

AML CLASS DISCUSSIONS

September 26, 2019

Bayes rule

$$P(c|x) = \frac{P(x|c) P(c)}{P(x)} \rightarrow P(c|x) \underset{\substack{\uparrow \\ P(c|x) \\ P(x)}}{=} \sum_c P(x,c)$$

Example Question.

① Data consisting records of the form (x_1, x_2, x_3, c)

How to use Naive Bayes classification model approach for obtaining a probability of class c given x_1 and x_2 , stating assumptions.

Show that max. likelihood estimate of $P(c=1)$ is given by the proportion of times the class att. c is 1 in the data.

model. $P(D)$.

$$D = \{(x_{1i}, x_{2i}, x_{3i}, c_i) \mid i \in 1, \dots, N_D\}$$

IID

Assumption: data is (conditioned on the model) is independent and identically distributed for each data item

→ each data item is assumed to be from the data model, and is independent of each other.

$$\text{Therefore, } P(D) = \prod_{i=1}^{N_D} P(x_{1i}, x_{2i}, x_{3i} \mid c_i) P(c_i)$$

$$P(D) = \prod_{i=1}^{N_D} \left[\prod_{j=1}^3 P(x_{ji} \mid c_i) \right] P(c_i)$$

Assumption: attrs are conditionally idpd. given c .

Approach (NB algorithm)

① Estimate max. likelihood parameters for probs in the model.

$$P(C|X_1, X_2) = \frac{p(X_1|C) p(X_2|C) P(C)}{P(X_1, X_2)}$$

$$\text{where } P(X_1, X_2) = p(X_1|C=1) p(X_2|C=1) P(C=1) + p(X_1|C=0) p(X_2|C=0) P(C=0)$$

② compute $P(C|X_1, X_2)$ for test values X_1, X_2 .

③ choose highest probability class for the label.

* NB is very good for missing data because of the independence assumption.

② Training data

$$(1, 0, \sqrt{3.5}, 1), (0, 1, 0, 1), (0, 1, -2, 0), (1, 0, 2, 0), \\ (1, 1, -\sqrt{3.5}, 1).$$

2. Using NB, what is $p(C=1 | X_1=1, X_2=1)$?

$$P(C=1 | X_1=1, X_2=1) = \frac{p(X_1=1 | C=1) p(X_2=1 | C=1) P(C=1)}{P(X_1=1, X_2=1)}$$

$$\downarrow \\ p(X_1=1 | C=1) p(X_2=1 | C=1) P(C=1) + p(X_1=1 | C=0) \\ p(X_2=1 | C=0) P(C=0)$$

$$P(C=1) = 3/5 \quad P(C=0) = 2/5$$

does not
have to
add up
to 1

$$\left\{ \begin{array}{ll} p(X_1=1 | C=1) = 2/3 & p(X_2=1 | C=1) = 2/3 \\ p(X_1=1 | C=0) = 1/2 & p(X_2=1 | C=0) = 1/2 \end{array} \right.$$

$$P(C=1 | X_1=1, X_2=1) = \frac{\frac{2}{3} \times \frac{2}{3} \times \frac{3}{5}}{\frac{2}{3} \times \frac{2}{3} \times \frac{3}{5} + \frac{1}{2} \times \frac{1}{2} \times \frac{2}{5}} = \frac{8}{11} \approx 0.72$$

③ Discuss problems if training data:

$$(1, 0, \sqrt{35}, 1), (1, 1, 0, 1), (9, 1, -2, 0), (1, 0, 2, 0), \\ (1, 1, -\sqrt{35}, 1)$$

Solution?

this item changed from 0 to 1

$$P(X_1=1 | C=1) = 3 \quad \text{and} \quad P(X_1=0 | C=1) = 0$$

↓
count is 0

Have prior to all numbers \rightarrow add a 1 to all the counts to normalise.

You prefer probabilities that are central, not extreme.

④ How to include X_3 in the model?

X_3 is real valued, so model using Gaussian dist.

$$\mu_{X_3} = 0 \quad \sigma_{X_3}^2 = \frac{27}{3} = 9$$

$$P(X_3 | C=1) = N(x; 0, 9)$$

$$P(X_3 | C=0) = N(x; 0, 4) \quad x-\mu = 0$$

$$x_3 = 0. \quad P(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \cdot \frac{(x-\mu)^2}{\sigma^2}\right) = \frac{1}{\sqrt{2\pi \cdot 9}} \quad \downarrow$$

$$P(X_3=0 | C=1) = \frac{1}{\sqrt{2\pi \cdot 9}} = \frac{1}{\sqrt{2\pi}} \cdot \frac{1}{3} = \frac{2}{5} \cdot \frac{1}{3} = \frac{2}{15}$$

⑤ What happens if we have outliers?

- We have probabilities on the 'tail' of the Gaussians
- We end up multiplying all our values with this really small number that's at the tail of the Gaussian
- Can affect it greatly \rightarrow mixing discrete + real

September 30, 2019

Decision Tree - allows splitting things to quadrants
→ linear DB doesn't allow you to do this

Entropy - amount of uncertainty

$$H = - \sum_{\text{node}} p(\text{node}) \sum_{\text{label}} p(\text{label} | \text{node}) \log p(\text{label} | \text{node})$$

- additive down the tree
 - adding a new node allows us to compute the reduction in entropy locally for each decision.

Sensible modelling - entropy on a node on the training set can only reduce when we split.
→ less uncertainty.

October 3, 2019

Why is optimising likelihood NOT ALWAYS good idea?

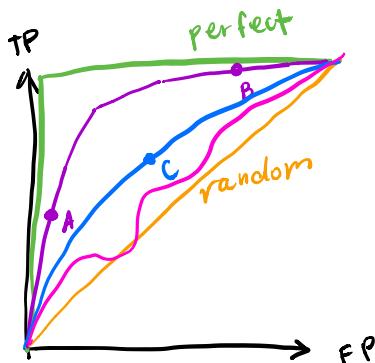
↓ this is the problem

- adding extra parameters to a model \propto better E_{train} .
 - more parameters = more ways to fit same data
 - you want tighter fit!
- optimisation picks one very good solution
 - may be many good sols.
 - uncertainty in the model
- presence of noise → overfitting, low E_{tr} , v high E_{test}

Improving generalisation

- go Bayesian
 - use very simple models → reduce overfitting
 - but they might not catch some functions we might care about
 - choose models w/ diff. complexities to min. overfitting
 - compare them → cross-validation
 - penalise flexible models, so only use the flexibility if they really have to → regularisation with a parameter
 - simpler model might be better than using complex + regularisation
 - muffling sound vs. turning the volume down
- can't use max. likelihood as always prefers more complex models.
→ hence cross-validation

ROC curves



- a. A/B/C more important in a screening programme for a rare illness over a large population?
 ↓
 positive cases are rare
 Low false positive rate.
 Find the gradient.

- b. Making same bet one one of two options in a prediction market?
 A and B

- c. Predict whether a safety measure prevents injury or death.
 High true positive rate. → B
 False positives are okay.
 - minor inconvenience.

- d. Is the pink line a valid ROC curve?
 Yes.
 However this is bad.
 Norm = have something convex.

Exam Questions.

- ① Describe a method to regularise NB to prevent overfitting.
 Any similar method for DT?

Gaussian NB → add threshold to the variances of the Gaussian
 → not allowing G_s to 'blow up'
 → add a soft thing
 → penalise variances so things are not too large or too small
 → stop things to fit too tight
 → tight variance → good for training
 → terrible for unseen data
 → adding counts (everything on one class, none on other)
 e.g. $P(x_1=1 | c=1) = 5/5$ $P(y_1=1 | c=0) = 0/5$
 → move everything to the middle inst. of extremes.

DT → adding additional counts to nodes
 → pulls things closer to $1/2$

Regularisation → more data, the more data overwhelms regularisation
 → negligible in infinite data

- ② cross-val. to choose min. no of dp per leaf in DT
 $2, 3, \dots k$ folds.
 Build a pruned DT.
 → stop when the data items K in a leaf node.
- Across folds set K.
 Pick the ones with the best average performance.
- ③ Contingency tables
- | | |
|------------|----------|
| $TP = 400$ | $FN = 0$ |
| $FP = 60$ | $TN = 0$ |
- "I've got a big bucket.
 Everything goes to
 the positive bucket."
- | | |
|------------|------------|
| $TP = 200$ | $FN = 200$ |
| $FP = 20$ | $TN = 40$ |
- huge overlaps in negatives
 ↓
 more useful
 useful in identifying a region of negatives.
- low acc.
 but gives information

October 7, 2019

Linear Regression

Example Question

- * consider a linear regression from vectors $x \in \mathbb{R}^d$ to scalars y
- $$y = w^T \phi(x) + \epsilon \rightarrow \text{Gaussian noise w/ zero mean and } \sigma^2.$$
- column vector of parameters
- $$\phi(x) = (\phi_1(x), \phi_2(x), \dots)^T$$
- where each ϕ_i is a feature mapping points in \mathbb{R}^d to scalars.
- Given training dataset $P = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$, the loss L
- $$L(w) = \frac{1}{2} \sum_i (y_i - w^T \phi(x_i))^2.$$

- ① Show that minimising $L(w)$ is maximising likelihood, irrespective of the fixed choice for the value σ ?

$$\epsilon \sim N(0, \sigma^2)$$

$$y \sim N(w^T \phi(x), \sigma^2)$$

$$P(y) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{1}{2} \frac{(y_i - w^T \phi(x))^2}{\sigma^2} \right]$$

What if I have all the datapoints? $x_1, y_1 \dots x_N, y_N$

$$P(D|y) = \prod_i \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - w^\top \phi(x_i))^2}{2\sigma^2}\right)$$

$$-\log P(D|y) = \sum_i \frac{(y_i - w^\top \phi(x_i))^2}{2\sigma^2} + \log \sqrt{2\pi\sigma^2}$$

minimise loss = maximise likelihood

\approx no w here
(constant)

$$-\log P(D|y) = \sum_i (y_i - w^\top \phi(x_i))^2 \cdot \frac{1}{2\sigma^2} + \log \sqrt{2\pi\sigma^2}$$

* finding minima is the same under scalar multiplication and addition.

* we don't need to know the noise to find optima!

② By computing derivatives show that the parameter vector that minimises the training is given by

$$\underline{w} = (\phi^\top \phi)^{-1} \phi^\top \underline{y}$$

where $\underline{y} = (y_1, y_2, \dots, y_N)^\top$ and ϕ is the design matrix $\phi^\top = (\phi(x_1), \dots, \phi(x_N))$.

$$L(w) = \frac{1}{2} \sum_i (y_i - \sum_j w_j \phi_{ij})^2$$

* $\phi \rightarrow$ row corresponds to different items
column corresponds to the features (i)
(j)

Calculate derivatives

$$\begin{aligned} \frac{\partial L}{\partial w_k} &= \frac{1}{2} \sum_i (y_i - \sum_j w_j \phi_{ij}) \cdot \frac{d(y_i - \sum_j w_j \phi_{ij})}{d w_k} \\ &= -\sum_i (y_i - \sum_j w_j \phi_{ij}) (\phi_{ik}) \quad + \frac{d(y_k - w_k \phi_{ik})}{d w_k} = \phi_{ik} \\ &= -\phi^\top \underline{y} + \phi^\top \phi \underline{w} \end{aligned}$$

set derivatives to 0

$$\phi^\top \underline{y} = \phi^\top \phi \underline{w}$$

$\phi^\top \phi$ is square matrix (assumption)

$$\underline{w} = (\phi^\top \phi)^{-1} \phi^\top \underline{y}$$

Hence proved!

③ Assumptions about the rel. between M and N using

$$\underline{w} = (\phi^\top \phi)^{-1} \phi^\top \underline{y}$$

$\phi^\top \phi$ would be a square.

N would be large.

- more data the better!

④ Is LR sensitive to outliers?

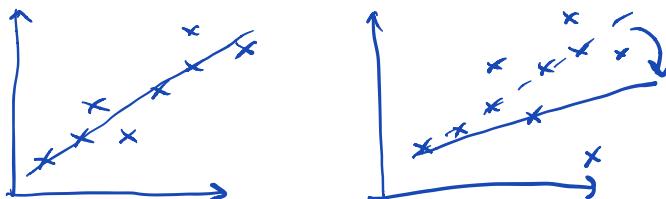
value that is far from the line \rightarrow dramatic effect on the reg'n.
optimisation of $L(\underline{w})$ v sensitive to values on the tail of G .
line deviates dramatically.

Why?

$$L(\underline{w}) = \frac{1}{2} \sum_i (y_i - w^\top \phi(x_i))^2$$

squared of the distance.

Result: adjust whole line to compensate.



⑤ Model accident rate on a road as a fn of traffic density, outside air temp, date, time.

Discuss what features ϕ you might choose and why.

- More cars \rightarrow more accidents, even tho acc. rate constant
- speeds go down with high density
- temperature \rightarrow special value at 0 (frozen)
- date/time \rightarrow rush hour, weekends vs. weekdays

October 10, 2019

Logistic Regression

Generative vs. conditional

- Naive Bayes: $P(c) P(x|c) \rightarrow P(c|x)$
- Log. reg.: $P(x) P(c|x) \rightarrow P(c|x)$ \rightarrow smooth, infinitely differentiable
→ implicitly has understanding of $P(x)$.
→ benefit: map $x \rightarrow c$ is high dim to low dim.
→ easier to build a flexible map.
- * Either is fine when modeling a stationary system.
If things might change, consider causal structure. → Robust!
- Need to know why things change

Logistic Regression

- Lin Reg + Log. trick

push linear outputs through a squashing function to get probs.

- Binary case $P(c=1|x) = \sigma(f(x))$

$$\sigma(x) = \frac{\exp(x)}{1 + \exp(x)} = \frac{1}{1 + \exp(-x)}$$

- Consequence $P(c=0|x) = 1 - \sigma(f(x))$

→ no closed form solution

→ instead convex optimization: optimal solution exists

Loss

$$P(c|x) = [\sigma(w^T x + b)]^c \times [1 - \sigma(w^T x + b)]^{1-c}$$

$$L(w, b) = -\log P(c|x) = -c \log \sigma(w^T x + b) - (1-c) \log (1 - \sigma(w^T x + b))$$

↳ individual single data item loss

$$* \text{Important info } \frac{\partial \sigma}{\partial x} \sigma(x) = \sigma(x)(1 - \sigma(x))$$

Dervatives

\uparrow negative log probability of the whole data

$$L(w, b) = \sum_i -c_i \log \sigma(w^T x_i + b) - (1-c_i) \log (1 - \sigma(w^T x_i + b))$$

* maximise log probability, minimise loss (\Leftrightarrow log prob)
→ go down hill, not uphill.

$$\begin{aligned}
 \frac{\partial L(w, b)}{\partial w_j} &= \sum_i -c_i \left[\frac{\sigma(\underline{w}^T \underline{x}_i + b)(1 - \sigma(\underline{w}^T \underline{x}_i + b))}{\sigma(\underline{w}^T \underline{x}_i + b)} x_{ij} \right] \\
 &\quad - (1 - c_i) \left[\frac{-\sigma(\underline{w}^T \underline{x}_i + b)(1 - \sigma(\underline{w}^T \underline{x}_i + b))}{(1 - \sigma(\underline{w}^T \underline{x}_i + b))} x_{ij} \right] \\
 &= -\sum_i c_i (1 - \sigma(\underline{w}^T \underline{x}_i + b)) x_{ij} + (1 - c_i) \sigma(\underline{w}^T \underline{x}_i + b) x_{ij} \\
 &= \sum_i (c_i - \sigma(\underline{w}^T \underline{x}_i + b)) x_{ij}
 \end{aligned}$$

Multivariate distributions \rightarrow output for each class

$$P(y=k | \underline{x}) = \text{softmax } (f) = \frac{\exp(f_k(\underline{x}))}{\sum_k \exp(f_k(\underline{x}))} \rightarrow \text{to normalize things}$$

- Log. Reg. can be used for more than 1 class.
 \hookrightarrow Logistic sigmoid $\sigma(z) = \frac{1}{1 + \exp(-z)}$

Binomial and Multinomial

- replace bernoulli and multivariate w/ binomial and multinomial

Features

- can replace \underline{x} w/ features $\phi(\underline{x})$
- compute features ahead of time
- there is curse of dimensionality

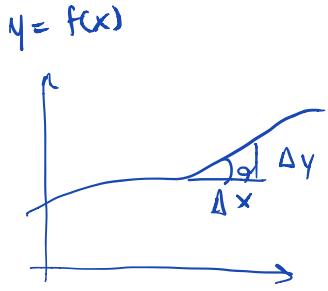
Radial Features

- \hookrightarrow distance from some reference points
- methods
 - distance from key locations
 - use each data as radial feature - non param.
 - grows when data grows
- $\begin{cases} \oplus & \text{very flexible} \\ \ominus & \text{scaling effects} \end{cases}$
- things can only happen locally
 - rent/buy w.r.t income.
 - kind of nearest neighbor
 - RBF confused with really different data
 - extremes should just be similar to normal data.

October 14, 2019

Optimisation

General optimisation



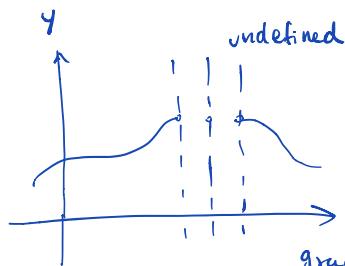
$$\alpha = \tan\left(\frac{\Delta y}{\Delta x}\right)$$

if $\Delta y \oplus$ and $\Delta x \oplus$, increasing fn.
 $\Delta y \ominus$ and $\Delta x \ominus$, decreasing fn.

$$\frac{\Delta y}{\Delta x} = \tan(x)$$

$$\lim_{\Delta x \rightarrow 0} \frac{\Delta y}{\Delta x} = \frac{df}{dx} = f'(x)$$

tells us how
fast we are
incr/dec.



NO slope at undefined regions
 \rightarrow can trip up optimization

Generalise $\rightarrow \nabla f = \left\langle \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_N} \right\rangle^T$

$f = E(\underline{w})$ \rightarrow error of the weights
 step towards the decrease of that error.

$$\underline{w}^{t+1} = \underline{w}^t - \eta \nabla E(\underline{w})$$

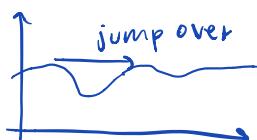
↓ learning rate + greedy

$$E(\underline{w}) = \underset{\text{MSE}}{\frac{1}{n}} \sum_i (y_i - f(x_i))^2 \longrightarrow \text{convex!}$$

Problem 1



Problem 2



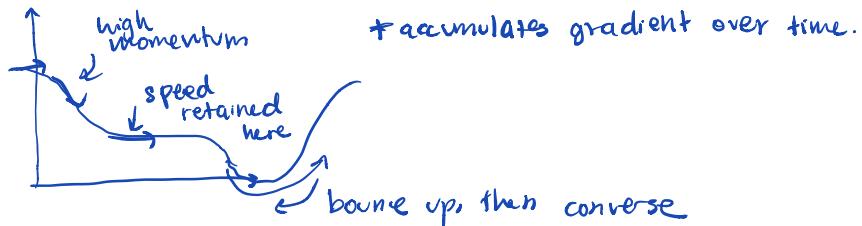
$$\underline{w}^{t+1} = \underline{w}^t - \eta \nabla E(\underline{w}^t)$$

↓ step size too high
 'skip' over global minimum

Momentum

$$z^{t+1} = \beta \cdot z^t + \nabla E(w_t)$$

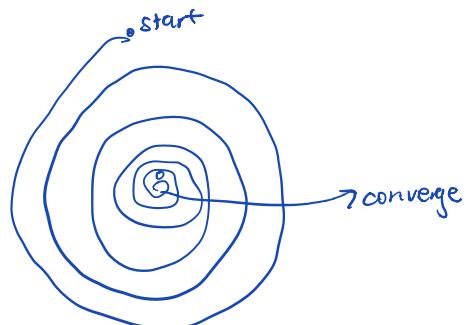
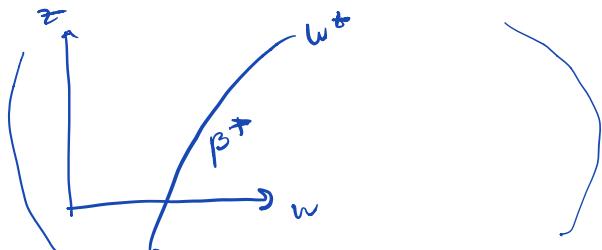
$$w^{t+1} = w^t + \alpha z^{t+1}$$



$$\underline{z^{t+1}} = \underline{\beta z^t} + \nabla E(w_t)$$

velocity damping time step

$$w^{t+1} = w^t + \alpha z^{t+1}$$

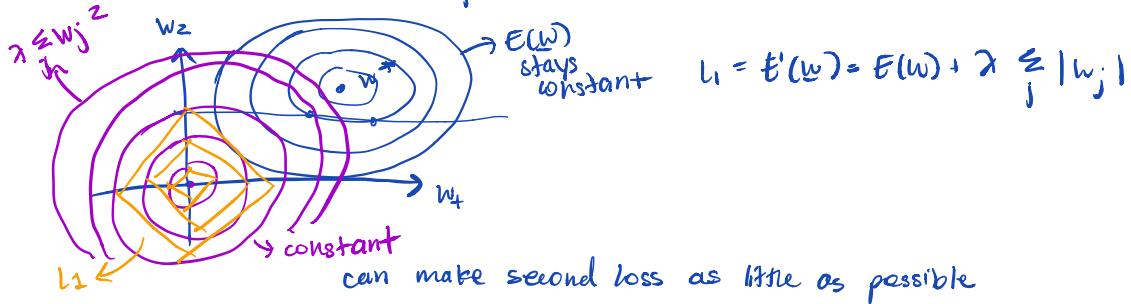


if $\beta > 1$
will swing and go out, rotate and then converge

Regularisation

$$E(w) = \frac{1}{n} \sum_i (y_i - f(x_i; w))^2$$

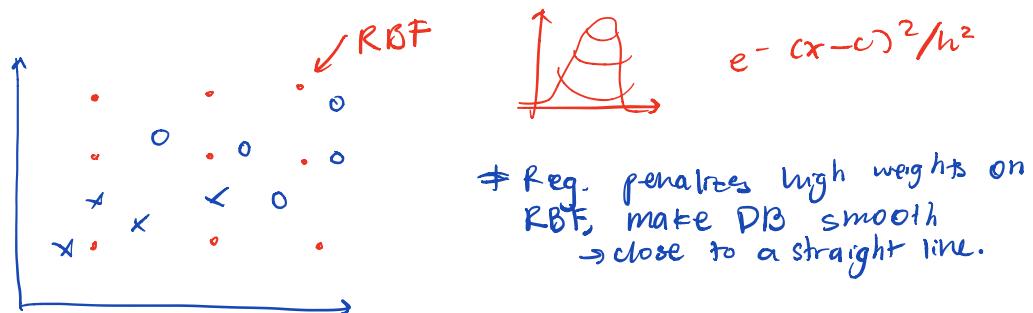
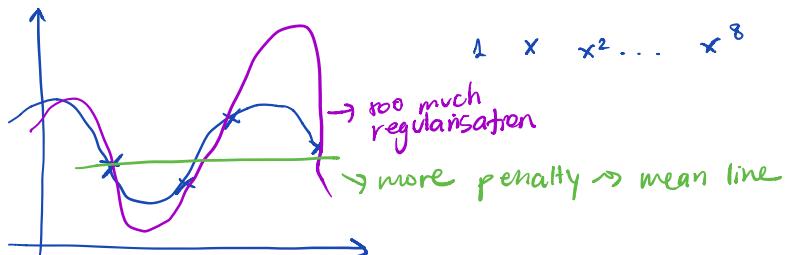
$$E'(w) = E(w) + \lambda \sum_j w_j^2 \quad E' \text{ big when } w_j \text{ big}$$



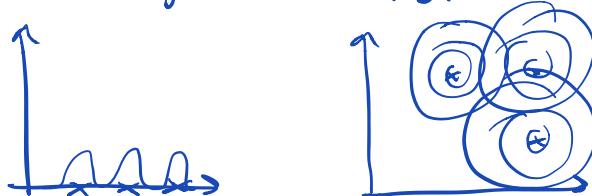
* Along w_1 , the value of the error does not change much, hence L_2 can be large (smaller w_1)

* Along w_2 , $E(w)$ varies more, hence less L_2

Regularisation - add info in order to solve an ill-posed problem or to prevent overfitting.
 ↓
 not enough data to solve.
 - less freedom in parameters



If we don't regularise on RBF → you get 'islands'



October 21, 2019

SVM

Q1. maximising $\frac{1}{2} \|\mathbf{w}\|$ & minimising $\|\mathbf{w}\|^2$.

Why is this minimising the squared of $\|\mathbf{w}\|$ instead of just $\|\mathbf{w}\|$ itself?

this is solving for α_i in $\mathbf{w} = \sum_i \alpha_i \mathbf{x}_i$

solution is quadratic programming problem.

But we know that it is $\frac{1}{2} \|\mathbf{w}\|^2$ minimising $\|\mathbf{w}\|^2$.

Q2. Projection of \mathbf{X} onto \mathbf{w}

$$\mathbf{w}^\top \mathbf{x}$$

assume \mathbf{w} is unit vector

$$\text{w/o assumption} \rightarrow \frac{1}{\|\mathbf{w}\|} \mathbf{w}^\top \mathbf{x}$$

Q3. How can we deal w/ infinite dimensions because of kernel trick?

* point of kernel trick

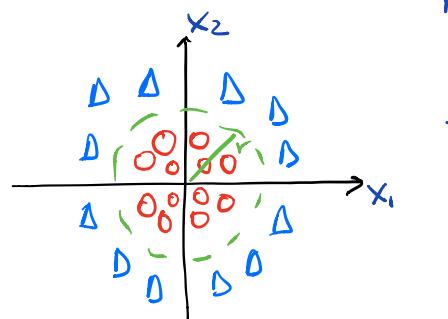
- all we have to compute is $\mathbf{x}_i \cdot \mathbf{x}$ new test point

$$- \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j)$$

- only ever need these dot products

- we can operate in kernel space.

Kernel trick

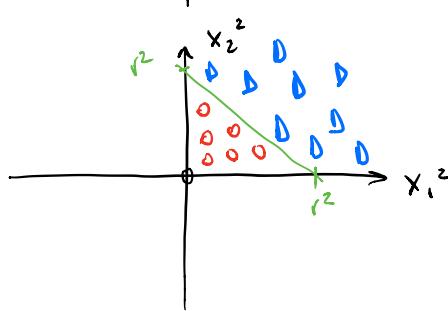


Are these linearly separable? No.

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / \sigma^2)$$

To separate, draw a circle.

$$\mathbf{x}_1^2 + \mathbf{x}_2^2 = r^2$$



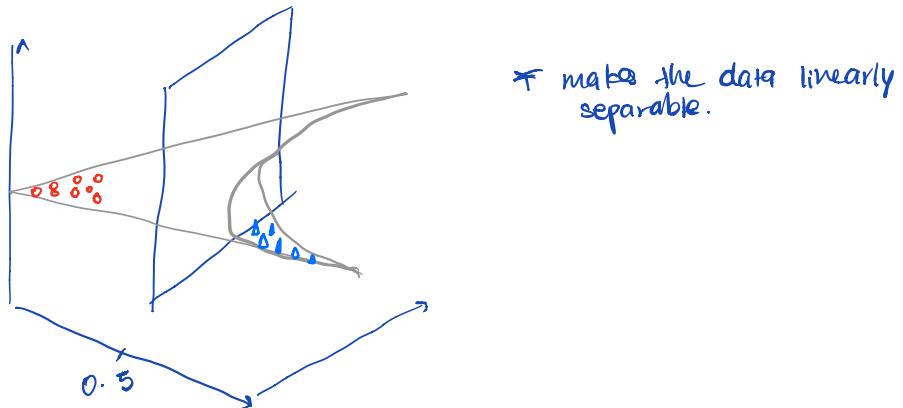
New feature space, $(\mathbf{x}_1^2, \mathbf{x}_2^2)$.

$$\mathbf{x} \rightarrow \begin{pmatrix} \mathbf{x}_1^2 \\ \mathbf{x}_2^2 \end{pmatrix}$$

Example of a kernel for 2D input space

$$\phi(x_i) = \begin{pmatrix} x_{i,1}^2 \\ \sqrt{2}x_{i,1}x_{i,2} \\ x_{i,2}^2 \end{pmatrix} \quad \text{then } k(x_i, x_j) = (x_i^\top x_j)^2$$

How would this transformation look like?



SVM vs. Logistic Regression

SVM optimization

$$\min_w \|w\|^2 + C \sum_i \xi_i$$

$$\text{subject to } y_i f(x_i) \geq 1 - \xi_i, \quad i=1, \dots, n.$$

$$\text{where } \xi_i \geq 0 \text{ for all } i, \text{ and } f(x_i) = w^\top x_i + w_0$$

The optimization problem can be re-written as

$$\min_{w \in \mathbb{R}^d} \|w\|^2 + C \sum_i \max(0, 1 - y_i f(x_i))$$

$$y_i f_i > 1 \Rightarrow 1 - y_i f_i < 0 \quad \text{glunge}$$

$$* \text{glunge} = \max(0, 1 - y_i f_i) \rightarrow \text{glunge} = 0$$

$$y_i f_i = 1 \geq 1 - y_i f_i \Rightarrow 1 - 1 = 0 \Rightarrow \text{glunge} = 0$$

$$y_i f_i < 1 \Rightarrow 1 - y_i f_i > 0 \Rightarrow \text{glunge} = 1 - y_i f_i.$$

Logistic regr. w/ ridge penalty

$$P(Y=+1 | x) = \frac{1}{1+e^{-f(x)}} \quad P(Y=-1 | x) = \frac{1}{1+e^{f(x)}}$$

combining these we have

$$P(y_i|x_i) = \frac{1}{1 + e^{-y_i f(x_i)}}$$

$$\text{Loss} \rightarrow -\log P(y_i|x_i) = \log(1 + e^{-y_i f(x_i)})$$

$$\min_{w \in \mathbb{R}^d} \|w\|^2 + c \sum_i \log(1 + e^{-y_i f(x_i)})$$

\downarrow squared penalty

$$\text{Define } g_\sigma(z) = \log(1 + e^{-z})$$

SVM has a similar op. prob. but with hinge(x)

* Both are convex

* use g_σ is not sparse, but log regr. gives probability output

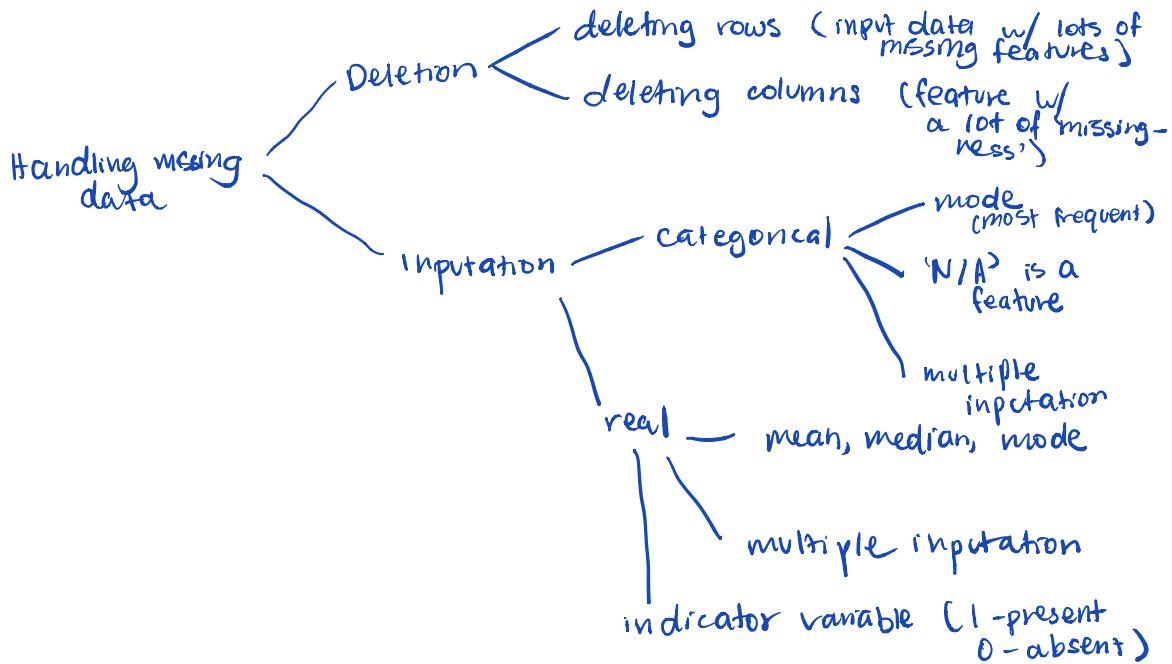
October 24, 2019

Merits of distance measure

- validation (hyperparameter = distance)

Missing values

- understand why the data can be missing



K - Nearest Neighbors

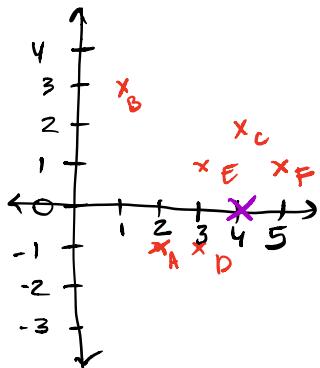
- Voronoi tessellation
- partitions the space that surrounds each dp.
- Odd K \rightarrow tie-breaker
- Increase K \rightarrow smooth out decision surface.

Example Question

① consider the ff. dataset.

$$\begin{array}{lll} A(2, -1, -2) & B(1, 3, -3) & C(4, 2, 0) \\ D(3, -1, -1) & E(3, 1, -1) & F(5, 1, 1) \end{array}$$

Use kNN to predict x_3 from $x_1=4$, $x_2=0$. Assume $k=3$



$$k=3 \rightarrow \text{NN} = D, E, F$$

$$\hat{x}_3 = \frac{-1 + -1 + 1}{3} = -\frac{1}{3}$$

Introduce new point G(10, 1, -1000)

How would the outlier affect the preds of lin. reg and kNN reg.
Be specific about which inputs would be affected.

↑ fit a plane

- With LR, the entire plane would be pushed because of G's value
 - plane tries to go as close to ALL
- With kNN, affects points that are close to G.
 - if G is a neighbor $\rightarrow x_3$ value will be near G's x_3
 - local effect

Suppose H(5, 1, 0) and $k=3$.

- You have 2 NNs!
- add all, divide by two
- noise: diff y values, same x value.

Can you do with integers? Yes. Use mode instead of average.

kernels in KNN

$$p(Y=1 | X) = \frac{\sum_{i=1}^n I(Y_i=1) k(x_i, x)}{\sum_{i=1}^n k(x_i, x)}$$

$\sum_{i=1}^n w_i = 1$
1 > $w_i \geq 0$
weight on
the datapoint
★ point which are nearby has more weight

SVM

$$\hat{y}(x) = \operatorname{sgn} \left[\sum_{i=1}^n y_i k(x, x_i) + w_0 \right]$$

↑ uses of to define support vectors

Parzen

$$p(Y=-1 | X) = \frac{\sum_{i=1}^n I(Y_i=-1) k(x, x_i)}{\sum_{i=1}^n k(x, x_i)}$$

decision boundary

$$p(Y=1 | X) = p(Y=-1 | X) = 0.5$$

$$\begin{aligned} \sum_{i=1}^n I(Y_i=1) k(x, x_i) &= \sum_{i=1}^n I(Y_i=-1) k(x, x_i) \\ &= \sum_{i=1}^n y_i k(x, x_i) = 0. \end{aligned}$$

dB

October 28, 2019

Intrinsic Evaluation

- look at pairs

- should they be in the same group?

$$\text{Rand Index} = \frac{TP + TN}{TP + TN + FP + FN} \quad (\text{accuracy})$$

- overlapping clusters \rightarrow use AMI

→ very rare case

- use table \rightarrow how good is this assignment?

- AMI - information one variable has on the other?

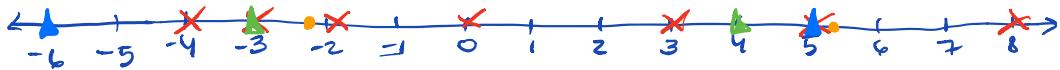
	y_1	y_2	y_3
x_1			x
x_2	x		
x_3		x	

→ perfect information

Example

$\{ -4, -3, -2, 0, 3, 5, 8 \}$.

Run k-means. $K=3$, $\mu_1 = -6$, $\mu_2 = 5$.



$$\textcircled{1} \quad C_1 = \{-4, -3, -2\} \quad \mu_1' = \frac{-4-3-2}{3} = -3$$

$$C_2 = \{0, 3, 5, 8\} \quad \mu_2' = \frac{0+3+5+8}{4} = 4$$

$$\textcircled{2} \quad C_3 = \{-4, -3, -2, 0\} = \mu_1' = \frac{-4-3-2+0}{4} = -2.25$$

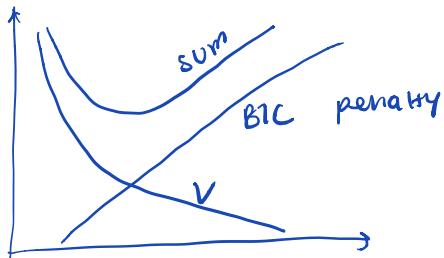
$$C_4 = \{0, 3, 5, 8\} \quad \mu_2' = \frac{3+5+8}{3} = 5.33$$

\textcircled{3} updates do not change!

Choosing the value of K

- cross-validation with hyperparameters
- optimise V for $K=2, 3, \dots$

$$V = \sum_j \sum_{x_i \rightarrow c_j} D(c_j, x_i)^2$$



\Rightarrow BIC - Bayesian Information Criterion

* penalise no. of cluster centres

$$BIC = \max (L - \frac{K}{2} \log n)$$

$$V = -cL$$

$$\max \rightarrow -\frac{V}{c} - \frac{f}{2} \log n$$

$$\min \rightarrow \frac{V}{c} + \frac{(p)}{2} \log n \rightarrow \text{no. of parameter } KD$$

October 31, 2019

Q: Prior Probability for EM

We use $P(b|x)$ to label the data instance. By using Bayes rules in the second line, we need to know prior prob. How can you get prior prob. if we DK which class the instances belong to? Randomly estimate / plot?

$$b_i = P(b|x_i) = \frac{P(x_i|b)P(b)}{P(x_i|b)P(b) + P(x_i|a)P(a)}$$

what is the initial $P(b)$?

We need to have an initial setting (guess)

$$P(a) = 0.5 \quad P(b) = 0.5$$

then they get updated

$$P(b) = (b_1 + b_2 + \dots + b_n) / n$$

$P(a) = 1 - P(b)$ * total law of probs \rightarrow sum up to 1.

$$* P(a|x_i) + P(b|x_i) = 1.$$

We can choose to re-estimate after each iteration.

- We might want to fix variance (σ^2)

A Gaussian MM with K components

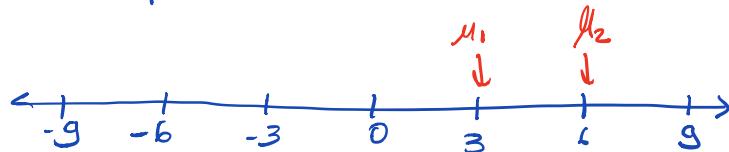
$$P(X) = \sum_{k=1}^K P_k P(X|k)$$

$$P(X|k) = \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp \left\{ -\frac{(X-\mu_k)^2}{2\sigma_k^2} \right\}$$

Consider dataset $\{-9, -6, -3, +3, +6, +9\}$.

Initial means $\mu_1 = 3$, $\mu_2 = 6$, $\sigma_1 = \sigma_2 = 1$. Assume prior $= \frac{1}{2}$

- ① Compute post. probs of the two Gaussians for the obs. $x=6$. Report the results to 2 DP.



$$p(x=6 | k=1) = \frac{1}{\sqrt{2\pi \cdot 1}} \exp \left(-\frac{(6-3)^2}{2 \cdot 1} \right) = \frac{1}{\sqrt{2\pi}} e^{-9/2}$$

$$= 0.0044$$

$$p(x=6 | k=2) = \frac{1}{\sqrt{2\pi \cdot 1}} \exp \left(-\frac{(6-6)^2}{2 \cdot 1} \right) = \frac{1}{\sqrt{2\pi}} = 0.3990$$

$$\begin{aligned} p(k=1 | x=6) &= \frac{p(k=1)p(x=6 | k=1)}{p(k=1)p(x=6 | k=1) + p(k=2)p(x=6 | k=2)} \\ &= \frac{y_2 \times 0.0044}{y_2 \times 0.0044 + y_2 \times 0.3990} = 0.01 \end{aligned}$$

$$p(k=2 | x=6) = 1 - 0.01 = 0.99$$

* Think how fast the Gaussian decays.

- ② Estimate the posterior probability of the 1st Gaussian for every point in ds.

$$p(k=1 | x=-9) = 1.00$$

$$p(k=1 | x=-6) = 1.00$$

$$p(k=1 | x=-3) = 1.00$$

$$p(k=1 | x=3) = 0.99$$

$$p(k=1 | x=6) = 0.01$$

$$p(k=1 | x=9) = 0.00$$

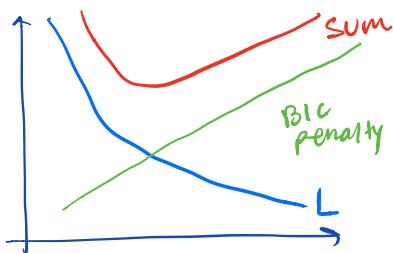
on the left

on the right

* Gaussian decay \rightarrow squared of the distance

How to pick $k \rightarrow$ Occam's razor

BIC: minimize $-L + \frac{k}{2} \log n$.
As k increases, log likelihood increase.



means μ_k \longrightarrow FD parameters

Covariance matrices Σ_k \longrightarrow $K \frac{(D+1)}{2}$

Mixing properties P_k \longrightarrow $K-1$ for each param

K Gaussian in D dimensions

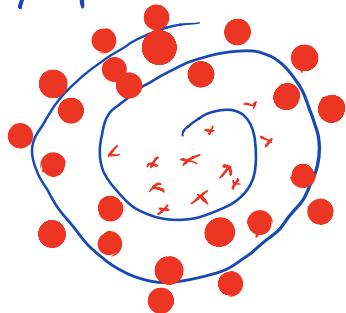
$$\begin{aligned} & \text{Σ symmetric} \\ & D \times D \text{ matrix} \\ & \frac{D^2 - D}{2} + D = \frac{D(D+1)}{2} \end{aligned}$$

Real-world applications

- real-valued vector data
- handwritten digits, normalized 0,1.

November 4, 2019

Q1. Dimensionality reduction can convert non-linearly separable problem into a linearly separable one?



* if we unfold the spiral,
x one side • one side

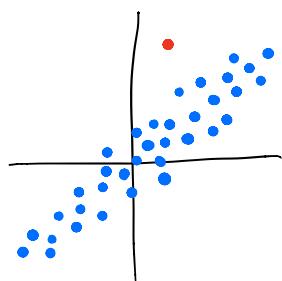
* non-linear dim. red.

Q2: $\Sigma e = \lambda e$ "length"

gives you back same vector, but scaled

if $\|e\|^2 = 1$, then λe has length λ .

Q3. Sensitivity to outliers is not a problem for PCA



covariance matrix will be affected by outliers

* SOME effect on the resulting principal component.

Q4: Lagrange Multiplier Formulation from slides
 Eigen vector = Direction of max variance

$$V = \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^d x_{ij} e_j \right)^2 - \lambda \left(\sum_{j=1}^d e_j^2 - 1 \right)$$

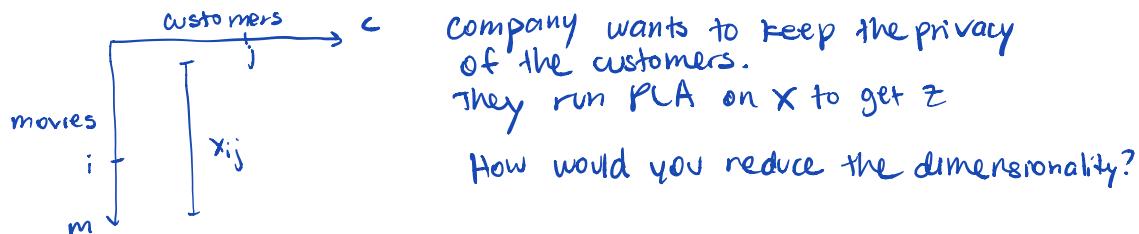
In the Lagrange multiplier part of the eqn λ is the length of the vector.

Why is there no square root?

$|e_j|$ is unit vector, so $\sqrt{1^2} = 1$
 * easier to differentiate \rightarrow 'convenience'

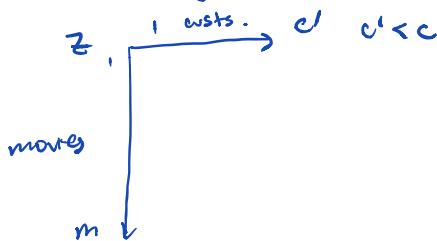
Example Question

- ① MovieFlix - recommend movies to customers



* Reduce dimensionality of the customers.

$Z \rightarrow$ same rows, fewer columns



* Take those diff. vectors, do PCA on them, pick eigenvectors.

* PCA - not taking out columns

- each column is some combination of the customers

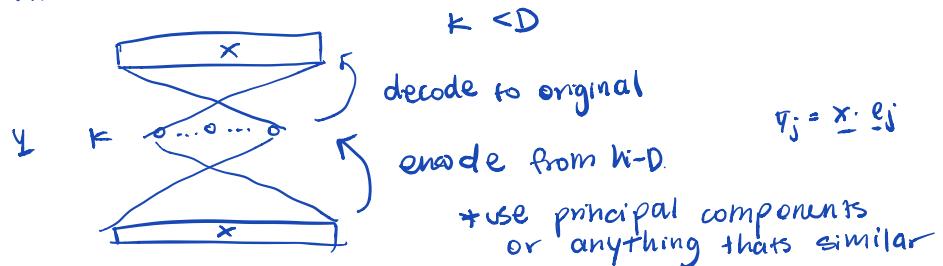
* Movies data is sparse

- PCA on missing data \rightarrow it is handled

- good on data that has a lot of 'missingness'

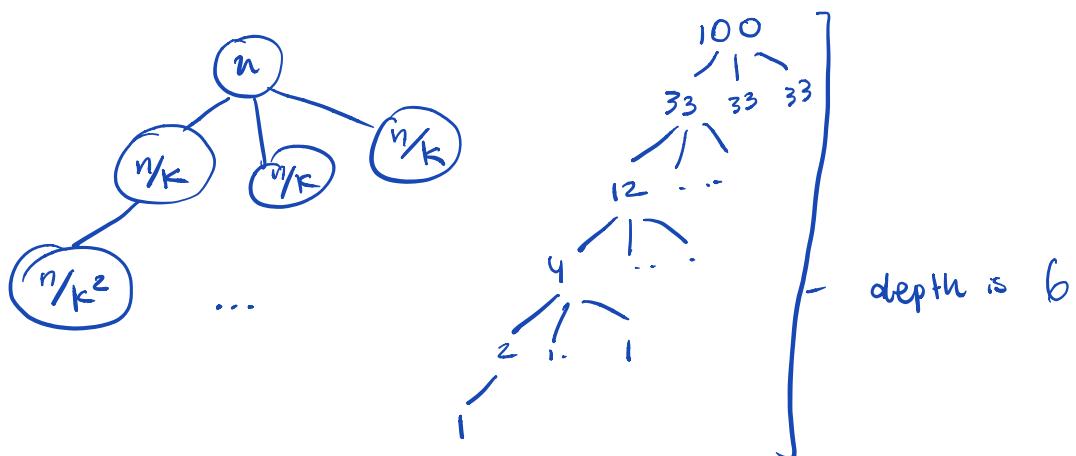
- low-rank decomposition

Auto-encoder



November 11, 2019

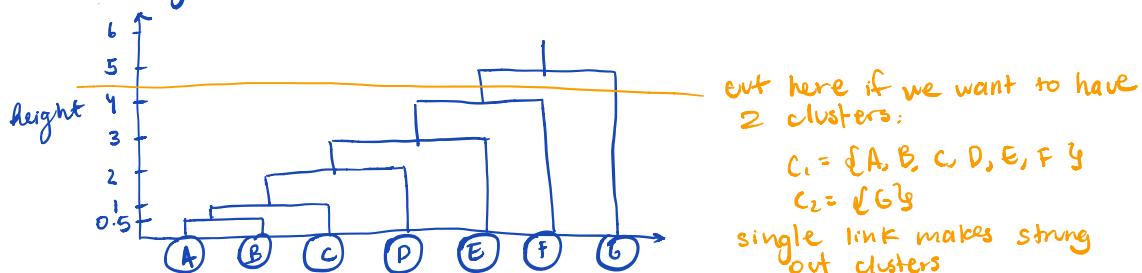
Q1. How many levels are there in a tree created by top-down divisive clustering with a ds of 100 entries, setting $k=3$.
 → clusters may not be of equal size



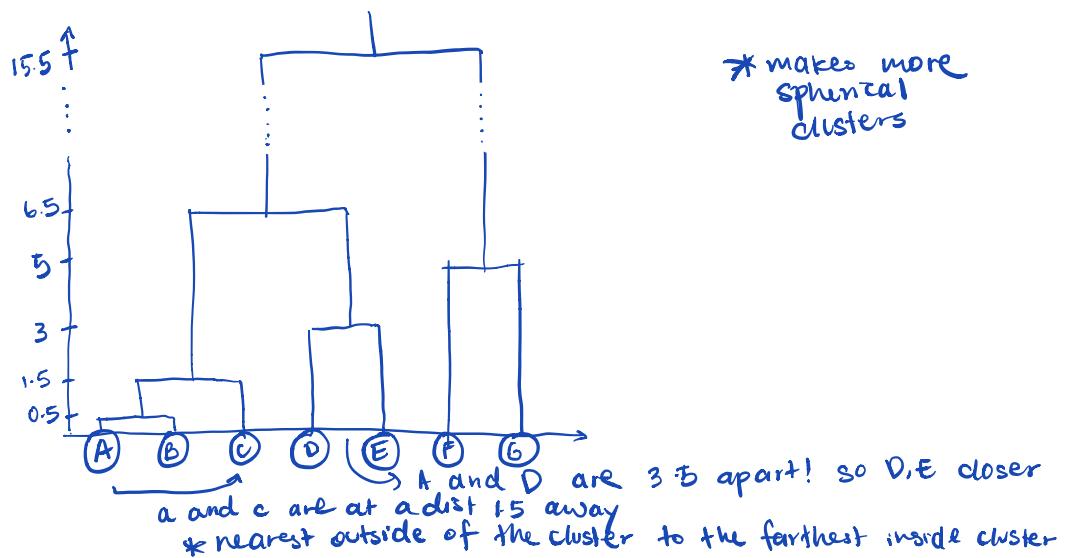
Agglomerative clustering

Dataset: $\{ -3.5, -3, -2, 0, 3, 7, 12 \}$

a. single link

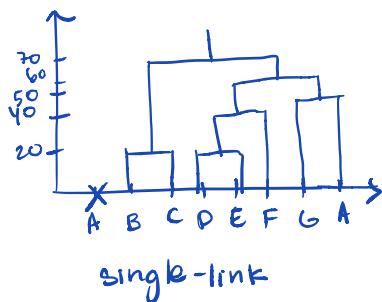
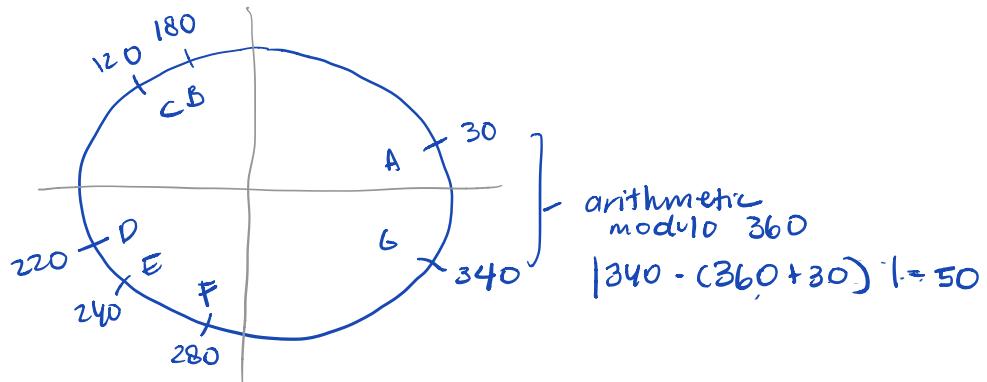


b. complete link



Example Question

Dataset in degrees = {30, 100, 120, 220, 240, 280, 340}



November 18, 2019

Rosenblatt's rule for training perceptrons

$$\hat{y} = \text{sign}(\underline{w}^T \underline{x}) \begin{cases} 1 & \text{if } \underline{w}^T \underline{x} \geq 0 \\ -1 & \text{if } \underline{w}^T \underline{x} < 0 \end{cases}$$

Perception Learning Rule

```

repeat
  for i in 1, 2, ..., n
     $\hat{y}_i \leftarrow \text{sign}(\underline{w}^T \underline{x}_i)$ 
    if  $\hat{y}_i \neq y_i$ 
       $\underline{w} \leftarrow \underline{w} + y_i \underline{x}_i$  ← 'update', dependent on true label
  until all training examples correctly classified
  *if no separating hyperplane, will run forever.

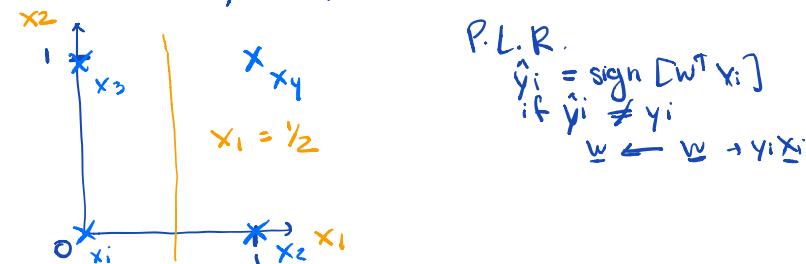
```

Example Question

Data:

	x_0	x_1	x_2	y
x_1^T	1	0	0	-1
x_2^T	1	1	0	1
x_3^T	1	0	1	-1
x_4^T	1	1	1	1

- Initialise weight vector $\underline{w}^T = (0, 0, 0)$
- Run PLR on dataset to convergence. What is final \underline{w} ?
- Consider the dp $(\underline{x}_1, \underline{x}_2, \underline{x}_3, \underline{x}_4)$ in order, and repeat if not ALL correctly classified.



First update

$$\text{sign}(z) = \Theta(z)$$

$$\Theta(w \cdot x_1) = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} = \Theta(0) = 1 \neq y_1 = -1$$

$$w \leftarrow \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} - 1 \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} = (-1, 0, 0)$$

$$\Theta(w \cdot x_2) = \Theta(-1) \neq y_2 = 1$$

$$w \leftarrow (-1) + \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} = (0, 1, 0)$$

$$\Theta(\underline{w} \cdot x_3) = \Theta(0) = 1 \neq y_3 = -1$$

$$w \leftarrow \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} - \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} = \begin{pmatrix} -1 \\ -1 \\ 0 \end{pmatrix}$$

$$\Theta(\underline{w} \cdot x_4) = \Theta(-1) = -1 \neq y_4 = 1$$

$$w \leftarrow \begin{pmatrix} -1 \\ -1 \\ 0 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = (0, 2, 0)$$

NOW run through all of it again (second update)

$$\Theta(\underline{w} \cdot x_1) = \Theta(0) = 1 \neq y_1 = -1$$

$$w \leftarrow (0, 2, 0) - (1, 0, 0) = (-1, 2, 0)$$

$$\Theta(\underline{w} \cdot x_2) = \Theta(1) = 1 = y_2 \text{ correct}$$

$$\Theta(\underline{w} \cdot x_3) = \Theta(-1) = -1 = y_3 \text{ correct}$$

$$\Theta(\underline{w} \cdot x_4) = \Theta(1) = 1 = y_4 \text{ correct}$$

Finally check x_1

$$\Theta(\underline{w} \cdot x_1) = \Theta(-1) = -1 = y_1 \leftarrow -1 \text{ correct.}$$

$$\text{correct } \underline{w} = (-1, 2, 0)$$

$$\text{PB} \rightarrow \underline{w} \cdot \underline{x} = 0$$

$$\begin{bmatrix} -1 \\ 2 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = -1 + 2x_1 = 0$$

$$x_1 = \frac{1}{2}$$

BIAS!