

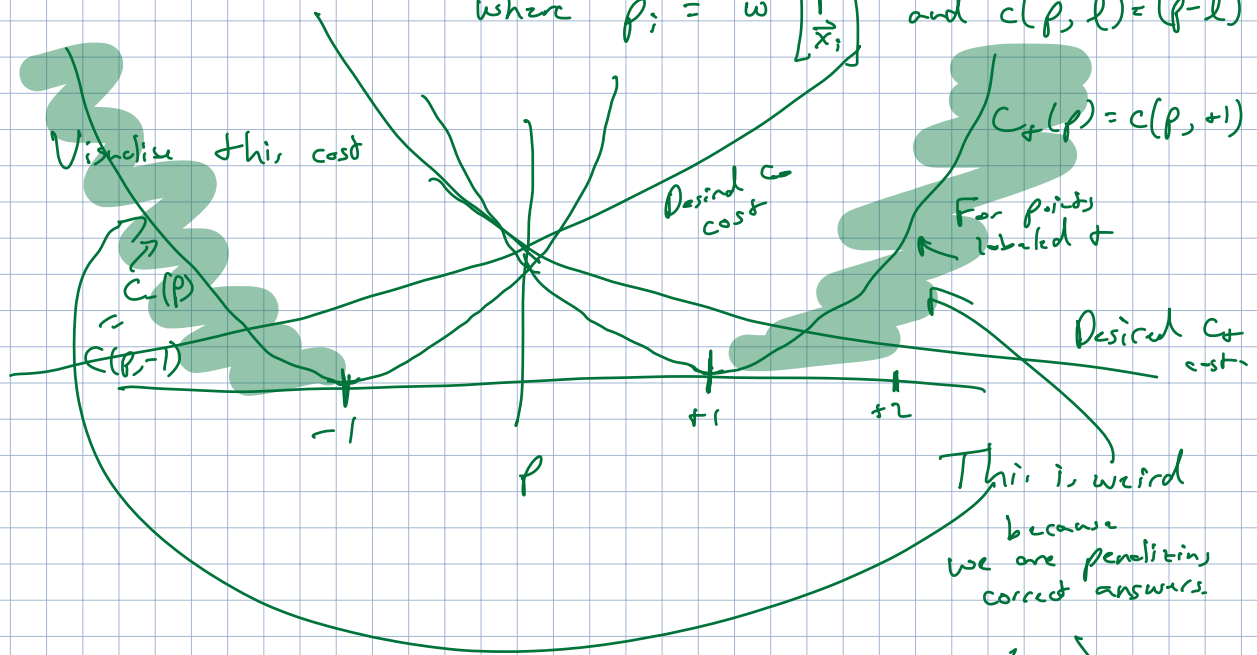
16B Prof. ANAND SAHAI

Today: Finish classification-inspired optimization story
Loose ends: Complex inner products

Announce: Happy Thanksgiving!
No Dis W
No Lec Th
HW B due next Fri (Not this Fri)
HW 14 will release before Wed
& due Wed RFL week.

Recall: Have labeled points: (\vec{x}_i, l_i) $i=1, \dots, m$
where $l_i = \pm 1$

What went wrong? Least-squares is minimizing $\sum_{i=1}^m c(p_i, l_i)$
where $p_i = \vec{w}^T \begin{bmatrix} 1 \\ \vec{x}_i \end{bmatrix}$ and $c(p, l) = (p - l)^2$



What kind of cost looks like

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

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

We want $\arg \min_{\vec{w}} \sum_{i=1}^m \exp(-l_i \vec{w}^T \begin{bmatrix} 1 \\ \vec{x}_i \end{bmatrix})$

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

$$\arg \min_{\vec{w}} \sum_{i=1}^m \exp(-l_i \vec{x}_i^T \vec{w})$$

$$\vec{w} \quad \underbrace{\quad \quad \quad}_{C(\vec{w})} \quad c_i(\vec{x}_i^T \vec{w})$$

How do I optimize this?

- 3 approaches:
- (a) Pretend this is locally quadratic and massage it into least-squares form for solve. around \vec{w}_0

$$C(\vec{w}) \approx C(\vec{w}_0) + \left. \frac{d}{d\vec{w}} C \right|_{\vec{w}_0} (\vec{w} - \vec{w}_0) + \frac{1}{2} (\vec{w} - \vec{w}_0)^T H_{\vec{w}_0}(C) (\vec{w} - \vec{w}_0)$$
 - (b) Look for a local min/max by seeing when the system of nonlinear equations $\frac{\partial}{\partial \vec{w}} C(\vec{w}) = \vec{0}^T$ is solved.
 (Use Newton's Method)
 - (c) Pretend it is linear and just iterate \Rightarrow go down the slope a little.

It turns out that these two give the same algorithm.

SUPER IMPORTANT

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

$$C(\vec{w}) \approx C(\vec{w}_0) + \left. \frac{d}{d\vec{w}} C \right|_{\vec{w}_0} (\vec{w} - \vec{w}_0) + \frac{1}{2} (\vec{w} - \vec{w}_0)^T H_{\vec{w}_0}(C) (\vec{w} - \vec{w}_0)$$

When $\frac{d}{d\vec{w}} C = \left[\frac{\partial C}{\partial w_1} \quad \frac{\partial C}{\partial w_2} \quad \dots \quad \frac{\partial C}{\partial w_n} \right]$

$$H_{\vec{w}} C = \begin{bmatrix} \frac{\partial^2 C}{\partial w_1^2} & \frac{\partial^2 C}{\partial w_1 \partial w_2} & \frac{\partial^2 C}{\partial w_1 \partial w_3} & \dots \\ \vdots & \ddots & \ddots & \ddots \\ \frac{\partial^2 C}{\partial w_n^2} \end{bmatrix}$$

How to minimize something like this? (exploit local quadratics)

1) Choose some initial \vec{w}_0 (e.g. $\vec{w}_0 = \vec{0}$)
Set $i=0$

2) Take a local quadratic approximation around \vec{w}_i

3) Find \vec{w}_{i+1}^{\sim} = minimizer of this quadratic

$$\text{Set } \vec{w}_{i+1} = (1-\eta) \vec{w}_i + \eta \vec{w}_{i+1}^{\sim} \quad \text{of } 0 \leq \eta \leq 1$$

↑
"etc"

4) See if "done" if so stop
else goto (2).

- ↳
- A) Is the $c(\vec{w})$ not changing too much.
 - A') Is w_i close to w_i^* ?
 - B) Have you run out of time.

How to minimize this quadratic?

$$\frac{1}{2} \vec{w}^T H \vec{w} + \vec{b}^T \vec{w} + d$$

\uparrow $n \times n$ matrix Hessian
 \uparrow n -dim row $\frac{d}{d\vec{w}} c(\vec{w}) \Big|_{\vec{w}_i}$

HOPE: Can we just match terms to least squares?

Do we care about constant term d ?

No. It doesn't matter for argmin.

Hope $\exists A$ so that

$$\frac{1}{2} H = A^T A \Rightarrow H = 2A^T A$$

What do we know?
Least Squares!

$$A \vec{w} \approx \vec{y}$$

Find \vec{w} that minimizes $\|A\vec{w} - \vec{y}\|^2$

$$\|A\vec{w} - \vec{y}\|^2 = (\vec{w}^T A^T - \vec{y}^T)(A\vec{w} - \vec{y})$$

$$= \vec{w}^T A^T A \vec{w} - 2\vec{y}^T A \vec{w} + \vec{y}^T \vec{y}$$

$$\text{Solved by } \hat{\vec{w}} = (A^T A)^{-1} A^T \vec{y}$$

Also could use SVD

Not worried about invertibility

↳ Can threshold to zero very small singular values

$$\text{If } A = U \Sigma V^T$$

$$\text{Then } (A^T A)^{-1} A^T = V \Sigma^{-1} U^T$$

Hope $\exists \vec{y}$ so that
 $\vec{b}^T = -2 \vec{y}^T A$

\leftarrow I need A to be invertible.
 $\vec{y} = -\frac{1}{2} A^{-T} \vec{b}$

With these hopes, I will get the solution: $\hat{\vec{w}} = (A^T A)^{-1} A^T \vec{y}$

$$= \left(\frac{1}{2} H\right)^{-1} \left(-\frac{1}{2}\right) A^T A^T \vec{b}$$

$$= -H^{-1} \vec{b}$$

$$= -H^{-1} \left(\frac{d}{d\vec{w}} c\right)^T$$

Still hanging on hope...

Called the Gradient $\nabla_{\vec{w}} c$

Is there an A s.t. $\frac{1}{2}H = A^T A$??

$$\frac{d}{d\vec{w}} c(\vec{w}) = \sum_{i=1}^M \frac{\partial}{\partial \vec{w}} c_i(\vec{x}_i^T \vec{w})$$

$$= \sum_{i=1}^M c_i'(\vec{x}_i^T \vec{w}) \begin{bmatrix} x_{i[0]} & x_{i[1]} & \dots & x_{i[n]} \end{bmatrix}$$

$$\frac{\partial}{\partial \vec{w}} [\vec{x}_i^T \vec{w}] = \frac{\partial}{\partial \vec{w}} \sum_{k=0}^n x_{i[k]} w[k]$$

$$= \sum_{i=1}^M c_i'(\vec{x}_i^T \vec{w}) \vec{x}_i^T$$

$$H_{\vec{w}} c = \frac{d}{d\vec{w}} \left[\frac{d}{d\vec{w}} c(\vec{w}) \right]^T$$

$$= \sum_{i=1}^M \frac{d}{d\vec{w}} c_i'(\vec{x}_i^T \vec{w}) \vec{x}_i$$

$$= \sum_{i=1}^M \vec{x}_i \frac{d}{d\vec{w}} c_i'(\vec{x}_i^T \vec{w})$$

$$= \sum_{i=1}^M \vec{x}_i c_i''(\vec{x}_i^T \vec{w}) \vec{x}_i^T$$

} By same logic as

? = $2A^T A$ for some A ??

Notice that if $c_i''(\vec{x}_i^T \vec{w}_0) \geq 0$, then

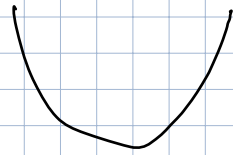
$$H_{\vec{w}} c = \sum_{i=1}^M \vec{x}_i \sqrt{c_i''} \cdot \sqrt{c_i''} \vec{x}_i^T \Rightarrow$$

$$= 2 A^T A$$

So $A = \frac{1}{\sqrt{2}} \begin{bmatrix} \sqrt{c_1''} & \vec{x}_1^T \\ \sqrt{c_2''} & \vec{x}_2^T \\ \vdots & \vdots \\ \sqrt{c_m''} & \vec{x}_m^T \end{bmatrix}$

$\left(\sqrt{c_i''(\vec{x}_i^T \vec{w}_0)} \right) \vec{x}_i^T$ ← positive

Aside: What are ^{Scalar} functions whose second derivative is always ≥ 0 called?



Convex ← From calculus.

Note: e^x and e^{-x} are both convex 😊

This is why convexity is something people can care about.

See 12.7 for more....

Let's check to make sure things make sense.

Let's use squared error: $c_i(p) = (p - l_i)^2$
 $c_i'(p) = 2(p - l_i)$
 $c_i''(p) = 2$

$$\frac{1}{2}H = \sum_{i=1}^m \vec{x}_i \vec{x}_i^T = A^T A \quad \text{if } A = \begin{bmatrix} \vec{x}_1^T \\ \vdots \\ \vec{x}_m^T \end{bmatrix} \quad \checkmark$$

$$\frac{d}{d\vec{w}} c = \sum_{i=1}^m 2(\vec{x}_i^T \vec{w} - l_i) \vec{x}_i^T = 2(A\vec{w} - \vec{l})^T A$$

Recall:
 In a least-sq. problem
 $A = \begin{bmatrix} \vec{x}_1^T \\ \vdots \\ \vec{x}_m^T \end{bmatrix}$ since
 $\|A\vec{w} - \vec{l}\|^2 = \sum_{i=1}^m (\vec{x}_i^T \vec{w} - l_i)^2$

Put everything together.

$$\begin{aligned} \|A(\vec{w}_0 + \delta\vec{w}) - \vec{l}\|^2 &= (A\vec{w}_0 - \vec{l} + A\delta\vec{w})^T (A\vec{w}_0 - \vec{l} + A\delta\vec{w}) \\ &= \|A\vec{w}_0 - \vec{l}\|^2 + \delta\vec{w}^T A^T A \delta\vec{w} + 2(A\vec{w}_0 - \vec{l})^T A \delta\vec{w} \\ &\approx \|A\vec{w}_0 - \vec{l}\|^2 + 2(A\vec{w}_0 - \vec{l})^T A \delta\vec{w} + \delta\vec{w}^T A^T A \delta\vec{w} \end{aligned}$$

Passes For something actually quadratic, the local quadratic approximation is exact.

Understanding "transpose trick" i.e. Gradient Descent.
 Just expand out to first derivative:

$$c(\vec{w}_0 + \delta \vec{w}) \approx c(\vec{w}_0) + \underbrace{\left. \frac{dc}{d\vec{w}} \right|_{\vec{w}_0}}_{\text{row vector}} \delta \vec{w}$$

Want to minimize or maximize? What direction should $\delta \vec{w}$ be?

We know by Cauchy-Schwarz $|\langle \vec{a}, \vec{b} \rangle| \leq \|\vec{a}\| \cdot \|\vec{b}\|$
 with equality only when $\vec{a} = \alpha \vec{b}$

$$\vec{b}^T \vec{a}$$

So $\delta \vec{w} = \left[\left. \frac{dc}{d\vec{w}} \right|_{\vec{w}_0} \right]^T$ is the direction that maximizes the change in c

$$\delta \vec{w} = - \left[\left. \frac{dc}{d\vec{w}} \right|_{\vec{w}_0} \right]^T$$

is the direction that minimizes

$\nabla_{\vec{w}} c(\vec{w})$: Gradient of cost w.r.t. \vec{w}

For Gradient Descent: $\vec{w}_{i+1} = \vec{w}_i - \eta \nabla_{\vec{w}} c(\vec{w}_i)$

scalar step size aka learning rate

For squared error (Least-Squares)

$$\nabla_{\vec{w}} c(\vec{w}_i) = \left[\left. \frac{d}{d\vec{w}} c(\vec{w}) \right|_{\vec{w}_i} \right]^T$$

$$= 2 A^T (A \vec{w}_i - \vec{l})$$

$$\vec{w}_{i+1} = \vec{w}_i + \eta \cdot 2 \cdot A^T (\vec{l} - A \vec{w}_i)$$

residual.

Source of A^T in HW Problem

Remember: η can't be too big or dynamics go unstable.

For generic losses. $\nabla_{\vec{w}} c(\vec{w}) = \sum_{i=1}^n c_i'(\vec{x}_i^T \vec{w}) \vec{x}_i$

$$\vec{w}_{i+1} = \vec{w}_i + \eta A^T \begin{bmatrix} -c_1'(\vec{x}_1^T \vec{w}_i) \\ \vdots \\ -c_m'(\vec{x}_m^T \vec{w}_i) \end{bmatrix}$$

residual
These derivatives of the losses take the place of the residual.

Loose End.

When we did upper-triangularization, we restricted to real eigenvalues/eigenvectors.

\Rightarrow We need a complex inner-product.

For real vectors $\langle \vec{a}, \vec{b} \rangle = \vec{b}^T \vec{a} = \vec{a}^T \vec{b} = \sum_{i=1}^n a_i b_i$

Doesn't work for complex vectors, ...

Want $\|\vec{a}\| = \sqrt{\sum_{i=1}^n |a_i|^2} = \langle \vec{a}, \vec{a} \rangle$

For complex vectors makes sense but squaring entries won't give us magnitude squared

$$\langle \vec{a}, \vec{b} \rangle \stackrel{?}{=} \vec{b}^T \vec{a} = \sum_{i=1}^n a_i b_i$$

$$\stackrel{?}{=} \vec{a}^T \vec{b} = \sum_{i=1}^n \bar{a}_i b_i$$

We choose this one.

i.e. $\langle \vec{a}, \vec{b} \rangle = \vec{b}^* \vec{a}$

where $\vec{b}^* = (\overline{\vec{b}})^T$

called conjugate transpose.

Next time: finish up complex inner products. (Also discussions next week)
& Review.