

Course Contents

- 11 lectures
- 1-3 labs (mandatory)
- programming challenge
- written Exam

- ① Nearest Neighbour Classifier (Memory-based)
- ② Decision Trees (Logical inference)
- ③ Challenges in Machine Learning
- ④ Regression
- ⑤ Probabilistic Methods
- ⑥ Learning as Inference
- ⑦ Learning with Latent Variables
- ⑧ Support Vector Machines
- ⑨ Artificial Neural Networks
- ⑩ Ensemble Methods
- ⑪ Dimensionality Reduction
- ⑫ Mini lectures (beyond the scope of DD2421), exam Q&A
- ⑬ Decision Trees
- ⑭ Support Vector Machines
- ⑮ Bayes Classifier & Boosting
- labs are carried out by students and examined by TAs
- use Canvas to book time slots **before the deadlines**

Classification

We would like to enable a computer to learn from data to answer a question "what is it?"
You are given sample data (to find pattern)

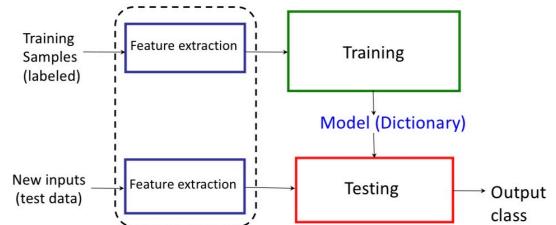
The framework of classification

- 1) Training Phase: To give the concept of classes to a machine using labeled data
- 2.) Testing Phase: To determine the class of new unseen (unlabeled) data

Examination:

- It is **your** task to convince the examiner that you have done the assignment and understood the results.
- Strongly encouraged to work+report by pairs of 2 students (not 3!).
- 10 minutes, **be there on time**.
- No programming code to be shown
- Bring your ID (tell the TA if you are yet to be registered)

Schematic of classification

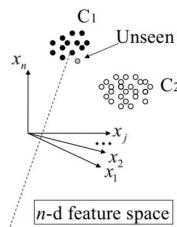


Nearest Neighbour Methods

Nearest Neighbour methods

• Binary classification

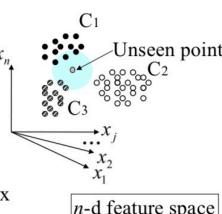
- N_1 samples of class C_1
- N_2 samples of class C_2
- Unseen data x
- Compute distances to all the $N_1 + N_2$ samples



- Find the nearest neighbour
→ classify x to the same class

• k -nearest neighbour rule

- Compute the distances to all the samples from new data x
 - Pick k neighbours that are nearest to x
- Majority vote to classify point x
(Nearest Neighbour is 1-NN)



Pros and cons of k -NN

• k -NN / 1-NN comparison summary

- the boundary becomes smoother as k increases
- lower computational cost for lower k
- k -NN better **generalizes** given many samples

• Pros:

- simple; only with a single parameter k
- applicable to multi-class problems
- good performance, effective in **low dimension data**

• Cons:

- costly to compute distances to search for the nearest
- memory requirement: must store all the training set

Lecture Contents

Decision trees

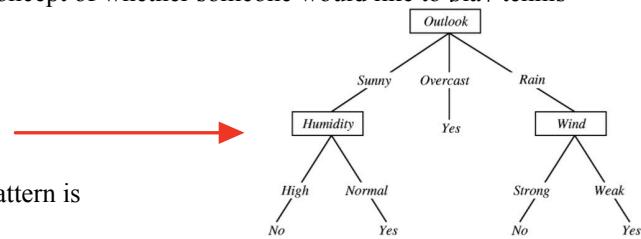
↳ The representation } Right now
↳ Training }

Basic idea: Test attributes (features) sequentially = Ask questions about the target/status sequentially

Example: Building a concept of whether someone would like to play tennis



Each leaf node bears a category label, and the test pattern is assigned the category of the leaf node reached



What does this tree encode?

How do we grow the tree?

(Sunny \wedge Normal Humidity) \vee (Cloudy) \vee (Rainy \wedge Weak Wind)

1. Choose the best question, according to the information gain, and split the input data into subsets
2. Terminate: Call branches with a unique class labels leaves (no need for further questions)
3. Grow: Recursively extend other branches with subset bearing mixtures of labels

Unpredictability

↳ Entropy
↳ Information gain } Right now
↳ Gini impurity

How to measure information gain?

↳ The Shannon information of an outcome is

$$\log_2 \frac{1}{p_i} \rightarrow (p_i: \text{probability for event } i)$$

Example of Entropy

↳ The entropy - Measure of uncertainty (unpredictability)

$$\text{Entropy} = \sum_i -p_i \log_2 p_i$$

is a sensible measure of expected information content

Example: rolling a die
 $p_1 = \frac{1}{6}, p_2 = \frac{1}{6}, \dots, p_6 = \frac{1}{6}$

$$\begin{aligned} \text{Entropy} &= \sum_i -p_i \log_2 p_i = \\ &= 6 \times \left(-\frac{1}{6} \log_2 \frac{1}{6} \right) = \\ &= -\log_2 \frac{1}{6} = \log_2 6 \approx 2.58 \end{aligned}$$

Example: rolling a fake die
 $p_1 = 0.1; \dots, p_5 = 0.1; p_6 = 0.5$



$$\begin{aligned} \text{Entropy} &= \sum_i -p_i \log_2 p_i = \\ &= -5 \cdot 0.1 \log_2 0.1 - 0.5 \log_2 0.5 = \\ &\approx 2.16 \end{aligned}$$

The result of a die-roll has 2.58 bit of information

A real die is more unpredictable (2.58 bit) than a fake (2.16 bit)

Unpredictability of a dataset (think of a subset at a node)

- 100 examples, 42 positive

$$-\frac{58}{100} \log_2 \frac{58}{100} - \frac{42}{100} \log_2 \frac{42}{100} = 0.981$$

- 100 examples, 3 positive

$$-\frac{97}{100} \log_2 \frac{97}{100} - \frac{3}{100} \log_2 \frac{3}{100} = 0.194$$

The result of a coin-toss has 1 bit of information

Example: tossing a coin
 $p_{\text{head}} = 0.5; p_{\text{tail}} = 0.5$



$$\begin{aligned} \text{Entropy} &= \sum_i -p_i \log_2 p_i = \\ &= -0.5 \log_2 0.5 - 0.5 \log_2 0.5 = \\ &= 1 \end{aligned}$$

Back to Decision trees

Smart Idea: Ask about the attribute which maximize the expected reduction of the entropy,

Information gain:

Ask about attribute A for a data set S that has the Entropy $\text{Ent}(S)$, and get subsets S_v according to the value of A

$$\text{Gain} = \underbrace{\text{Ent}(S)}_{\text{Before}} - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Ent}(S_v)$$

Weighted Sum After

What is the entropy of this binary dataset (attributes = $\{A, B, C, D\}$, $n = 25$)?

$$\text{Ent} = -\frac{12}{25} \log_2 \frac{12}{25} - \frac{13}{25} \log_2 \frac{13}{25} \approx 0.9988$$

$$A = \bullet: \frac{3}{6} \text{ positive} \rightarrow 1.0$$

$$A = \circ: \frac{9}{19} \text{ positive} \rightarrow 0.9980$$

$$\text{Expected: } \frac{6}{25} \cdot 1.0 + \frac{19}{25} \cdot 0.9980 \approx 0.9985$$

$$B = \bullet: \frac{9}{11} \text{ positive} \rightarrow 0.684$$

$$B = \circ: \frac{3}{14} \text{ positive} \rightarrow 0.750$$

$$\text{Expected: } 0.721$$

$$C = \bullet: \frac{6}{12} \text{ positive} \rightarrow 1.0$$

$$C = \circ: \frac{6}{13} \text{ positive} \rightarrow 0.9957$$

$$\text{Expected: } 0.9977$$

$$D = \bullet: \frac{3}{5} \text{ positive} \rightarrow 0.9710$$

$$D = \circ: \frac{9}{20} \text{ positive} \rightarrow 0.9928$$

$$\text{Expected: } 0.9884$$

| A | B | C | D | |
|---|---|---|---|---|
| ○ | ● | ● | ○ | + |
| ● | ● | ○ | ○ | + |
| ○ | ○ | ○ | ○ | |
| ○ | ○ | ● | ● | + |
| ○ | ● | ○ | ○ | + |
| ● | ○ | ● | ○ | |
| ○ | ● | ● | ○ | + |
| ○ | ○ | ○ | ○ | |
| ● | ○ | ○ | ○ | |
| ○ | ○ | ● | ○ | + |
| ○ | ● | ● | ○ | + |
| ○ | ○ | ● | ● | |
| ○ | ● | ○ | ○ | + |
| ○ | ○ | ○ | ○ | |
| ● | ● | ● | ● | |
| ○ | ● | ○ | ○ | + |
| ○ | ○ | ● | ○ | |
| ● | ○ | ○ | ○ | |
| ○ | ○ | ● | ● | + |
| ● | ● | ○ | ○ | + |
| ○ | ○ | ○ | ○ | |
| ○ | ○ | ● | ○ | + |

$$\text{Gain}(A) = 0.9988 - 0.9985 = 0.0003$$

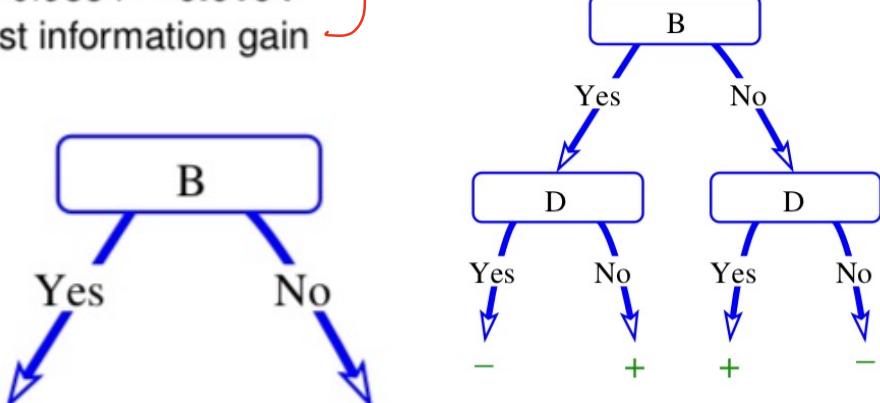
$$\text{Gain}(B) = 0.9988 - 0.7210 = 0.2778$$

$$\text{Gain}(C) = 0.9988 - 0.9977 = 0.0011$$

$$\text{Gain}(D) = 0.9988 - 0.9884 = 0.0104$$

Attribute B gives most information gain

Then D
afterwards



Greedy Approach: Choose the attribute which tells us most about the answer

In sum, we need to find good questions to ask. More than one attribute could be involved in one question.

Gini Impurity

↳ Another definition of predictability (impurity)

$$\sum_i p_i(1-p_i) = 1 - \sum_i p_i^2$$

(p_i : Probability for event i)

↳ The expected error rate at a node N, if the category label is randomly selected from the class distribution present at N

↳ Similar to the Entropy but more strongly peaked at equal probabilities

Overfitting

Overfitting: When the learned models are overly specialised for the training samples

↳ Good results on training data but generalizes poorly

When does this occur:

- Non representative sample
- Noisy examples
- Too complex models

What can be done about it?

Choose a simpler model and accept some errors for the training examples

Separate the available data into two sets of examples

- **Training set:** To form the learned model
- **Validation set:** To evaluate the accuracy of this model

The motivations:

- The training may be misled by random errors but the validation set is unlikely to exhibit the same random fluctuations
- The validation set to provide a safety check against overfitting the spurious characteristics of the training set

Validation Set needs to be large enough to provide statistically meaningful instances

Reduced Error Pruning

Split the data into training and validation set

Do until further pruning is harmful:

- Evaluate impact on validation set of pruning each possible node
- Greedily remove the one that most improves validation set accuracy

Produces smallest version of most accurate subtree

Possible ways of improving/extending the decision trees:

- Avoid overfitting
 - Stop growing when the data split not statistically significant
 - Grow full tree, then post-prune
- A collection of trees
 - Bootstrap aggregating (bagging)
 - Decision Forests

Challenges in Machine Learning

Basic idea

↳ How should we select/determine the right model f from data?

Given training data: $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$

of inputs $x_i \in \mathbb{R}^d$ and their labels y_i

Compute the missclassification rate on D } note: $\text{Ind}(x) = 1$ if $x = \text{TRUE}$ otherwise
 $\text{err}(f, D) = \frac{1}{n} \sum_{i=1}^n \text{Ind}(f(x_i) \neq y_i)$ $\text{Ind}(x) = 0$

Overfitting: Models are overly specialized for the training samples

↳ Good results on training data, but generalizes poorly due to

↳ non-representative sample

↳ Noisy examples

↳ If we are in a data-rich situation:

→ Partition the data into three sets: Training set, Validation set and Test set for assessment of the generalization error of the final chosen model.

Curse of Dimensionality

Image: inputs represented by 30 features but some of them are less relevant to target function.
Will you use all of them?

- Easy problems in low dimensions are harder in high-dimensions
- In High dimensions everything is far from everything else
 - ↳ Issues in Nearest Neighbors
- Any method that attempts to produce locally varying functions in small (stropic neighborhood cell) run into problems in high dimensions

Intuitions in low dimension do not apply in high dimensions. Real world $\rightarrow 3D$ but we deal with data for instance in 1000-d

- Uniform distribution on hypercube
- Volume of hypersphere

Techniques for dimensionality reduction / feature Selection exists

The Bias-Variance Trade-off

Let us imagine we could repeat the modelling for many times - each time by gathering new set of training samples, D

The resulting models will have a range of prediction due to randomness in the underlying data set

- **Bias:** The difference between the average prediction of our model and the correct value
- **Variance:** The variability of a model prediction for a given point between realizations of the model

Let us consider

$f(x)$: True function

$\hat{f}_0(x)$: Prediction function (= model) estimated with D

and a conceptual tool:

$E_D[\hat{f}_0(x)]$: average of models due to different sample sets

The mean square error (MSE) for estimating $f(x)$

$$E_D[(\hat{f}_0(x) - f(x))^2] = E_D[(\hat{f}_0(x) - E[\hat{f}_0(x)])^2] + (E[\hat{f}_0(x)] - f(x))^2$$

Variance + (Bias)²

Characterization of a classifier: Bias

Bias of a classifier is the discrepancy between its average estimated and true function

$$E[\hat{f}_0(x)] - f(x)$$

Regression

Function approximation

- How do we fit this dataset D ?

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

→ of n pairs of inputs x_i and targets $y_i \in \mathbb{R}$ D can be measurements in an experiment

- Task of Regression:

→ to predict target associated to any arbitrary new input

Linear Regression

- Linear regression tries to estimate the function $f(x)$ (x is d-dimensional) and predict output by

$$\hat{f}(x) = \sum_{i=0}^d w_i x_i = w^T x$$

- How to measure the error for N samples:

Mean Squares is used to approximate $\hat{f}(x)$ to $f(x)$

- Mean Square Error

$$E_{in}(\hat{f}) = \frac{1}{N} \sum_{n=1}^N (\hat{f}(x_n) - y_n)^2$$

- Minimizing in-Sample MSE

$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^N (w^T x_n - y_n)^2 = \frac{1}{N} \|Xw - Y\|^2$$

Where: $X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$, $Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

Residual Sum of Squares

- The sum of squared errors is a convex function of w

$$E_{in}(w) = \|Xw - Y\|^2$$

- The gradient with respect to the weights is:

$$\frac{d}{dw} E_{in}(w) = 2X^T(Xw - Y)$$

- The weight vector that sets the gradient to zero minimises the errors

$$X^T Xw = X^T Y$$

$$\rightarrow w = (X^T X)^{-1} X^T Y$$

RANSAC: Random Sampling Consensus

Repeat M times

- Sample two points to estimate the line
- Calculate the number of inliers or posterior likelihood for relation
- Choose relation to maximize the number of inliers
(Last two in Least Median of Square (LMedS))
- Calculate error of all data
- Choose relation to minimize median of error

Ridge Regression

- Similar to least squares but minimizes different quantity:

$$RSS + \lambda \sum_{i=1}^d w_i^2$$

The second term is called Shrinkage penalty)

- Shrinkage Penalty: Small when w_i are close to zero

- The Parameter λ : Control the relative impact of the two terms, the selection is critical!
 - gradient decent
 - Coordinate descent

The Lasso (Least absolute Shrinkage and Selection Operator)

- Similar to ridge regression but with slightly different term:

$$RSS + \lambda \sum_{i=1}^d |w_i|$$

- The lasso could be proven mathematically that some coefficients end up being set to exactly zero

- Variable Selection
- Yield Sparse model

Another formulation

- For every value of λ there is some S such that the equations will give the same coefficients estimate

- Ridge regression: Minimizing $RSS + \lambda \sum_{i=1}^d w_i^2$:

$$\text{Minimizing } RSS, \text{ s.t. } \sum_{i=1}^d w_i^2 \leq S$$

- Lasso: Minimizing $RSS + \lambda \sum_{i=1}^d |w_i|$

$$\text{Minimizing } RSS, \text{ s.t. } \sum_{i=1}^d |w_i| \leq S$$

RANSAC: RANDOM SAMPLING Consensus

Objective

Robust fit of model to data set S which contains outliers

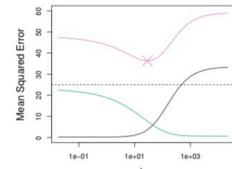
Algorithm

- (i) Randomly select a (minimum number of) sample of s data points from S and instantiate the model from this subset.
- (ii) Determine the set of data points S_i which are within a distance threshold t of the model. The set S_i is the consensus set of samples and defines the inliers of S .
- (iii) If the subset of S_i is greater than some threshold T , re-estimate the model using all the points in S_i and terminate
- (iv) If the size of S_i is less than T , select a new subset and repeat the above.
- (v) After N trials the largest consensus set S_i is selected, and the model is re-estimated using all the points in the subset S_i

(in Hartley and Zisserman, adapted from Fischler '81)

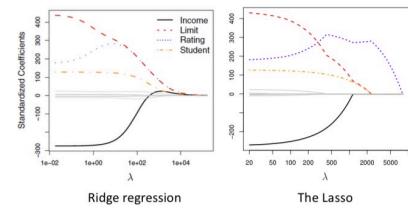
Ridge Regression Bias/Variance

- Green: Variance
- Black: Bias
- Purple: MSE



Increased λ decreases variance while increasing bias

Comparison of estimated coefficients



K-NN Regression (non-Parametric)

- Similar to the k-NN classifier
- To regress Y for a given value of X , consider k closest points to x in training data and take the average of the responses $\Leftrightarrow f(x) = \frac{1}{k} \sum_{x \in N_k} y_i$
- Larger values of k provide a smoother and less variable fit (lower variance)

Parametric or non-Parametric

← Important

- if the parametric form is close to the true form of f , the parametric approach will outperform the non-parametric
- Parametric methods tend to outperform non-parametric when there is a small number of observations per predictor (high dimension)
- Interpretability Stand Point: Linear regression preferred to KNN if the test MSEs are similar or slightly lower.

Summary

- KNN regression

- Ridge regression → Gradient descent, Coordinate descent
- Lasso

- RANSAC

Another formulations

- MSE, RSS

For every value of λ there is some s such that the equations will give the same coefficient estimates:

• Ridge regression: Minimizing $RSS + \lambda \sum_{i=1}^d w_i^2$

Minimizing $RSS, \text{sub.to } \sum_{i=1}^d w_i^2 \leq s$

• Lasso: Minimizing $RSS + \lambda \sum_{i=1}^d |w_i|$

Minimizing $RSS, \text{sub.to } \sum_{i=1}^d |w_i| \leq s$

Probabilistic Reasoning

Probability theory in ML

- ↳ Incorporate probabilistic thinking at all levels
 - We start with incomplete knowledge
 - We then reduce uncertainty by incorporating observations
 - We then propagate belief and update our estimation

Advantages of probability based Methods

- ↳ Interpretability: More transparent and mathematical rigorous than other ML models
- ↳ Transparency: Assumptions can be made more explicit than in other methods
- ↳ Efficiency: Works with poor data
- ↳ Flexibility: Easy to merge different parts of a complex system and update current theories
- ↳ Encompassing: Aspects of learning and inferences can be cast under the same theory

Disadvantages of Probabilistic Methods

Complicated: Often hard to derive closed solutions. Need to resort to computation and heuristic approximations

Scalability: Not computationally scalable to large datasets but many argue that the need for large data sets is a disadvantage

Notes from slides

Probability Theory in ML

- We start with incomplete knowledge
- We then reduce uncertainty by incorporating observation
- We then propagate belief and update our estimate

Probabilistic Machine Learning provides something like a unified theory for ML

Statistic Crash Course

Random Variables:

- Sample space: set of outcomes are not numbers.
- We need to "convert" outcomes to numbers
- Let's map the set of outcome to a set of numbers!

Probability distribution:

- The probability distribution function PDF of a r.v maps its range to positive numbers
 $\Pr(x) : X \rightarrow \mathbb{R} > 0$

- $\Pr(x)$ describes how probability density is distributed over the range of X , e.g. the probability of $0 \leq x \leq 1$ is given by:

$$\Pr[0 \leq x \leq 1] = \int_{x=0}^1 \Pr(x) dx$$

- X is distributed $\Pr(x)$ written more compactly
 $X \sim \Pr(x)$

Independence

- Independence if the joint distribution can be factorised $\Pr(A \cap B) = \Pr(A) \Pr(B)$

$$\Rightarrow \Pr(A \cap B) = \frac{\Pr(A \cap B)}{\Pr(B)} = \frac{\Pr(A) \Pr(B)}{\Pr(B)} = \Pr(A)$$

B happening says nothing about A

- Bayes theorem ← Memorise its derivation

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

Common Distribution

- Bernoulli: Domain: binary variables ($x \in \{0, 1\}$)

Parameter: $\lambda = \Pr[x=1]$, $\lambda \in [0, 1]$

$\Pr[x=0] = 1 - \lambda$ and

$$\Pr(x) = \lambda^x (1-\lambda)^{1-x} = \begin{cases} \lambda, & \text{if } x=1 \\ 1-\lambda, & \text{if } x=0 \end{cases}$$

- Categorical: Domain: discrete variables ($x \in \{x_1, \dots, x_k\}$)

Parameter: $\lambda = [\lambda_1, \dots, \lambda_k]$ $\Pr[x_{ik}=1] = \lambda_{ik}$

with $\lambda_{ik} \in [0, 1]$ and $\sum_{k=1}^K \lambda_{ik} = 1$

- Gaussian distribution: univariate normal distribution

Domain: Real number ($x \in \mathbb{R}$)
 $\Pr(x; \mu, \sigma^2) = N(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$

Central Limit theorem: The distribution of a large number of independent, identically distributed variables will tend to normal distribution, regardless of the underlying distribution.

Gaussian Distribution: D Distribution

- Known as multivariate normal distribution
- Domain: Real numbers ($x \in \mathbb{R}^D$)

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} \quad \nu = \begin{bmatrix} \nu_1 \\ \nu_2 \\ \vdots \\ \nu_D \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1D} \\ \sigma_{21} & \sigma_2^2 & \dots & \vdots \\ \vdots & & \ddots & \\ \sigma_{D1} & \dots & \dots & \sigma_D^2 \end{bmatrix}$$

$$Pr(x; \nu, \Sigma) = \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \cdot E\left[-\frac{1}{2}(x-\nu)^T \Sigma^{-1}(x-\nu)\right]$$

Eigen value decomposition of Σ matrix

$$\Sigma = Q \Lambda Q^T = [U V] \text{diag} [\lambda_1, \lambda_2] [U V]^T$$

Expectations

- $E[x] = \nu_x = \int x \Pr(x) dx \Rightarrow$ the center of gravity of a distribution
- $E[x] = \bar{\nu}_x = \sum x_i \Pr(x_i) dx \quad \text{Sampled Expected Value} = \text{Mean} \quad \bar{\nu}_x = \frac{1}{N} \sum_i^N x_i$

Variance

$$\text{Var}[x] = \sigma_x^2 = E[(x - E[x])^2] \quad \bar{\sigma}_x^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{\nu}_x)^2$$

The spread of a distribution

Covariance

- $\sigma_{x,y} = E[(x - E[x])(y - E[y])] \quad \bar{\sigma}_{xy} = \frac{1}{N-1} \sum_i^N (x_i - \bar{\nu}_x)(y_i - \bar{\nu}_y)$
- $\Sigma_x = E[(x - E[x])(x - E[x])^T]$

\rightarrow Unbiased Sample Co-Variance

Probabilistic Machine Learning

General ML Outline

Data: $\{(x_i, y_i), (x_1, y_1), \dots, (x_n, y_n)\}$

\rightarrow Where x are features (independent variables) and y is the answer (dependent variable)

- $i \in Y$ is discrete, finite: Classification
- if Y is continuous: Regression

Learning: We want to learn how variables are related

- We estimate $\Pr(x, y)$ from observations $\{(x_i, y_i)\}$

Inference: We want to predict an answer given an observation

- We estimate $\Pr(y|X=x)$ from observation $\{(x_i, y_i)\}$

Bayes' Rule

$$\Pr(y | X=x) = \frac{\Pr(x | Y=y) \Pr(Y=y)}{\Pr(X=x)}$$

- $\Pr(x | Y=y) \leftarrow$ Likelihood represents the probability density of observing data x given the hypothesis $Y=y$.
- $\Pr(Y=y) \leftarrow$ Prior represents the knowledge about Y before any observation.
- $\Pr(y | X=x) \leftarrow$ Posterior represents the probability density of hypothesis y given observation $X=x$.
- $\Pr(X=x) \leftarrow$ Evidence describes how well the model fits the evidence.

$$\Pr(X=x) = \begin{cases} \sum_y \Pr(x | Y=y) \Pr(Y=y) & \text{classification} \\ \int_y \Pr(x | Y=y) \Pr(Y=y) & \text{regression} \end{cases}$$

Probabilistic Regression

- Regression via Conditional Probability
- Compute Posterior of Y : $\Pr(Y|X=x) = \Pr(X=y)/\Pr(X=x)$
- Compute Conditional Expectation: $E[Y|X=x]$

Explicit Regression Model:

- Define a deterministic model $y = f(x) + \epsilon$
- Define probability distribution of the error $= \epsilon = Y - f(x)$
- Explicit parameters in $f(x)$

SKIPPING

Most Probable Hypothesis

• Maximum A posterior (MAP)

choose hypothesis from Y with the highest probability given observed data x :

$$\rightarrow Y_{\text{map}}(x) = \arg \max_{y \in Y} \Pr(Y=y|X=x)$$

$$\arg \max_{y \in Y} \frac{\Pr(X=x|Y=y) \Pr(Y=y)}{\Pr(X=x)}$$

$$\arg \max_{y \in Y} \Pr(X=x|Y=y) \Pr(Y=y)$$

• Maximum Likelihood (ML):

If we do not know prior distribution, then choose hypothesis with highest likelihood of generating the observed data:

$$\rightarrow Y_{\text{mlp}} = \arg \max_{y \in Y} \Pr(X=x|Y=y)$$

Before I took a COVID test, the doctor said 99% of the people in the area have COVID, and 90% of those with COVID are testing positive. A few days later the doctor called and said my test was positive, and that the probability I have COVID given this positive test is $p\%$ — I can't remember because I was in shock. Find the minimum value of p such that I can compute the probability I got a positive test but don't have COVID, and then compute the maximum probability I don't have COVID given my positive test.
Show your work.

Lecture 6 - Learning as Inference

Probabilistic Classification and Regression

$$\Pr(Y|X=x) = \frac{\Pr(x|Y=y) \Pr(Y=y)}{\Pr(x=x)}$$

Likelihood Prior
Pr(x=x) Evidence

- Classification: Y is discrete, finite
- Regression: Y is continuous

How do we obtain these distributions

Given:

- the training data $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$
- New observation x

Estimate the posterior probability of y $\Pr(y|x, D)$

Discriminative vs Generative Model

Discriminative model:

- Models $\Pr(Y, x, D)$ directly \rightarrow calculates prior directly
- Logistic regression

Generative Modelling

- Models $\Pr(x|y, D)$
- Naïves Bayes

Parametric vs Non-Parametric Inference

$$\Pr(Y|x) = \Pr(Y|x, \theta)$$

The distribution is characterized by parameters θ

Parametric Inference

- Estimate θ using D
- Compute $\Pr(Y|x, \hat{\theta})$ to make inference

Learning corresponds to estimating θ

Non-Parametric Inference

- Estimate $\Pr(\theta, D)$
- Compute $\Pr(Y|x, D)$ from $\Pr(Y|x, \theta, D)$ $\Pr(\theta|D)$ by marginalizing out θ
- grows with Data

Fundamental Assumption iid

- Observations are independent and identically distributed

$$D = \{o_1, o_2, \dots, o_n\} \quad o_i = (x_i, y_i)$$

- The likelihood of the whole data set can be factorised

$$\Pr(D) = \Pr(o_1, o_2, \dots, o_n) = \prod_{i=1}^n \Pr(o_i)$$

- Taking the log creates the log-likelihood

$$\text{Log } \Pr(D) = \sum_{i=1}^n \text{Log } \Pr(o_i)$$

Maximum Likelihood Estimate (ML)

$$\Pr(x|y) \equiv \Pr(x|y, \theta) \quad \text{or} \quad \Pr(y|x) \equiv \Pr(y|x, \theta)$$

Find the parameter values that makes the data most likely

- ML Optimality is defined as maximizing the likelihood of D:

$$\theta_{ML} = \arg \max_{\theta} P(D, \theta) = \arg \max_{\theta} \log P(D|\theta) \quad D = \text{Distribution} \quad \theta = \text{Maximize this}$$

- We can then approximate distributions given the data:

$$\Pr(x|y, D) \approx \Pr(x|y, \theta_{ML}) \quad \text{or} \quad \Pr(y|x, D) \approx \Pr(y|x, \theta_{ML})$$

Lecture 7 - Prior and Latent Variable

Check your understanding!

Goal of lecture: Answer both

Consider a dataset $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$ of N independent and identically distributed observations where each x_n is a p -dimensional real vector. Assume the random variable Y_n is distributed Laplacian with a mean $\beta^T x_n$ and known scale parameter $b > 0$. In other words, $Y_n | x_n, \beta \sim \mathcal{L}(\beta^T x_n, b)$. Define the a priori distribution of the parameters $\beta = (\beta_1, \beta_2, \dots, \beta_p)$ multivariate Gaussian with parameters mean 0 and variance $\sigma^2 \mathbf{I}$.

- ① Derive the maximum likelihood (ML) estimate of β .
- ② Derive the maximum a posteriori (MAP) estimate of β .

$$\text{Information given: } \Pr(Y_n | X_n, \beta) = \frac{1}{2b} \exp\left[-\frac{|Y_n - \beta^T x_n|}{b}\right]$$

$$\Pr(\beta) = N(0, \sigma^2 \mathbf{I})$$

Maximum likelihood (ML) estimate:

- $\Pr(Y_n | X_n, \beta) = \frac{1}{2b} \exp\left[-\frac{|Y_n - \beta^T x_n|}{b}\right]$

$$\Pr(\beta) = N(0, \sigma^2 \mathbf{I})$$

Log likelihood

- $\log \Pr(D|\beta) = \sum_{n=1}^N \log \Pr(Y_n | X_n, \beta)$

Maximum likelihood estimate is then

$$\begin{aligned} \beta_{\text{ML}} &= \underset{\beta}{\operatorname{argmax}} \log \Pr(D|\beta) = \underset{\beta}{\operatorname{argmax}} \sum_{n=1}^N \log \Pr(Y_n | X_n, \beta) \\ &= \underset{\beta}{\operatorname{argmax}} \sum_{n=1}^N \log \left(\frac{1}{2b} \exp\left[-\frac{|Y_n - \beta^T x_n|}{b}\right] \right) \\ &= \underset{\beta}{\operatorname{argmax}} \sum_{n=1}^N -\log(2b) - \frac{1}{b} |Y_n - \beta^T x_n| \\ &= \underset{\beta}{\operatorname{argmin}} N \log(2b) + \frac{1}{b} \sum_{n=1}^N |Y_n - \beta^T x_n| \\ &= \underset{\beta}{\operatorname{argmin}} \|Y - \beta^T X\|_1 \quad \leftarrow \text{L1 Sparsing the solution} \end{aligned}$$

Where $Y = [y_1, y_2, \dots, y_N]^T$ and $X = [x_1, x_2, \dots, x_N]^T$

Maximum a Posterior estimation

- ML estimates θ to maximize probability of D

- MAP Estimation chooses the most likely θ given D

$$\theta = \underset{\theta}{\operatorname{argmax}} \Pr(\theta | D)$$

$$= \underset{\theta}{\operatorname{argmax}} \frac{\Pr(\theta) \Pr(D|\theta)}{P(D)}$$

$$= \underset{\theta}{\operatorname{argmax}} \Pr(\theta) \Pr(D|\theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \left[\Pr(\theta) \sum_{n=1}^N \Pr(x_n | \theta) \right] =$$

$$\left. \begin{aligned} &\rightarrow \underset{\theta}{\operatorname{argmax}} \left[\log \Pr(\theta) + \sum_{n=1}^N \log \Pr(x_n | \theta) \right] \\ &\rightarrow \text{Ridge Regression} \rightarrow \text{Fit the data and keep the weights small} \end{aligned} \right\}$$

- 1) argmax
- 2) Bayes theorem without nominator
- 3) Add dog
- 4) $\sum_{n=1}^N$ and start logging
- 5) Arg min

Impact of different Priors

$$\hat{\omega}(\lambda) = \arg \min_{\omega} \left[\sum_{n=1}^N \underbrace{(y_n - \omega^T x_n)^2}_{\text{Fit data}} + \lambda \sum_{d=1}^D |\omega_d|^p \right]$$

- $\lambda = 0$ or $\omega_d \sim \text{Uniform}$ \rightarrow Fit the data
- $\lambda > 0$ \rightarrow Fit the data and make ω short (Ridge Regression) \leftarrow assuming Normal
- $\lambda > 0$ \rightarrow Fit the data and make ω sparse (Lasso) \leftarrow assuming Laplace

Check your understanding!

Consider a dataset $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$ of N independent and identically distributed observations where each \mathbf{x}_n is a p -dimensional real vector. Assume the random variable Y_n is distributed Laplacian with a mean $\beta^T \mathbf{x}_n$ and known scale parameter $b > 0$. In other words, $Y_n | \mathbf{x}_n, \beta \sim \mathcal{L}(\beta^T \mathbf{x}_n, b)$. Define the a priori distribution of the parameters $\beta = (\beta_1, \beta_2, \dots, \beta_p)$ multivariate Gaussian with parameters mean $\mathbf{0}$ and variance $\sigma^2 \mathbf{I}$.

- Derive the maximum likelihood (ML) estimate of β .
- Derive the maximum a posteriori (MAP) estimate of β .

Derive the maximum a posteriori (MAP) estimate of β .

- $Pr(y_n | \mathbf{x}_n, \beta) = \frac{1}{2b} \exp \left[-\frac{|y_n - \beta^T \mathbf{x}_n|}{b} \right]$
- $P(\beta) = \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$

The MAP estimate of β is defined

$$\begin{aligned} \beta_{\text{MAP}} &= \arg \max_{\beta} Pr(\beta | \mathcal{D}) = \arg \max_{\beta} Pr(\mathcal{D} | \beta) Pr(\beta) \\ &= \arg \max_{\beta} \log Pr(\beta) + \sum_{n=1}^N \log Pr(y_n | \mathbf{x}_n, \beta) \\ &= \arg \min_{\beta} \frac{1}{2\sigma^2} \|\beta\|_2^2 + \frac{1}{b} \|\mathbf{y} - \mathbf{X}^T \beta\|_1 \end{aligned}$$

$$\begin{array}{lll} \text{ML: } & \mathcal{D} & \rightarrow \theta_{\text{ML}} \rightarrow Pr(y | \mathbf{x}, \theta_{\text{ML}}) \\ \text{MAP: } & \mathcal{D}, \text{Pr}(\theta) & \rightarrow \theta_{\text{MAP}} \rightarrow Pr(y | \mathbf{x}, \theta_{\text{MAP}}) \\ \text{Bayes: } & \mathcal{D}, \text{Pr}(\theta) & \rightarrow \text{Pr}(\theta | \mathcal{D}) \rightarrow Pr(y | \mathbf{x}, \mathcal{D}) \end{array}$$

- ➊ consider θ as a random variable (same as MAP)
- ➋ characterize θ with the posterior distribution $\text{Pr}(\theta | \mathcal{D})$ given the data
- ➌ compute new predicting posterior $Pr(y | \mathbf{x}, \mathcal{D})$ marginalizing over θ (predictive posterior)

$$Pr(y | \mathbf{x}, \mathcal{D}) = \int_{\theta \in \Theta} Pr(y | \mathbf{x}, \theta) \text{Pr}(\theta | \mathcal{D}) d\theta$$

Occam Razor

Choose the simplest explanation for the data

- number of model Parameters
- number of data Points
- Model fit to the data

AnSupervised Learning

clustering vs Classification

- We don't know their classes

K-means Algorithms

- Describes each class with a centroid
- favors spherical classes
- Euclidean distance is isotropic

If no class independence \rightarrow Expectation Maximization

$$\Pr(x|\theta) = \sum_{k=1}^K \pi_k \Pr(x|\theta_k)$$

(with $\Theta = \{\pi_1, \dots, \pi_K, \theta_1, \dots, \theta_K\}$)

EM for two Gaussians

For each sample x_i introduce a *hidden variable* h_i

$$h_i = \begin{cases} 1 & \text{if sample } x_i \text{ was drawn from } \mathcal{N}(x|\mu_1, \sigma_1^2) \\ 2 & \text{if sample } x_i \text{ was drawn from } \mathcal{N}(x|\mu_2, \sigma_2^2) \end{cases}$$

and come up with initial values

$$\Theta^{(0)} = (\pi_1^{(0)}, \mu_1^{(0)}, \sigma_1^{(0)}, \mu_2^{(0)}, \sigma_2^{(0)})$$

for each of the parameters.

EM is an *iterative algorithm* which updates $\Theta^{(t)}$ using the following two steps...

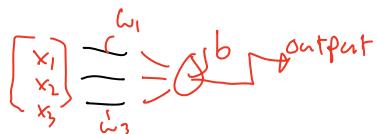
Lecture 7 - SVMs

Concept learning: Supervised learning of Boolean valued functions.
Learn from positive and negative examples to classify yes/no corrections

Examples of concepts: Is there a dog present, mammal, vehicle and so on. **Input is an array of attribute of values.** You could use abstract concepts such as "is there a crime committed here" and so on.

Artificial
Neural
Neuron

Gets a input in vector format. Weights in vector format. B is the threshold for it to output either yes or no. And there's a output.



$$y = \text{sign} \left(\sum_i x_i w_i - b \right)$$

Training
ANNs

- Find the best weights and biases. Done by Delta Rule and Perception learning

Perception
Learning ↗

Perceptron Learning (Binary output):

- Incremental Learning
- Weights only change when the output is wrong
- Update rule
- Always converge if the problem is solvable

} normal of the
wrong

2) Delta
Rule

Delta rule (LMS, continuous output)

- Incremental Learning
- Weights always change
- Converges only in the mean
- Will find an optimal solution even if the problem cannot be fully solved.

Structural Risk Minimization

Margins: A safe distance d from the margin to the datapoints. Wide margins restrict the possible hyperlanes to choose from. Less risk to choose a bad hyperplane. Goal is to find the largest margin from all data points. This sill most likely generalise the best.

$$\vec{w}^\top \vec{x} = 0$$

↳ Separating Hyperplane

$$\vec{w}^\top \vec{x} \geq 1$$

↳ Hyperplane with a Margin
 $t=1$ ↳ Positive target
 $t=-1$ ↳ Negative target

$$t\vec{w}^\top \vec{x} \geq 1$$

↳ Combined

How wide is the Margin

1) Select 2 points \vec{p}, \vec{q} on the two margins:
 $\vec{w}^\top \vec{p} = 1$ $\vec{w}^\top \vec{q} = -1$

2) Distance between \vec{p} and \vec{q} along \vec{w}

$$2d = \frac{\vec{w}^\top (\vec{p} - \vec{q})}{\|\vec{w}\|}$$

→ Support Vectors

- Moving a low dimensional data to a high-dimensional makes stuff linearly separable.
 ↳ Two problems
 - Many free parameters → Bad generalization
 - Extensive Computations.

Ensemble Learning

- What makes a crowd wise?

- ↳ Diversity of opinion:
- ↳ Independence
- ↳ Aggregation
- ↳ Decentralization

We will exploit Wisdom of crowd ideas for specific task by combining classifiers and aim to combine independence and diverse classifiers.

Bagging

- Use Bootstrap replicates of training set by Sampling with replacement.
On Each replicate, learn one Model

High variance: Classifiers produce differing decision boundaries which are highly dependent on the training data

Low bias: Classifiers produce decision boundaries which are on average good approximations to the true decision boundary.

Ensemble prediction using diverse high-variance, low bias classifiers reduce the total variance of the ensemble classifier.

Ensemble Method / Random Forests

- ↳ Randomness?
 - Sampling training data (Same as bagging)
 - Feature Selection at each node
- ↳ Trees are less correlated in example higher variance between weak learners
 - A classifier suited to multi-class problem'

Ensemble Method / Boosting

- ↳ You apply weights to certain data points. Thereafter apply learner to weighted samples. Increase weights of misclassified examples.

Ensemble Method / Dropout

- ↳ Similar to decision trees cutting out leaves/nodes. Here you cut out parts of your neural network.

Ensemble / Summary

- Combine many (high bias) weak classifier/regressors into a strong classifier; boosting
 - ↳ if weak classifiers are chosen and combined using knowledge of how well they and others performed on the task on training data
 - ↳ The selection and combination encourages the weak classifiers to be complementary and de-correlated