

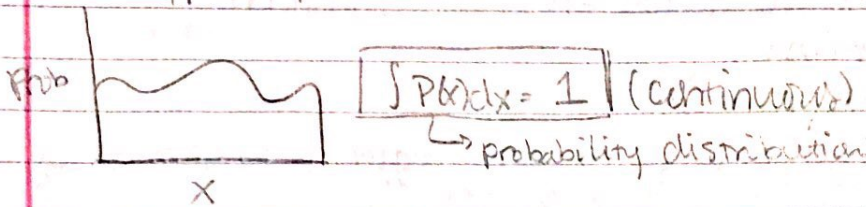
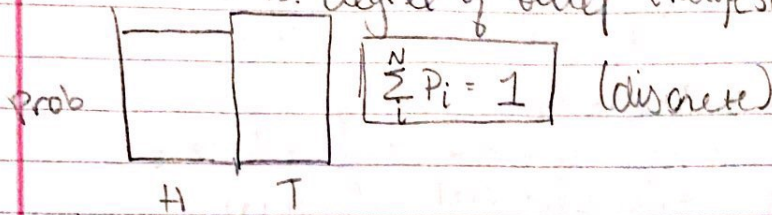
9/4/19 Learning Machine Learning

ex) Heads v. Tails

Probability of heads measured: $P_{meas}(H) = \lim_{T \rightarrow \infty} \frac{\#H}{\#T}$
 T : # of tails

→ frequentist approach

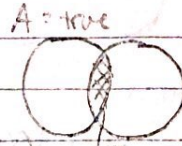
• vs. degree of belief (Bayesian)



Sum rule

• $A, B = \text{true/false}$

• $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
 $A \cup B$ ($A \vee B$) $A \cap B$



• $P(A \cap B) = P(A \cap B) = P(A)P(B|A) = P(B)P(A|B)$ (product rule)

↓ prob. of B assuming A (conditional)

• Bayes' thm.: $P(A|B) = P(A)P(B|A) / P(A \cap B)$

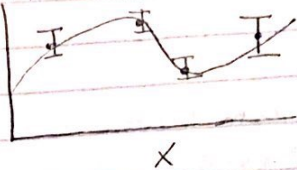
$P(\alpha|X) = \frac{P(\alpha)P(X|\alpha)}{P(X)}$
 posterior ← event → data → prior → likelihood → normalization/evidence

• likelihood: how consistent is the evidence w/ possibility α

• prior: prior belief in possibility (common sense)

LML 9/4/19

ex) fake data fitting



• assume Gaussian errors
→ likelihood? $\mathcal{L}_{\text{DATA}}(\vec{\alpha}) = \prod_i P(x_i | \vec{\alpha}) = \prod_i \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(x_i - \hat{y}(x_i))^2}{2\sigma_i^2}}$
↳ parameters ↳ variance

• x_i : data points
• $\hat{y}(x_i)$: predicted (model) value
→ log likelihood: $\log \mathcal{L} = -\sum \frac{(x_i - \hat{y}(x_i))^2}{2\sigma_i^2} + cte$
↳ chi-squared distribution: used to evaluate goodness of fit

→ maximize likelihood = minimize chi-squared
↳ can minimize χ^2 further by making model more flexible (more parameters/clof)
⇒ more complex model → more volume of possible

• minimize expectation value (avg.) $\mathbb{E}[(\hat{y}(x_i) - y_i)^2]$ (RSS, MSE)
model data label (supervised learning)

↳ can rephrase to get "variance" + "bias" terms

- bias: systematic shortcomings that make model incapable of fitting data
- variance: model fits data^(train) well, but high variance leads to low generalizability
↳ more parameters

ex) Slide 3: Decision Problem (Classifier)

- training a classifier to classify orange/blue
- Bayes' decision boundary = optimal boundary

Slide 4: math (Simple least squares model)

• linear least squares model: $\hat{Y} = \hat{\beta}_0 + \sum_{j=1}^p x_j \hat{\beta}_j$

↳ soln.: $\beta = (X^T X)^{-1} X^T y$ (optimal parameters)

LML

(ex) Con't

- k-nearest neighbors (Slide 6)
 - Clustering: choose nearest neighbors & average
 - Slide 7: boundary for $k=15$ (hyperparameter)
 - ↳ low variance, high bias for high k (underfitting)
 - ↳ vice versa for low k (overfitting)
- Slide 8: algorithm performance