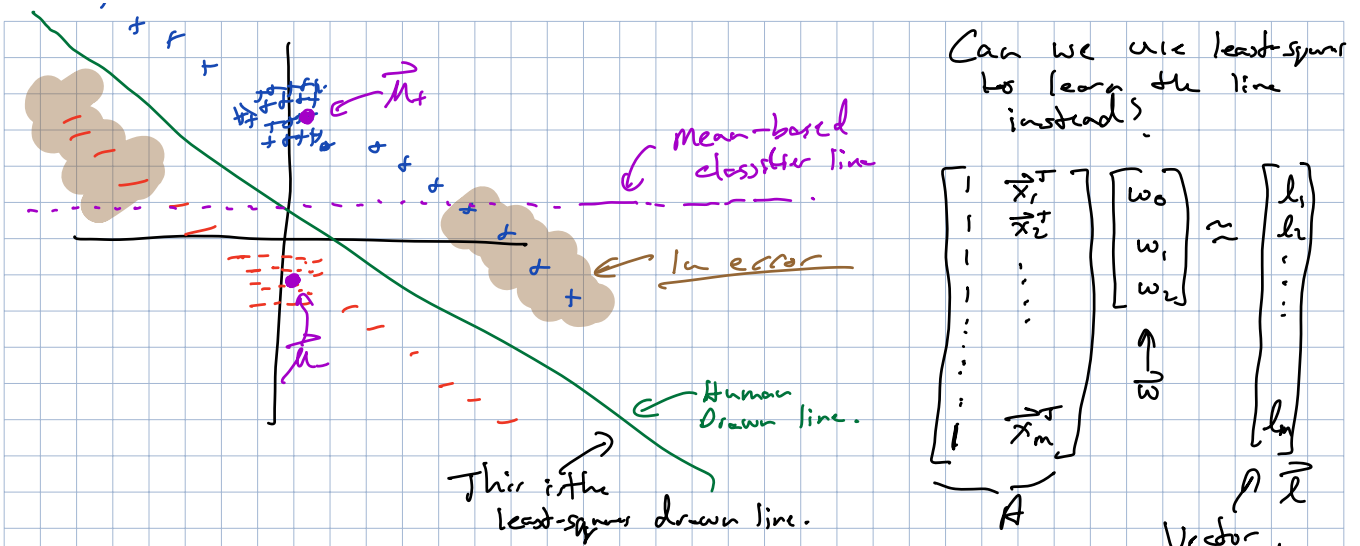
$$2(\vec{\mu}_- - \vec{\mu}_+)^T\vec{x} + \|\vec{\mu}_+\|^2 - \|\vec{\mu}_-\|^2 < 0 \implies \text{classify point } \vec{x} \text{ as } +$$

Alternatively $(\vec{\mu}_+ - \vec{\mu}_-)^T\vec{x} + \dfrac{\|\vec{\mu}_-\|^2 - \|\vec{\mu}_+\|^2}{2} > 0$

$$\text{Classify as the } \text{sign}\left( \underbrace{(\vec{\mu}_+ - \vec{\mu}_-)^T\vec{x}}_{\text{some row i.e. } [w_1, w_2]} + \underbrace{\dfrac{\|\vec{\mu}_-\|^2 - \|\vec{\mu}_+\|^2}{2}}_{\text{some constant. i.e. } w_0} \right)$$

Question: How else could we learn something of this form?

Motivation: When does mean-based classification get this wrong?!
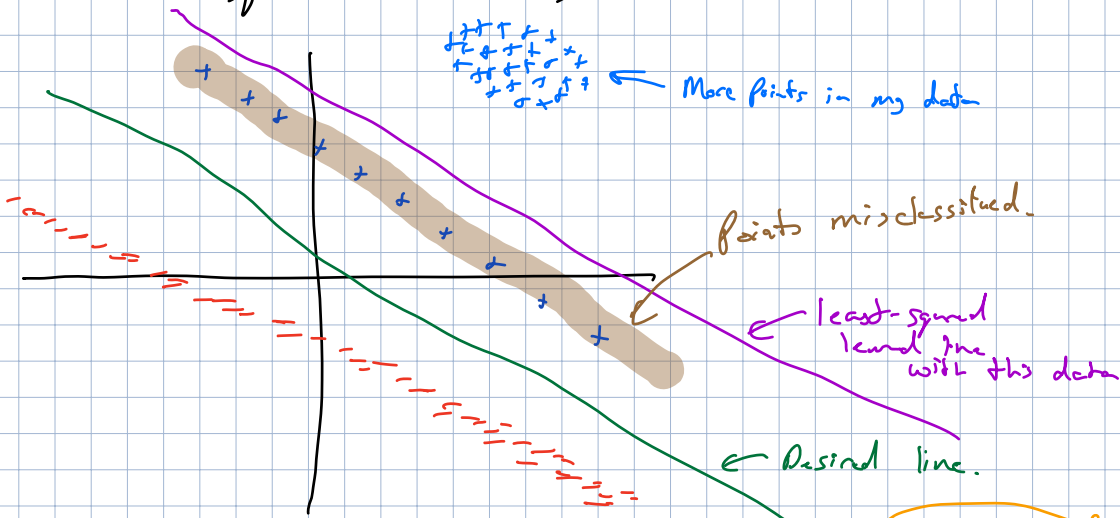
Mean-based classifier line

In error

$\vec{M}_+$

$\vec{M}$

Human Drawn line.

This is the least-squares drawn line.

Can we use least-squares to learn the line instead?

$$\underbrace{\begin{bmatrix} | & \vec{x}_1^T \\ | & \vec{x}_2^T \\ \vdots & \vdots \\ | & \vec{x}_m^T \end{bmatrix}}_{A} \underbrace{\begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}}_{\vec{w}} \simeq \underbrace{\begin{bmatrix} \ell_1 \\ \ell_2 \\ \vdots \\ \ell_m \end{bmatrix}}_{\vec{\ell}}$$

Vector is filled with $\pm 1$s.

Given data points in 2d, view them in 3d instead with a constant $1$ as the first entry.

Solving $A\vec{w} \simeq \vec{\ell}$    $\hat{\vec{w}} = (A^\sigma A)^{-1} A^T \vec{\ell}$
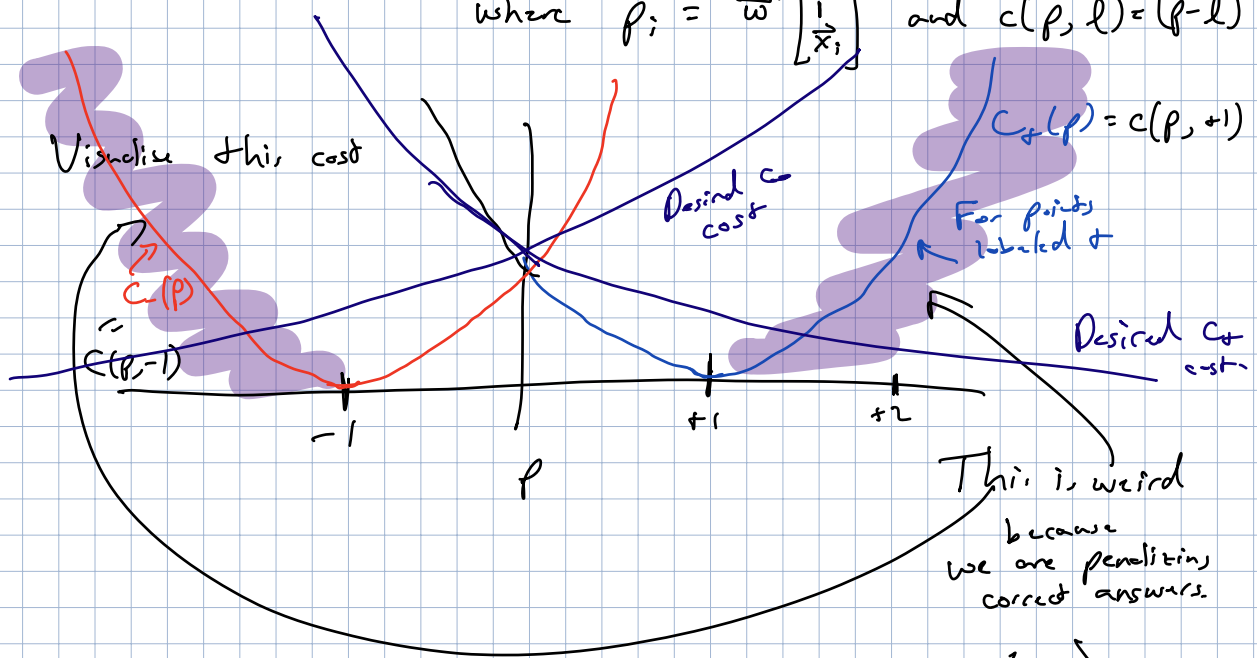
To classify a new point $\vec{x} \in$ 2d point.

Classify as $\text{sign}\left( \hat{\vec{w}}^T \begin{bmatrix} 1 \\ \vec{x} \end{bmatrix} \right)$.
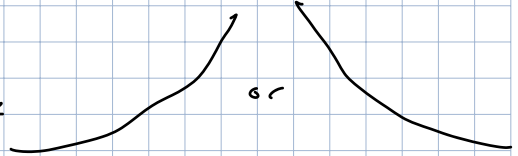
When does least-square based learning of classification not work?



More points in my data

Points misclassified.

Least-squared learned line with this data

Desired line.

Least Squares minimizes $\| A\vec{w} - \vec{\ell} \|^2 = \sum_{i=1}^{m} \left( \vec{w}^T \begin{bmatrix} 1 \\ \vec{x}_i \end{bmatrix} - \ell_i \right)^2$   Call this the cost associated with $i^{th}$ pt.

What went wrong?  Least-squares is minimizing $\sum_{i=1}^{m} c(p_i, l_i)$

where $p_i = \overline{w}^T \begin{bmatrix} 1 \\ \vec{x}_i \end{bmatrix}$ and $c(p, l) = (p - l)^2$

Visualise this cost

$C_+(p) = c(p, +1)$

For points labeled +

$C_-(p)$

$C(p, -1)$

Desired $c_-$ cost

Desired $c_+$ cost

$-1$     $p$     $+1$     $+2$

This is weird because we are penalizing correct answers.

What kind of cost looks like    or

e.g. exponential. Can we use a cost $C_+(p) = e^{-p}$ ?

& $C_-(p) = e^{+p}$ ?

$$c(p, l) = e^{-lp}$$

We want $\underset{\overline{w}}{\arg\min} \sum_{j=1}^{m} \exp\left(-l_i \, \overline{w}^T \begin{bmatrix} 1 \\ \vec{x}_i \end{bmatrix}\right)$

From now on assume $\vec{x}_i$ has $0^{th}$ component as $1$

$$\underset{\overline{w}}{\arg\min} \sum_{i=1}^{m} \underbrace{\exp(-l_i \, \vec{x}_i^T \overline{w})}_{c_i}$$

$C(\overline{w})$

How do I optimize this??

3 approaches:    ⓐ Pretend this is locally quadratic
and message it into least-squares form
to solve.

It turns
out
that these
two
give
the same
algorithm.

ⓑ Look for a local min/max by
seeing where the system of nonlinear
equations $\frac{\partial}{\partial \vec{w}} C(\vec{w}) = \vec{0}^T$ is solved.

(Use Newton's Method)

ⓒ Pretend it is linear and just iterate
⟹ go down the slope a little.

# SUPER IMPORTANT

Computers can take derivatives automatically,
eg. Pytorch _ (also tensorflow, JAX, etc...)

Will leverage this in discussion.