

Start with hypothesis testing for a few minutes for final projects

This class takes a strong optimization perspective.
- ultimately, our goal is to learn some

$$y \approx F(x)$$

- We take a parametric approach

$$y \approx F_w(x) = f(x; w) = f(x, w)$$

with parameter(s) w .

- What assumptions have we already made?
inductive bias

Example

$$f(x; w) = w_0 + w_1 x + \frac{1}{2} w_2 x^2$$

- We have a model of how we believe the world works but we don't know the specifics for our given problem.

- The above is the 2nd order kinematics eqn. relating position, velocity, acceleration, and time

$y = x_t$ position at time t .

$w_0 = x_0$ initial position

$w_1 = v_0$ initial velocity

$w_2 = a$ constant acceleration

- we assume newtonian physics and constant acc. and retrieve this function class but what function should we use? eg what values of $w = [w_0, w_1, w_2]$?
- could relax assumptions with

$$f(x; w) = w_0 + w_1 x + \frac{1}{2} w_2 x^2 + \epsilon$$

where ϵ is random and non-systematic (iid) maybe it captures "jerk" and "jounce" etc. and measurement error.

- Function Class

Example

class Kinematic:

def __init__(self, w0, w1, w2):

:

def predict(time):

vel = self.w0 + self.w1 * time + 1/2 self.w2 * time²

return vel

- we can have infinitely many of these

- uniquely determined by w_0, w_1, w_2

$$\text{Kinematic}(1, 2, 3) == \text{Kinematic}(1, 2, 3)$$

- Optimization: how do we find the best instance of our function class?

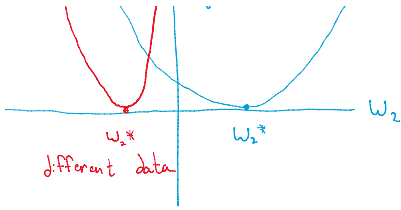
- what does "best" mean?

- objective function

Example: $\|f(x; w) - y\|^2$

assume we know w_0, w_1
Objective





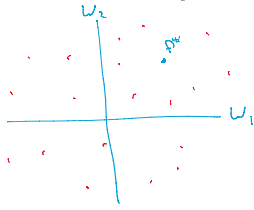
Question:

Should shape change with different data?
 is it always parabolic for this choice of F and objective?

- so how do we find the best:

$$F^* = \underset{F}{\operatorname{arg\,min}} \mathcal{L}(F(x), y) = \underset{w}{\operatorname{arg\,min}} \mathcal{L}(F(x; w), y)$$

- guess and check (random search / genetic algs etc.)



- gradient descent

$$F = F - \alpha \nabla_F \mathcal{L}$$

note this means we need to define what it means to add/multiply functions.
 because F is uniquely defined by w , this is easy.

- MLE / MAP

- MLE: $P(D|F)$ (if the world worked according to F , what are chances of seeing this (x, y) pair?)

$$\text{- MAP: } P(F|D) = \frac{P(D|F)P(F)}{P(D)}$$

Example

if my $x_0 = 0$ and velocity = 2 km/h
 acceleration = 0 km/h^2

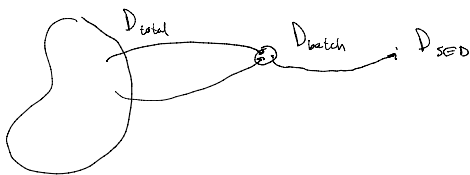
MLE: what is probability that I'm at position 100 at $t=1$?

MAP: given I'm at position 100 at $t=1$, what is the probability the above model is correct?

$$P(F) = \text{I'm pretty sure my velocity is } 5 \text{ km/h ish}$$

- Random Variables

- P(D) what is this?



If I ran my car long enough, I could collect all data about position vs. time
 or if I ran many exactly the same cars, etc.

Instead we get small batch of D .

When we deploy we reintroduce ourselves back to big \mathcal{D} .

- Gradients are random

- $\nabla_w(\mathcal{D})$ is a function of \mathcal{D}

Example $x \sim \mathcal{N}(0,1)$ is random

$x+2$ is random
 $2x$ is random
 x^2 is random
 $\mathbb{E}[x]$ is not

- how do I know ∇ is random?

sample many data points \uparrow and compute ∇ for each from \mathcal{D}_{total}

- Is $\nabla(\mathcal{D}_{set})$ random? yes
- Is $\nabla(\mathcal{D}_{total})$ random? no
- Is $\nabla(\mathcal{D}_{batch})$ random? yes according to \mathcal{D}_{total}
no according to \mathcal{D}_{batch}

- Bias / Variance

$$\text{Bias: } \mathbb{E}[\hat{x} - x]$$

- we want $\nabla(\mathcal{D}_{total})$, but don't have access

- is $\nabla(\mathcal{D}_{set})$ a biased estimate of $\nabla(\mathcal{D}_{total})$?

$$\mathbb{E}[\nabla(\mathcal{D}_{set}) - \nabla(\mathcal{D}_{total})] \stackrel{?}{=} 0$$

$$\mathbb{E}[\nabla(\mathcal{D}_{total})] = \nabla(\mathcal{D}_{total}) \checkmark$$

$$\begin{aligned} \mathbb{E}[\nabla(\mathcal{D}_s)] &= \frac{1}{n} \sum_{i=0}^n \nabla(\mathcal{D}_i) = \frac{1}{|\mathcal{D}|} \sum_{d \in \mathcal{D}} \nabla(d) \\ &= \nabla(\mathcal{D}_{batch}) \end{aligned}$$

- is $\mathbb{E}[\nabla(\mathcal{D}_{batch}) - \nabla(\mathcal{D}_{total})] = 0$?
hopefully!

- what about variance?

$$\begin{aligned} \text{var}(\nabla(\mathcal{D}_s)) &= \mathbb{E}[\nabla(\mathcal{D}_s)^2] - \mathbb{E}[\nabla(\mathcal{D}_s)]^2 \\ &= \mathbb{E}[\nabla(\mathcal{D}_s)^2] - \nabla(\mathcal{D}_b)^2 \end{aligned} \left. \vphantom{\begin{aligned} \text{var}(\nabla(\mathcal{D}_s)) &= \mathbb{E}[\nabla(\mathcal{D}_s)^2] - \mathbb{E}[\nabla(\mathcal{D}_s)]^2 \\ &= \mathbb{E}[\nabla(\mathcal{D}_s)^2] - \nabla(\mathcal{D}_b)^2 \end{aligned}} \right\} \text{element-wise}$$

can also consider $\text{Cov}(\nabla, \nabla)$.
considering on Var for simplicity

- So what?

- SURF: what if we use 2 estimates of $\nabla(\mathcal{D}_b)$ to have even lower variance

$$\begin{aligned} \nabla(\mathcal{D}_{surf}) &= \nabla(d_1) - (\nabla(d_2) - \mathbb{E}[\nabla(d_2)]) \\ &= \nabla(d_1) - [\nabla(d_2) - \nabla(\mathcal{D}_b)] \end{aligned}$$

why not just use $\nabla(\mathcal{D}_b)$ which is lowest variance?

maybe we can approximate it cheaply or compute it infrequently.

- Stepsizes:

$$f = f - \alpha \nabla(D)$$

$$\alpha = \frac{1}{\text{var}(\nabla(D))} \quad \text{"divide out variance"}$$

recall variance is second centered moment. That means

$$\mathbb{E}[\nabla(D)^2] - \nabla(D_0)^2 \quad \text{if we drop "center" } \nabla(D_0)^2$$

then we get

$$\begin{aligned} \mathbb{E}[\nabla(D)^2] &\approx \sum \nabla(D)^2 && \leftarrow \text{rms prop!} \\ &\approx (1-\beta)V + \beta \nabla(D)^2 \end{aligned}$$

- what if model is non-linear?

$$y = w_0 + w_1 x e^{-\frac{x}{D_3}} + w_2 [1 - e^{-\frac{x}{D_3}}] + \frac{1}{2} w_2 x^2$$

(kinematics + linear drag) with minor simplifications

$$\mathcal{L} = \|f(x) - y\|^2$$

$$\nabla_w (f(x) - y)^2 = 2 (f(x) - y) \nabla_w f(x)$$

for w_0, w_1, w_2 easy
for w_3 need product rule