

2) "I have 3 kids. Each of them rolled a pair of dice. I have at least one boy who rolled (1,1)."

What is the prob. that I have 3 boys?

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

BBB = all 3 kids are boys
 \tilde{B}_s = at least one is a boy who rolled (1,1)

$$P(BBB | \tilde{B}_s) = \frac{P(\tilde{B}_s | BBB) P(BBB)}{P(\tilde{B}_s)}$$

$$P(BBB) = \frac{1}{8} \quad P(\tilde{B}_s | BBB) = 1 - \underbrace{(1-\epsilon)^3}_{\text{prob. no one rolled (1,1)}}$$

$$\epsilon = \text{prob. of (1,1)} = \frac{1}{36}$$

$$1 - \epsilon = \text{prob. did not roll (1,1)}$$

trick to calc. denominator in Bayes' rule:

$$1 = \sum_A P(A|B) = \sum_A \frac{P(B|A)P(A)}{P(B)}$$

$$1 = \frac{1}{P(B)} \sum_A P(B|A)P(A)$$

$$\Rightarrow P(B) = \sum_A P(B|A)P(A)$$

$$P(\tilde{B}_s) = \sum_{\substack{xyz \\ \in \text{all 3} \\ \text{kids combos}}} P(\tilde{B}_s | XYZ) \underbrace{P(XYZ)}_{1/8}$$

$$P(\tilde{B}_s | GGG) = 0$$

$$P(\tilde{B}_s | BGG) = \epsilon$$

$$\text{" } \begin{matrix} G B G \\ G G B \end{matrix} = \epsilon$$

$$P(\tilde{B}_s | BBG) = 1 - (1 - \epsilon)^2$$

⋮

$$\text{algebra : } P(\tilde{B}_s) = \frac{1}{8} \epsilon (12 - 6\epsilon + \epsilon^2)$$

$$\Rightarrow P(BBB | \tilde{B}_s) = \frac{1 - (1 - \epsilon)^3}{\epsilon (12 - 6\epsilon + \epsilon^2)}$$

$$\epsilon = \frac{1}{36} \Rightarrow \approx 0.247 > \frac{1}{7}$$

$$\text{limit of small } \epsilon \quad \text{Taylor expand} \quad \approx \frac{1}{4} - \frac{\epsilon}{8}$$

more abstract interpretation

$$\mathcal{P}(BBB | \overset{\text{data}}{\tilde{B}_s}) = \underbrace{\left[\frac{\mathcal{P}(\tilde{B}_s | BBB)}{\mathcal{P}(\tilde{B}_s)} \right]}_{\text{"posterior":}} \underbrace{\mathcal{P}(BBB)}_{\text{"prior":}}$$

Knowledge after accounting for the data \tilde{B}_s

rule for updating our knowledge

knowledge before knowing \tilde{B}_s

$\mathcal{P}(\tilde{B}_s | BBB) \Rightarrow$ "likelihood":
prob. of data given prior knowledge

(hypothesis)

$\mathcal{P}(\tilde{B}_s) \Rightarrow$ normalization const.

model fitting:

dataset \mathcal{D}

model w/ some parameters

$$\vec{w} = (w_1, w_2, \dots)$$

$$\mathcal{P}(\vec{w} | \mathcal{D}) = \frac{\mathcal{P}(\mathcal{D} | \vec{w})}{\mathcal{P}(\mathcal{D})} \mathcal{P}(\vec{w})$$

posterior

prior

(reflects physical ranges of possible params)

note: to do MAP
you only need to
max. $\mathcal{P}(\mathcal{D} | \vec{w}) \mathcal{P}(\vec{w})$
& can ignore $\mathcal{P}(\mathcal{D})$

alternatively: $\mathcal{P}(\vec{w}) = \text{const.}$
 \Rightarrow no clue about what the params are

Two major applications:

1) find "best" model parameters:

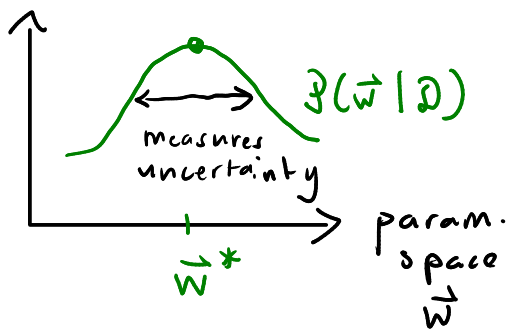
maximize $\mathcal{P}(\vec{w} | \mathcal{D})$ w/ respect to \vec{w}

$\Rightarrow \vec{w}^*$ is the "best" param. set.

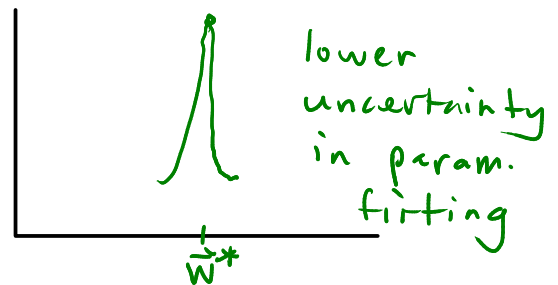
\Rightarrow MAP: maximum a posteriori fitting

(typical model fitting, i.e. training neural networks)

2) find $\mathcal{P}(\vec{w} | \mathcal{D})$ directly:



vs.



\Rightarrow gives info. about uncertainty in \vec{w}^* estimate

\Rightarrow generally hard: Bayesian neural network

\Rightarrow heat trick from stat. mech. that enables this (we will return to this later)