

Mid-term exam info

- In-class exam
- Closed books/notes, closed electronics, 2 sided A4 hand-written (not printed) cheat sheet allowed (upload somewhere by some date)
- Videos on
- Academic integrity violations severe consequences

- Multiple choice
- Via Gradescope or Canvas
- ~25 questions – 80 mins (3-4 mins/question)
 - ~40% are hard (think open book), ~60% are easy
- Grading will be curved
- Ask questions via private Zoom chat

Topics

log
neg. \log likelihood loss
- density est

- Basics - Probability, Matrix/vector calculus, Optimization (convexity etc)
- Basic ML concepts – training vs test data, overfitting, generalization, ML tasks, loss metrics, optimal classifier/regressor, decision boundaries

$$E[(f(x)-y)^2] \xrightarrow{\text{L2, 0/1}} E[\mathbb{1}_{f(x) \neq y}] \quad P_{XY} = \prod_{i=1}^n P(x_i, y_i) \xrightarrow{\text{density est}} P(f(x) \neq y) \approx \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{f(x_i) \neq y_i}$$

- Distribution/Density estimation – MLE, MAP, Histogram, Kernel density estimation
- Classification – Naïve Bayes, Logistic Regression, Neural Networks, k-Nearest Neighbors
- Regression – Linear Regression, Ridge, Lasso, Neural Networks, Kernel regression

Comparison chart (classification)

$\rightarrow P(X,Y) \sim P(X|Y) \sim P(X)$

Algorithm	Generative/ Discriminative	Assumptions	Decision boundary	Loss function	Training Algos
<u>Naïve Bayes</u>	G	<p>$P(X Y), P(Y)$</p> <p>Probable Gaussian</p> <p>Bel/Multin</p> <p>(conditional) ind</p> <p>decision boundaries</p>	<p>GNB:</p>		
Logistic Regression	D	$P(Y X) = \frac{e^{\sum w_i x_i}}{1 + e^{\sum w_i x_i}}$	$e^{\sum w_i x_i} \geq 1$		
Neural Networks	D				
k-Nearest Neighbors	D	$P(Y X) \approx \frac{k_y}{k}$			

Comparison (regression)

Algorithm	Generative/ Discriminative	Assumptions	Decision boundary	Loss function	Training
Linear Regression <i>(Lasso, Ridge)</i> <i>↳ l_1 l_2</i>					
Neural Networks					
Kernel regression					

Practice problems (basics)

$$P(A, B) = P(A)P(B|A)$$

$$P(A, B, C) \\ \text{''} \\ P(C)$$

- Which of the following expressions is equivalent to $P(A, B|C)$?

(a) $\frac{P(A, B)}{P(C)}$ ✗

(b) $P(A|C)P(B|A, C)$ ←

(c) $P(C|A, B)P(B|A)P(A)$ ← $P(A, B, C)$

(d) $\frac{P(A)P(B|A)P(C|A, B)}{\sum_{a \in A} P(A=a)P(C|A=a) + \sum_{b \in B} P(B=b)P(C|B=b)}$

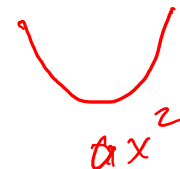
$$P(X, Y, Z) = P(X)P(Y|X)P(Z|X, Y) \\ = P(X)P(Y, Z|X)P(Z)$$

- The function $x^T A x$ where x is a d -dim vector and A is a $d \times d$ matrix is

- Convex if A is rank 1
- Convex because its quadratic in x
- Concave if A has all negative eigenvalues ✓

$$\begin{bmatrix} + & 0 \\ 0 & 0 \end{bmatrix}$$

$$x^T A x = -x_1^2$$



- Stochastic Gradient descent is faster but uses more memory than regular gradient descent.

Practice problems (ML intro)

- When the feature space is larger, overfitting is less likely.
- Which of the following are supervised learning tasks:
 - Predicting the proportion of students in class who slept < 7 hours $\equiv P(Y)$
 - Predicting rating of a movie given movie genre
 - Tagging tweets that are racially provocative
- A classifier that attains 100% accuracy on the training set and 70% accuracy on test set is better than a classifier that attains 70% accuracy on the training set and 75% accuracy on test set.

Practice problems (density/distribution estimation)

- You have received a shiny new coin and want to estimate the probability θ that it will come up heads if you flip it. A priori you assume that the most probable value of θ is 0.5. You then flip the coin 3 times, and it comes up heads twice. Which will be higher, your maximum likelihood estimate (MLE) of θ , or your maximum a posteriori probability (MAP) estimate of θ ? → 2/3

=

- The maximum ^{conditional} likelihood estimate of model parameter α for the random variable $y \sim N(\alpha x_1 x_2^3, \sigma^2)$, where x_1 and x_2 are random variables, can be learned using linear regression on n iid samples of (x_1, x_2, y)

↙

$$y = \alpha x_1 x_2^3 + N(0, \sigma^2)$$

$$= \alpha \cdot x$$

$$x = x_1 x_2^3$$

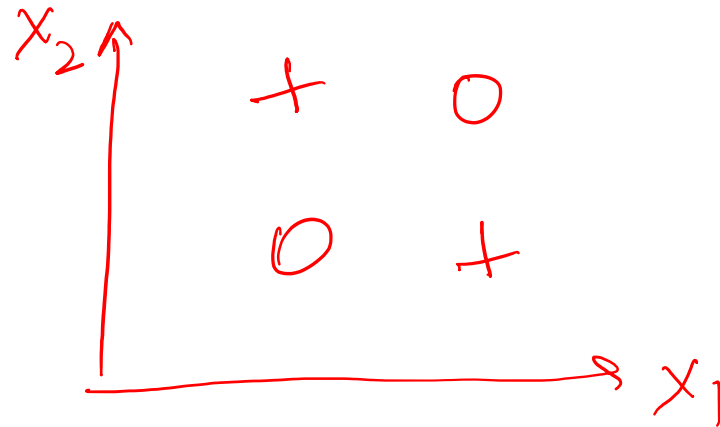
Practice problems (classification)

- To predict the chance that Steelers football team will win the Super Bowl Championship next year, you should prefer to use logistic regression instead of ~~decision trees~~. *k-nearest neighbor*

$$P(Y=1|X)$$

- Which of the following classifiers can perfectly classify the following dataset $\sum_i w_i x_i \geq 0$

- Naïve Bayes *No.*
- Logistic regression *No.*
- Neural network *✓*



$$\frac{e^{\sum w_i x_i}}{1 + e^{\sum w_i x_i}} \geq 1$$

Practice problems (regression)

- Suppose you wish to predict age of a person from his/her brain scan using regression, but you only have 10 subjects and each subject is represented by the brain activity at 20,000 regions in the brain. You would prefer to use least squares regression instead of ridge regression.

$$n = 10 \quad \text{dim} = 20,000$$

- When doing kernel regression on a memory-constrained device, you should prefer to use a box kernel instead of a Gaussian kernel.

