

Pitch set up

Pitch coordinates

The origin of the coordinate system is defined to be at the center spot with the y axis directed along the width of the pitch and the x axis along the length such that the boundaries are at ± 2100 cm and the goalkeeping point is at $00. + 6250$ cm.

[Score on the pitch line to end direction of play follows]

Standardisation

As we are working with a range of data recorded at different football grounds, the data needs to be standardised.

As not every pitch is the same size

We define the standardised pitch dimensions to be 105m in the x and 68m in the y.

These were chosen as they are the only stadium [REDACTED] they play half of their matches at home. So we chose this to be the baseline as it is the most common.

So every x coordinate is transformed by pitch x range(m)

105

and y coordinate by $\frac{\text{Pitch y range (m)}}{68}$

Normalisation

Each game also needs to be normalised \rightarrow [REDACTED] always attacks in the same direction

The opta data defines to the home team attacking left to right (+ve x direction) and the away team right to left (-ve x)

However as [REDACTED] is a constant in all of our data we define that [REDACTED] should always be attacking right to left (-ve x)

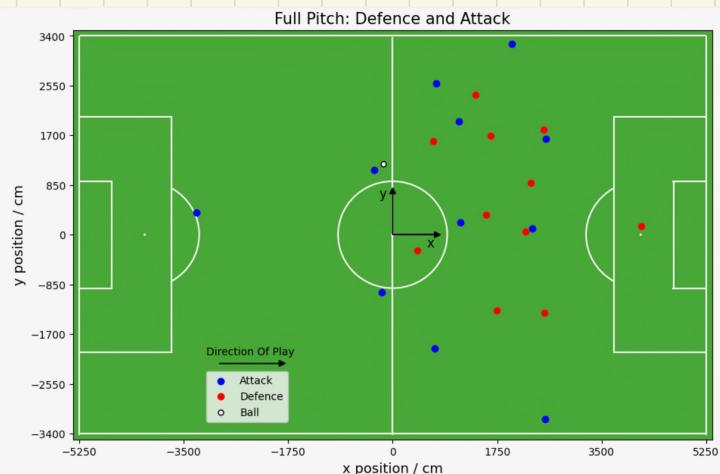
This is achieved by [REDACTED]

If they are they already are attacking R2L so no transformation is needed

If they are not home every player and ball coordinate is rotated around the origin (x and y multiplied by -1)

i.e. [REDACTED] (the home team) now plays right to left

A full representation of all 8 of these definitions



Data Structure

Inputs + output

Inputs : The model takes coordinates of the 11 attacking players as input + the coordinates of the ball

It still needs to be decided whether the coordinates of the ball adds any benefit to the model, if it does not conform it

because it should be decided whether the model benefits from the position of the attacking players

Outputs : The model returns the coordinates of the 11 defending players

Data Shape

Each time step is defined as a frame, the opta data tracks at 5 for so each time step is 0.2 seconds

The aim of the project is to accurately predict how  would strike effectively against a given attacking formation

It needs to be noted that the strike of the player does not just depend on the current attacking frame but also preceding frames

Therefore a simple neural network cannot be used, one has to find a relationship between the current defensive formation and the current attacking formation but also preceding attacking formations

→ thus, one needs time series data.

Therefore the frames need to be built into sequences [defined as a series of frames]

the length of a sequence was determined by modelling a game as a Markov chain (see transition matrix below)

This proved that the position of the ball is 2 frames completely independent if the frames are more than 10 frames apart

So the sequence length was defined at 9 seconds to try to predict a sequence that is long enough but also only contains frames that are not statistically independent
→ 45 frames

It is also worth noting that the data was also batched to train faster, i.e. train on multiple sequences simultaneously

So the input data had shapes (batchsize, sequence length, number of coordinates)

Batch size was later fixed as a hyperparameter, but sequence length = 9

and number of coords = 25 for players and 8 ball coords

Output shape was the same but had 22 coordinates instead

Sparse Matrix Data Shape

Another option is to structure the data as a matrix. The matrix has 1D score dimensions as the pitch \rightarrow grid and into a 10×10 grid.

The positions of players are then converted by a 1D in the corresponding grid cell, organization else is a score.

This (obviously) will be more efficient than coordinates but one home pitch function needs (detected in one function global) and pitch index (only take 1D inputs).
This is up to you for next semester.

Data filtering

The first filter on the data is on  not having the ball. Note you have  will be referred to as the 'defending team'.

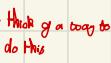
This means that the defending team is always facing an attack and we can investigate how they line up defensively.

The next filter is cutting off when the defence is in a 'defensive' position.

If the defending team have players pressing into the attacking team's half, this is not a defensive position, it's a pressing position.

So one recursive rule that has defending player in the attacking half.

This Selection Cut is somewhat arbitrary, as it is difficult to generalize a defensive position.

But the cut does not make any sense due to how much data   think of a way to do this

note this doesn't filter for the central formation not the wings

Another cut is also made for 4 defenders, 4 midfielders and 2 attackers according to be in the defending team.

One defender a 4-4-2 or a 'box standard' defensive line up.

Any 5 back formation would give too little data \rightarrow but could be 'more defensive'.

and any 3-5-2 variation is too attacking.

So any game / substitution that does not fit these criteria are omitted.

And finally, simple cut such as, there must be 8+ players within the pitch boundaries overall as the ball, and the ball must be in play (not dead) for a frame to be 'alive'.

So to summarize, the list of criteria that defines an 'alive' frame:

- 8+ players + ball in pitch boundaries
- Ball is alive
- A 4-4-2 line up
- 3-5-2 defences \rightarrow the one half
- Attacking team has the ball

The Continuous play is then split into episodes which are constant phases of play where each of these criteria are met.

As soon as one of these criteria isn't valid the episode will end.

Any episode shorter than 45 frames is deleted (maybe introduce padding later) the episodes longer than 45 frames are then quantized into sequences.

Note that this is done from the end of the episode backwards.

This helps to remove any 'messy' start such as a corner, free kick etc. as the discrete will be removed if the episode is not an integer multiple of 45.

Annotations:
Panis do not accept variable sequence length
If we want to use shorter sequences, need to pad inputs, however this caused problems.
initially padded with zeros than first frame
better be backfill than forward fill

Ordering output / Expected

The opta data does not order the players in the same order for each frame.

Initially, the output and the predicted data was being compared element wise:

$$\text{Pred} : [x_1, y_1, x_2, y_2, x_3, y_3, \dots, x_n, y_n]$$
$$\text{Exp} : [x_1, y_1, x_2, y_2, x_3, y_3, \dots, x_n, y_n]$$

However, the expected array is not always in the same order. Sometimes, x_i may be the coordinate of the 3rd, sometimes the 5th.

So the model will struggle to learn what to put in the first index if the expected is always changing and this is what it is being compared to.

So we have 2 options:

- make sure the expected is always in the same order [Gk, Df, Df, Df, Mid, ...]
- make the loss calculation non-index dependent

We first tried the first option by sorting: [Gk, Df, Df, Df, Mid, ...]

But there is no further play data to sort by position, and due to variations in the order of each group, we then sorted each group by their number.

However, there is no guarantee this order is constant between games \rightarrow this caused a spike in loss throughout the season.

So, making the cost independent of the order of the expected array was the next option.

i.e., the model doesn't have to predict index 2 to match index 2 in the expected, it just needs to predict the coordinate at some index in the predicted array.

This was done by defining a cost matrix for each frame, this describes the loss of each pairing in the pred and exp arrays:

$$\text{Cost-Matrix} = \begin{pmatrix} L_{11} & L_{12} & L_{13} & \dots \\ L_{21} & L_{22} & & \\ L_{23} & & \ddots & \\ \vdots & & \vdots & \end{pmatrix} \quad \text{where } L_{ij} = (r_i - s_j)^2$$

and r_i is an expected coordinate and s_j is a predicted coordinate

The Hungarian algorithm is then used to find the minimum of each row without repeating combinations.

Along the same vein, the model should be invariant to changes in the order of the inputs.

The idea is to sort the attracting coords (x only) in increasing order, then have the ball x and y at the end.

Transition matrix

Weeks 4-6

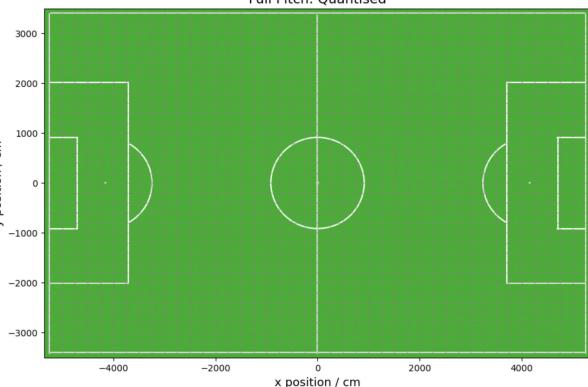
Here we focus on how long each sequence of play should be.

The aim is to find how long it takes for 2 events to be independent in a game of football

For this we modelled a game of football as a discrete-homogeneous Markov process.

The pitch was quantised into $0.25m^2$ boxes \rightarrow note each pitch has different dimensions. So this was normalised to be (136×210) \rightarrow most common size

Full Pitch: Quantised



the cell of the ball at each time was recorded along with

which cell the ball is in at the next time.

Defining transition matrix

This data was then used to define the transition matrix for the Markov process.

The transition matrix is defined as:

$$P_{i+1} = Q P_i \quad \text{where } P_i \text{ and } P_{i+1} \text{ are the probabilities of the location of}$$

the ball at i and $i+1$ respectively

$$Q = \begin{pmatrix} q_{11} & q_{12} & q_{13} & q_{14} & \dots & 0 \\ q_{21} & q_{22} & 0 & 0 & 0 & 0 \\ q_{31} & 0 & & & & \\ q_{41} & 0 & & & & \\ \vdots & \vdots & & & & \end{pmatrix}$$

In this case Q is $([68 \times 105] \times [68 \times 105])$ dimensional.

where q_{ij} represents the probability of the ball moving from cell $j \rightarrow i$ in one time step

So after all positions of the ball were found, each index of the Q matrix was filled in

i.e. if the ball was in cell $(1,1)$ at $t=0$ then $(1,2)$ at $t=1$ then $Q_{(1,2), (1,1)}$ was given a count.

To finally define Q it needed to be turned into probabilities as opposed to counts

As each column should sum to 1 (all possible movements from $i \rightarrow$ cell $j = 1$)

So each column was normalised by its sum

Stationary state

The stationary state of a markov process is defined as:

$$Q \underline{P}^{\text{st}} = \underline{P}^{\text{st}}$$

i.e. any further application of Q will have no effect on \underline{P} . The system will evolve no more

This corresponds to the eigenvector of Q that has an eigenvalue of 1.

Due to the large size of Q here, the eigenvectors could not be solved for analytically.

So power iteration was used to find $\underline{P}^{\text{st}}$.

The eigenvector with the largest eigenvalue can always be found by iteratively calculating

$$\underline{V}_{t+1} = \frac{Q \underline{V}_t}{\|Q \underline{V}_t\|}$$

Until \underline{V}_t converges. Note \underline{V}_0 is a random vector.

↳ below a chosen tolerance (we chose the modulus to be $< 10^{-6}$)

Note that for a given t , $\underline{P}(t) = \sum_i C_i(\lambda_i)^t \underline{V}_i$

so as $t \rightarrow \infty$, $\underline{P} \rightarrow C_n \underline{V}_n$ where n has eigenvalue 1.

Also note that $\|P\| = 1$ (probability) so C_n and thus $\underline{P}^{\text{st}}$ can be found by normalising \underline{V}_n .

Once $\underline{P}^{\text{st}}$ was found, the time taken for an initial condition to reach the steady state was found.

This was done by iteratively calculating $\underline{P}_{\text{tri}} = Q \underline{P}_t$ until $\underline{P}_{\text{tri}} = \underline{P}^{\text{st}}$ with a 5% tolerance on the magnitude.

This proves how long it takes two events to be independent as the stationary state is independent of the initial state.

So a sequence longer than this should not be used as it will contain frames that are statistically independent.

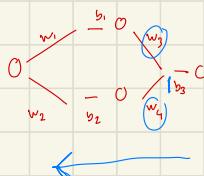
The time taken to reach the stationary state was found to be ~ 10 .

$$\underline{P}_{t+1} = Q \underline{P}_t$$

$$\underline{P}_{t+2} = Q^2 \underline{P}_t$$

$$\underline{P}_{\text{st}} = \boxed{Q^{\infty} \underline{P}_t}$$

$$\underline{P}_{\text{st}} = \boxed{Q^{\infty+1} \underline{P}_t}$$



Cost Function

Weeks 2-3 & 7

Initial definition

For the network to learn, there needs to be a quality of output measure \rightarrow between true data and the prediction

Initially, the idea was to simply calculate the mean squared error of each frame: $\text{Cost} = \frac{1}{N} \sum_{i=1}^N (r_i^{\text{exp}} - r_i^{\text{pred}})^2 \Rightarrow N=11$ reduce to 10 and ignore 6th?

This would then be averaged over the 45 frames of the sequence \rightarrow the sequence has one value for the cost (not one for each frame)

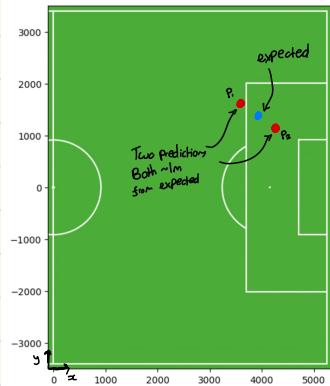
And if the data is batched, the average of each sequence would be calculated for the cost of the batch.

Angular Cost

It quickly became apparent that the model would never converge to a zero cost \rightarrow it would always be finite (especially for the testing data)

But the same modulus between the predicted and expected position can correspond to very different positions.

i.e. Consider a modulus of $\sim 1m$ between the predicted and expected:



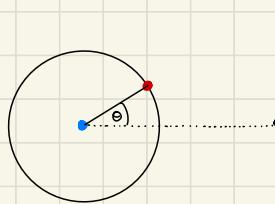
Despite both predictions having the same cost, P2 is obviously a better

Prediction from a footballing opinion.

So we needed to define a Cost function that would return a

Sensible grading position if the exact position is not predicted

If the expected position has the same y coordinate as the centre of the goal, this is straightforward



A circle is drawn to show the location points would have the same distance based cost.

Our new cost function leads to minimise the angle between the line connecting the centre of the goal and the expected and the line between the predicted and the expected.

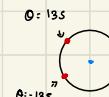
i.e. $\theta=0$ corresponds to the expected, predicted and goal being in line.

it is also worth noting that we defined the maximum cost to be 180

Anything larger than this would be negative. This speeds up the minimisation of the cost \Rightarrow

as it can be reduced by moving anticlockwise in the negative segment.

i.e. $0 < \theta < 360$, and $0-359^\circ$, to reduce it the position would have to move through 359° .

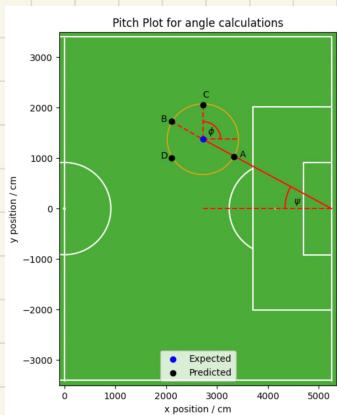
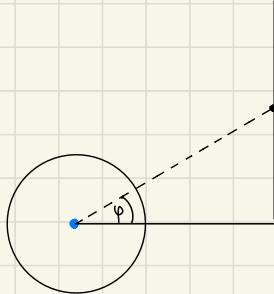
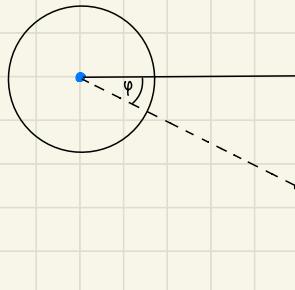


However, with $0 < |θ| < 180$, if $θ = -r$, it only has to be moved 1° to be zero. \Rightarrow much more efficient

Up back up with learning time graph?

Or simple enough to just claim?

When the expected isn't inline with the center of the goal, it is slightly more complicated.



The angle between the expected position and the center of the goal is defined as $θ$. This is $0 \rightarrow 90$ ° for above the $y=0$ axis

and $0 \rightarrow -90$ ° below $y=0$

The angle from the horizontal is the colatitude as before but the difference between $θ$ and $φ$ is minimized \Rightarrow zero angular cost leaves the prediction inline with the center of the goal.

Exponential loss

We have played with the idea of introducing an exponential decay to the importance of the distance \rightarrow i.e., when the prediction is far from the expected, the angular cost is negligible, only changing the distance will reduce the overall cost.

But when the distance $\sim 1m$, the angle becomes equally important as the distance to reduce the cost.

Mathematically, without an exponential, the total cost is:

$$\text{Cost: Distance Cost} + \text{Angular Cost}$$

As the distance cost is measured in cm^2 , the two are comparable when the avg. distance $\approx 200 \text{ cm}^2$

So the average distance needs to be $\sim 1\text{m}$ \rightarrow this is way too small

Initially we added a simple multiplying factor to the angular cost

$$\text{Cost: Distance Cost} + (50 \times \text{Angular Cost})$$

This means the angle cost becomes important at an average distance of $\sim 1\text{m}$

This is much better, however, we would like the distance cost to not be linear in hope to train the model faster

So the new cost will have the form:

Note that the loss for a sequence is calculated as the average of each frame
And a batch is the average of each sequence

Sparse matrix loss (issues)

Week 2

Decide the loss we will implement: a sparse matrix format.

Initially the idea was to simply extract the coordinates of each position from the matrix, then do all of the previous steps.

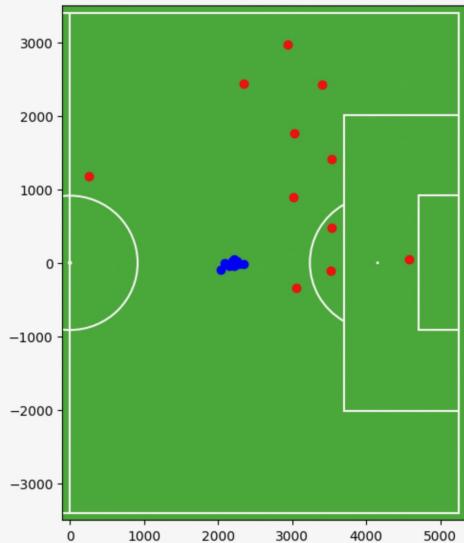
However, extracting the coordinates is a non-differentiable operation, there is no associated gradient function so backpropagation can't be done. \rightarrow model doesn't learn.

One method to fix this involves simply finding the difference between the matrices and summing the components.

However, this will be very step-like, the model will only be rewarded when predicting the exact position, it will not be guided in the right direction.
More thought is needed on this next semester.

For a long time (weeks 5-6) had problems with the model learning.

Test getting results like this:



Model wanted to predict

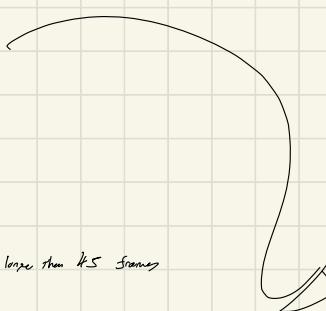
all players together

I think this was a issue with padding sequences

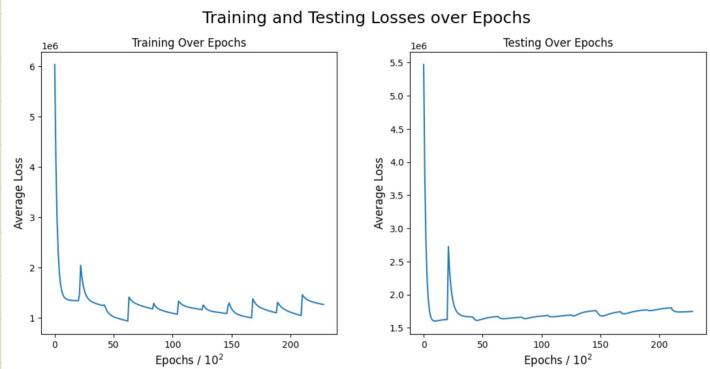
With the RNN \rightarrow RNN tended to predict

average positions.

If padding, \rightarrow predicting avg. position



So we removed padding (for now) and took sequences longer than 45 frames



loss never got below $\sim 10^6$

Network architecture

Week 8-9

Note for this section, the hyperparameters include:

- Batch size
- Epochs
- Layers
- Hidden dimensions
- Learning rate

most of these are tuned later for fixing overfitting
Apart from 1D and layers.

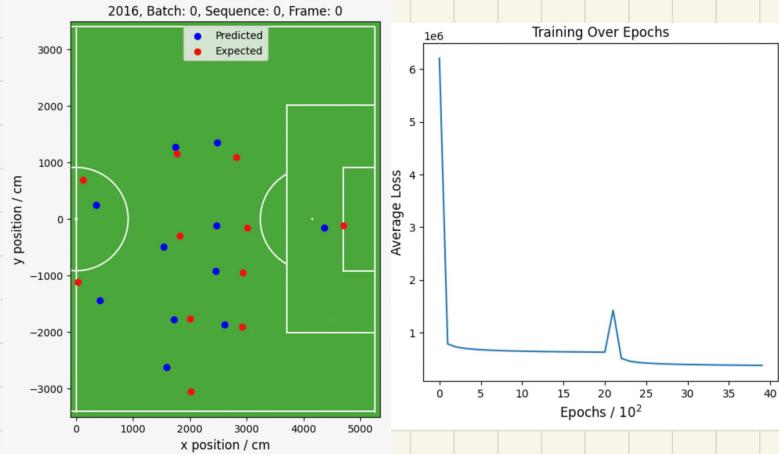
RNN

The first aim is to get a model that can accurately learn the training data \rightarrow not worried about testing

First we used a basic 1 layer RNN, however this struggled to capture the

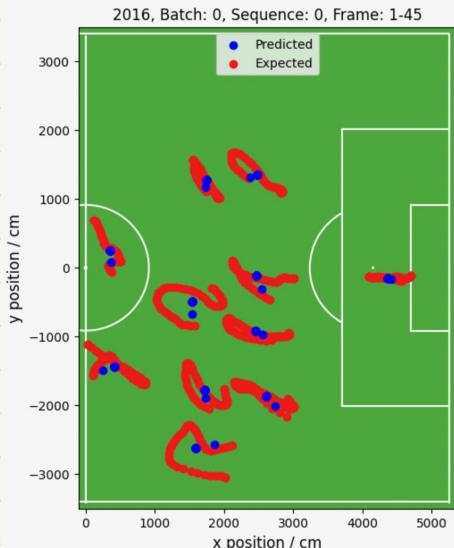
time series data, it tended to predict the average position of each player across the

Sequence.



Viewing the whole sequence shows that the prediction doesn't really move across the sequence, but is still being predicted in the right area.

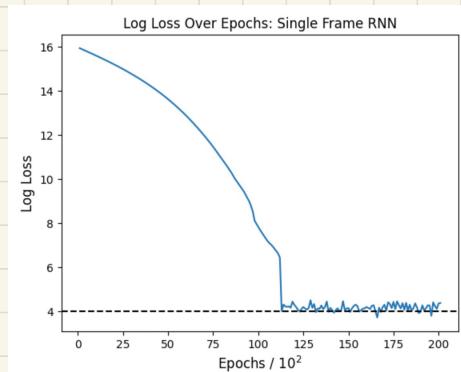
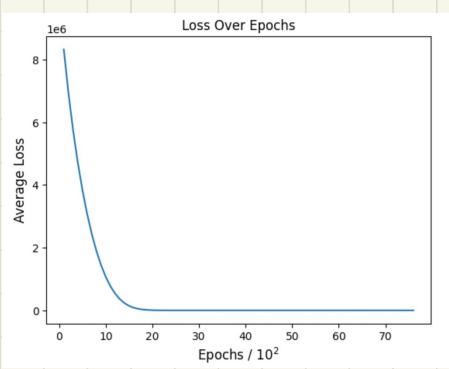
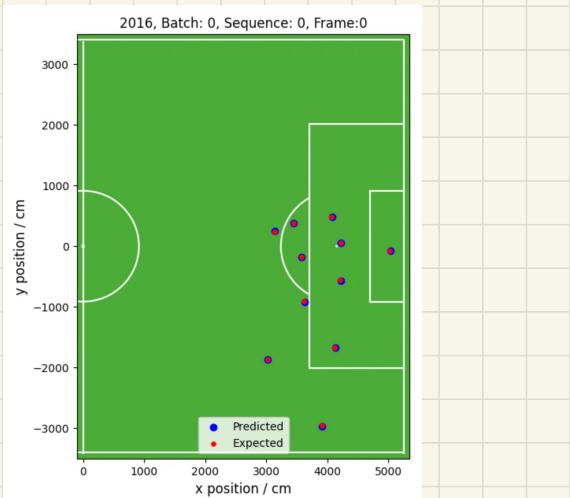
↳ each prediction is on the right expected path



It's clear to see from the loss graph that the model has the ability to learn, however, it asymptotes at around 50,000. This is an average distance per player of ~ 2 m. This is not bad, but it's clear that the time structure of the data isn't being captured.

Just to solidify this argument, we trained the model on a stationary sequence (every frame is the same)

If our previous hypothesis is right, this should be almost perfect.



It can be seen from the pitch plot that this gives a perfect

solution, this corresponds to a log loss of ~ 4 , ie

loss ~ 50 cm per player.

We define loss of order ~ 100 as optimum.

We will use this as a basis to build a model that can predict sequences, games and seasons.

LSTM

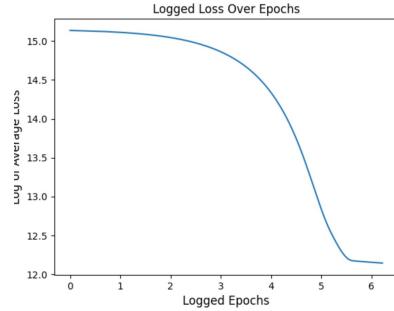
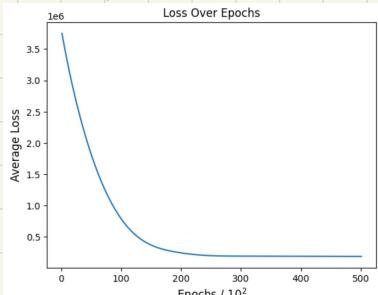
week 10 ~ 12

We now move towards an LSTM, these are more capable at capturing long term time dependencies.

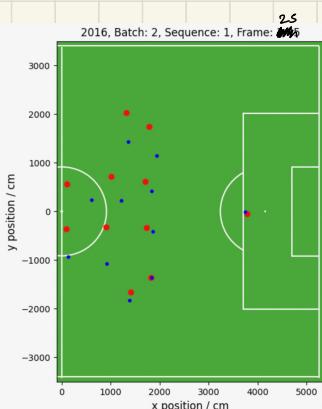
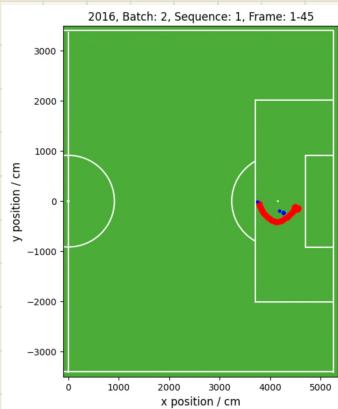
We initially start with just a sequence, we've accurately predicted a frame, the sequence is the next step.

We start with a 1 layer LSTM with 64 hidden dimensions.

For 1 Sequence:

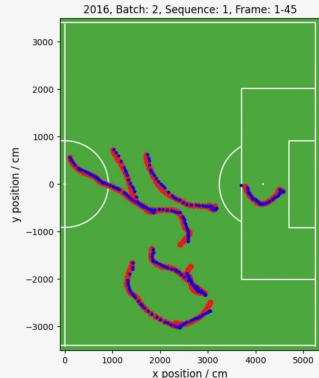
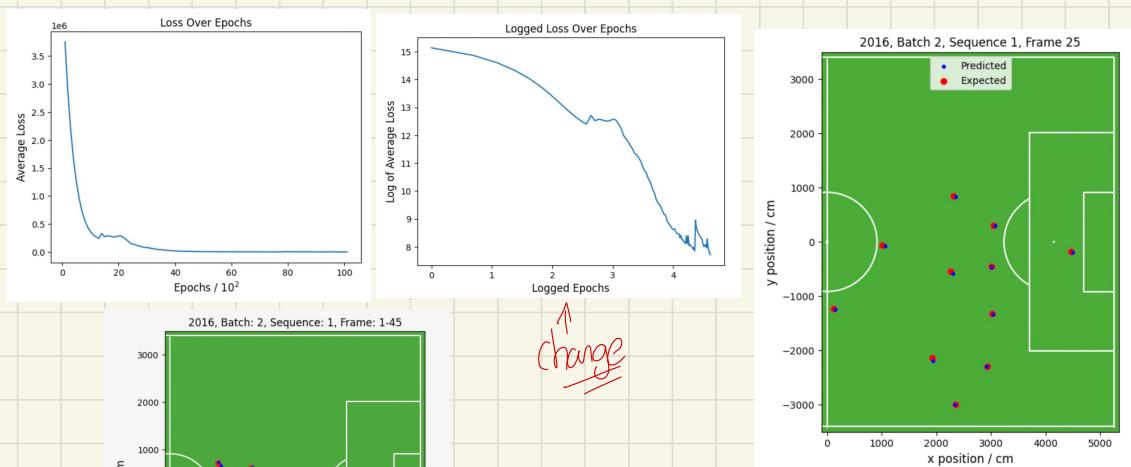


This loss asymptotes at $\sim 150,000$ (ie quite bad) but looking at the pitch:



it is starting to be able to predict some time evolution.

By simply increasing to 2 layers:



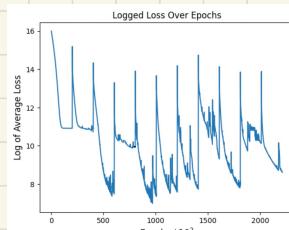
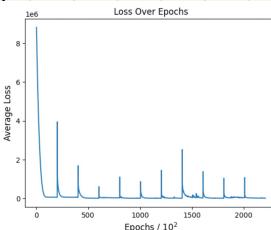
This is accurately predicting a game but also a whole sequence

← Plotting 6 player paths to demonstrate. Cleverly paths

loss asymptotes at ~ 2000 $\sim e^8$

However this takes $\sim 10^4$ epochs

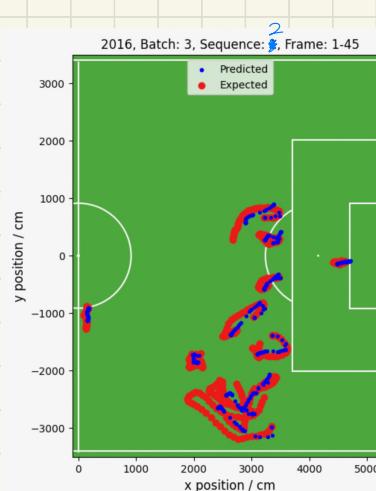
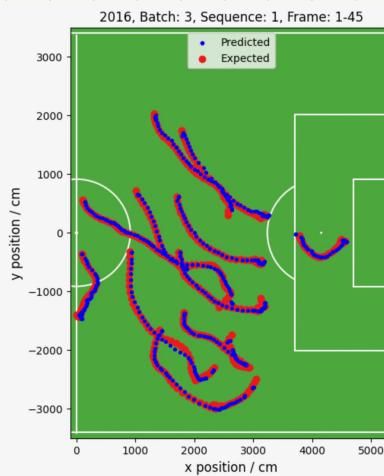
And then for one game:



It's clear to see for some sequences, it learns very well \rightarrow reaching a loss of ~ 2000

However, there are a couple that it cannot learn well, asymptoting at a loss of $\sim e^{10} \sim 20,000$

Investigating some of these batches:



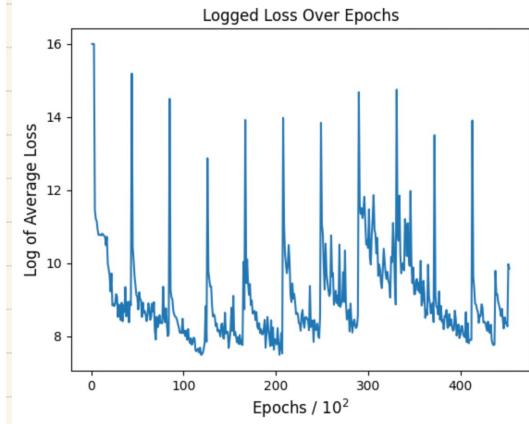
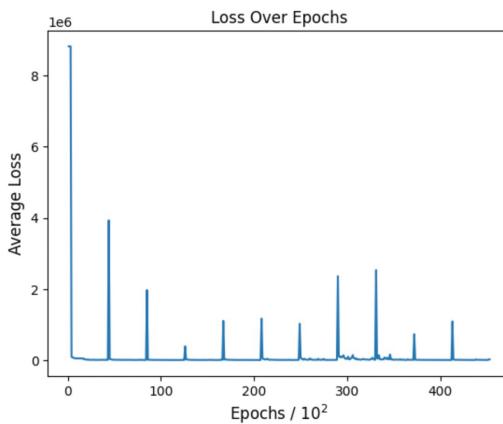
Batch 2 is roughly the same as before, but batch 3 which asymptotes at $\sim 20,000$ struggles to predict all of the player paths over the whole game.

So we look into expanding the LSTM to improve the predictions over a game before looking into multiple games.

After some trial and error with different activation functions and other layers

we settle on a 2 layer LSTM with a dropout layer, 2 fully connected linear layers and a lastly ReLU activation function

For the same game:



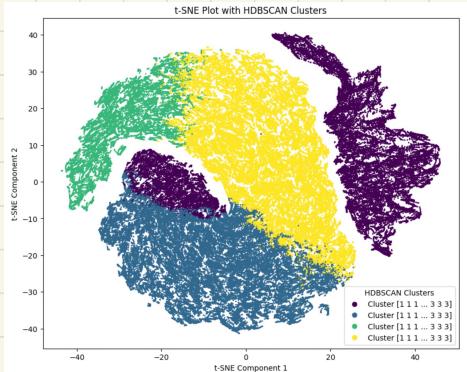
It's clear to see here, the batches that were previously asymptoting at $\sim 20,000$ are now $\sim 2000 \rightarrow$ huge improvement

The highest a batch asymptotes is $\sim e^9 \sim 8,000$.

This is not perfect ($\sim 3m$ between each player)

In an attempt to clean the data (remove frames that are difficult to predict)

we tried clustering:



This showed some promise at sorting corners etc...

But needs more work done.

Might be easier to just manually clean data

All files are saved on databricks

$$Q = \begin{pmatrix} q_{11} & q_{12} & q_{13} \\ q_{21} & q_{22} & q_{23} \\ q_{31} & q_{32} & q_{33} \end{pmatrix} \quad \text{Probability of } 2 \rightarrow 1 \text{ in 1 time step}$$

$$Q^2 = \begin{pmatrix} \cdot & \boxed{\cdot} & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{pmatrix} \quad \text{Probability of } 2 \rightarrow 1 \text{ in 2 time steps} \rightarrow P(X_2 = 1)$$

$$= \begin{pmatrix} q_{11} & q_{12} & q_{13} \\ q_{21} & q_{22} & q_{23} \\ q_{31} & q_{32} & q_{33} \end{pmatrix} \begin{pmatrix} q_{11} & q_{12} & q_{13} \\ q_{21} & q_{22} & q_{23} \\ q_{31} & q_{32} & q_{33} \end{pmatrix}$$

$$= \begin{pmatrix} q_{11}^2 + q_{12}q_{21} + q_{13}q_{31} & q_{11}^2 + q_{12}q_{21} + q_{13}q_{31} & q_{11}^2 + q_{12}q_{21} + q_{13}q_{31} \\ q_{21}q_{11} + q_{22}q_{21} + q_{23}q_{31} & q_{21}q_{11} + q_{22}q_{21} + q_{23}q_{31} & q_{21}q_{11} + q_{22}q_{21} + q_{23}q_{31} \\ q_{31}q_{11} + q_{32}q_{21} + q_{33}q_{31} & q_{31}q_{11} + q_{32}q_{21} + q_{33}q_{31} & q_{31}q_{11} + q_{32}q_{21} + q_{33}q_{31} \end{pmatrix}$$

$$= \begin{pmatrix} q_{11} & q_{12} & q_{13} \\ q_{21} & q_{22} & q_{23} \\ q_{31} & q_{32} & q_{33} \end{pmatrix} \begin{pmatrix} q_{11}^2 + q_{12}q_{21} + q_{13}q_{31} & q_{11}^2 + q_{12}q_{21} + q_{13}q_{31} & q_{11}^2 + q_{12}q_{21} + q_{13}q_{31} \\ q_{21}q_{11} + q_{22}q_{21} + q_{23}q_{31} & q_{21}q_{11} + q_{22}q_{21} + q_{23}q_{31} & q_{21}q_{11} + q_{22}q_{21} + q_{23}q_{31} \\ q_{31}q_{11} + q_{32}q_{21} + q_{33}q_{31} & q_{31}q_{11} + q_{32}q_{21} + q_{33}q_{31} & q_{31}q_{11} + q_{32}q_{21} + q_{33}q_{31} \end{pmatrix}$$

$$Q^2 = \begin{pmatrix} (q_{11}^2 + q_{12}q_{21} + q_{13}q_{31})q_{11} + (q_{21}q_{11} + q_{22}q_{21} + q_{23}q_{31})q_{12} + (q_{31}q_{11} + q_{32}q_{21} + q_{33}q_{31})q_{13} & \dots \\ \dots & \dots \end{pmatrix}$$

$$= \begin{pmatrix} (q_{11}^2 + q_{12}q_{21} + q_{13}q_{31})q_{21} + (q_{21}q_{11} + q_{22}q_{21} + q_{23}q_{31})q_{22} + (q_{31}q_{11} + q_{32}q_{21} + q_{33}q_{31})q_{23} & \dots \\ \dots & \dots \end{pmatrix}$$

$$\text{So independent if } P(X_1 = 2 \text{ and } X_2 = 1) = P(X_1 = 2) \times P(X_2 = 1)$$



Probability of being
in state 2 at time 1

Probability of being
in state 1 at time 2

Probability of being

in State 2 at time 2

and State 1 at time 2

Rough notes
From here

$$P_{111} = Q P_1$$

$$P_{112} = Q P_{11}$$

$$= Q^2 P_1$$

$$\vdots$$

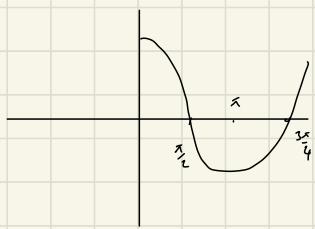
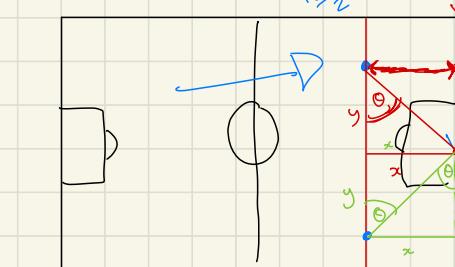
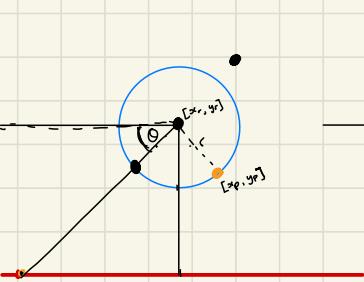
$$P_n = Q^n P_0$$

$$\# P'' = Q P''$$

$$P(X_3 = 1 \mid X_0 = 1)$$

$$\begin{pmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \\ 1 & 3 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 4 & & \\ & 4 & \\ & & 4 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 & 1 \\ 2 & 3 & 1 \\ 1 & 3 & 1 \end{pmatrix}$$

$$1 = L_{21}(4) + 1(0) + 0(0)$$



$$\Omega = \theta - \theta_0$$

Diagram showing a circle with center (x_r, y_r) and radius r . A point (x_p, y_p) is marked on the circle. A dashed line connects the center to the point. The angle θ is labeled.

$$\tan^{-1}\left(\frac{y}{x}\right)$$

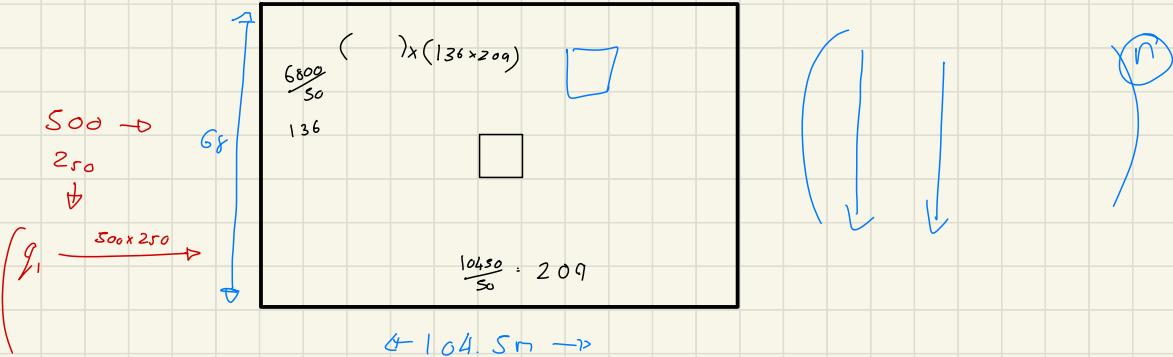
$$GK$$

$$34$$

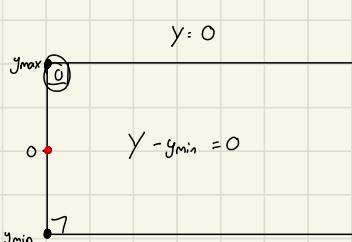
$$105$$

$$a^2 = b^2 + c^2 - 2bc \cos \theta$$

$$\cos^{-1}\left(\frac{b^2 + c^2 - a^2}{2bc}\right) = \theta$$



4-104.5m \rightarrow

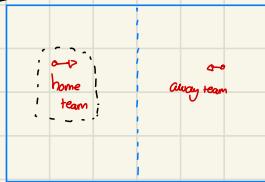


$$\underline{Q} = \begin{pmatrix} q_{11} & q_{21} \\ q_{12} & q_{22} \end{pmatrix} \quad \begin{pmatrix} q_{11} & q_{21} \\ q_{12} & q_{22} \end{pmatrix}$$

$$\underline{Q}^2 = \begin{pmatrix} (q_{11}q_{11} + q_{21}q_{12}) & q_{11}q_{21} + q_{21}q_{22} \\ q_{12}q_{11} + q_{22}q_{12} & q_{12}q_{21} + q_{22}q_{22} \end{pmatrix} \quad Y_{\max} - (y_{\min}) = \frac{6000}{50}$$

2169

3644



We want defense to always be playing R2L

So if TO I is away \rightarrow no change

If TO I is home \rightarrow reverse

Now A/I doesn't matter

we always want City to be L2R

So if home \rightarrow nothing,

if away \rightarrow reverse

$$\text{data} = \begin{bmatrix} x \\ y \\ z \\ a \end{bmatrix}$$

$$\underline{Q} = \begin{pmatrix} \frac{2}{3} & \frac{1}{2} \\ \frac{1}{3} & \frac{1}{2} \end{pmatrix} \Rightarrow \lambda_1 = \frac{1}{6}, \underline{v}_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$

$$\lambda_2 = 1, \underline{v}_2 = \begin{pmatrix} 3 \\ 2 \end{pmatrix}$$

↳ always has to have $\lambda_i = 1$

↳ otherwise as $t \rightarrow \infty$ $\underline{P}(t) = \underline{0}$

Suppose $\underline{P}(0) = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ (ie in state 1)

need to write this as a superposition of \underline{v}_i :

$$\underline{P}(0) = C_1 \begin{pmatrix} 1 \\ -1 \end{pmatrix} + C_2 \begin{pmatrix} 3 \\ 2 \end{pmatrix}$$

Q could have 2 or more $\lambda_i = 1$ but this is a reducible Markov chain.

$$C_1 = \frac{2}{5}, C_2 = \frac{1}{5}$$

↳ Steady state Solution

So $\underline{P}(t) = \frac{2}{5} \begin{pmatrix} 1 \\ -1 \end{pmatrix} \left(\frac{1}{6}\right)^t + \frac{1}{5} \begin{pmatrix} 3 \\ 2 \end{pmatrix} (1)^t$

If $\underline{P}(0) = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \lambda_1 = \frac{1}{6}, \underline{v}_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$

as $t \rightarrow \infty$: $\underline{P}(t) = \begin{pmatrix} \frac{2}{5} \\ \frac{3}{5} \end{pmatrix} + \begin{pmatrix} \frac{3}{5} \\ \frac{2}{5} \end{pmatrix}$

$$\lambda_2 = 1, \underline{v}_2 = \begin{pmatrix} 3 \\ 2 \end{pmatrix}$$

$$C_1 = -3C_2$$



$$C_1 + 3C_2 = 0$$

$$-C_1 + 2C_2 = 1$$

$$5C_2 = 1$$

$$C_2 = \frac{1}{5} \text{ and } C_1 = -\frac{3}{5}$$

So $\underline{P}(t) = -\frac{3}{5} \begin{pmatrix} 1 \\ -1 \end{pmatrix} \left(\frac{1}{6}\right)^t + \frac{1}{5} \begin{pmatrix} 3 \\ 2 \end{pmatrix} \quad \text{↳ Steady state}$

• Steady state Solution is independent of initial Conditions.

• Just need to: - find \underline{Q}

• Find $\underline{P}^{\text{st}}$

• Calculate at what t $\underline{P}(t) = \underline{P}^{\text{st}}$ (Say with 5% certainty)

• Remember that \underline{Q} has to be defined with probabilities over a constant time step

↳ ie, Probability of transition after 1 frame.

↳ then steps in t have to be the same as the time step.

Transition matrix report notes

• Largest eigenvalue has to be 1 (Can't be larger than any entry)

• Can use power iteration methods to find corresponding eigenvector

• Or Krylov subspace methods (Arnoldi method)

Finding Q1:

• Divide pitch into 0.25 m^2 areas:



• Find all instances of being in Square 1

• Find position of the ball at the next time step for every instance.

• Calculate all probabilities

• Repeat for over 28424 squares



68
104.5

28424×28424 matrix

transformation: $(+5250, -3400)$

$(-5250, 3400)$

$(0, 3400)$

$(0, 0)$

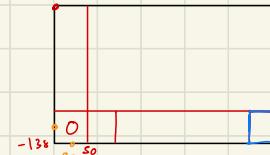
no. grid. y: 136

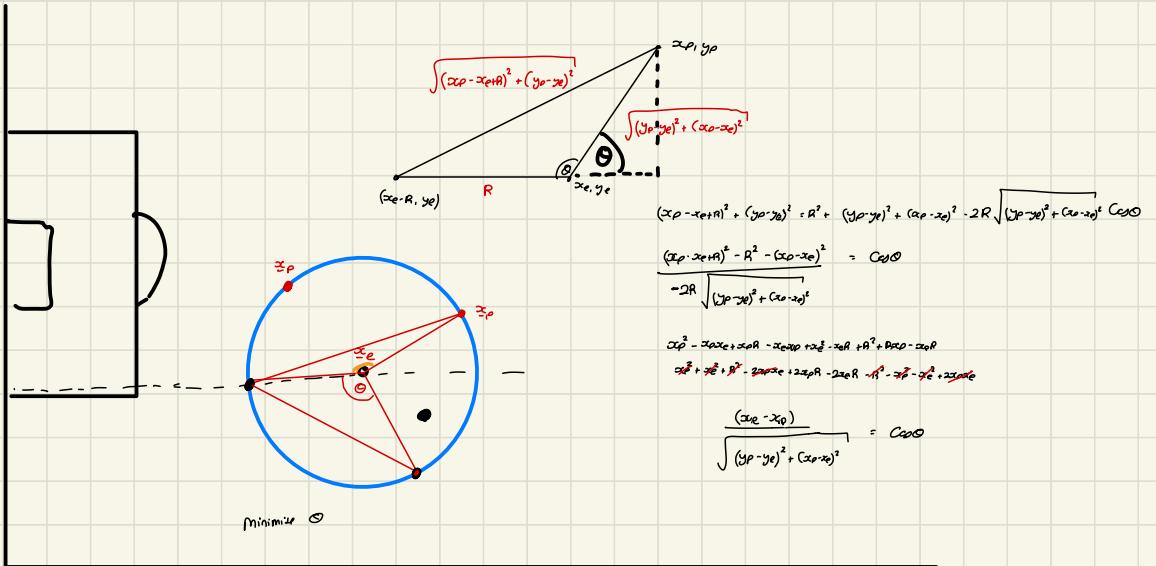


$f(-5250, -3400)$

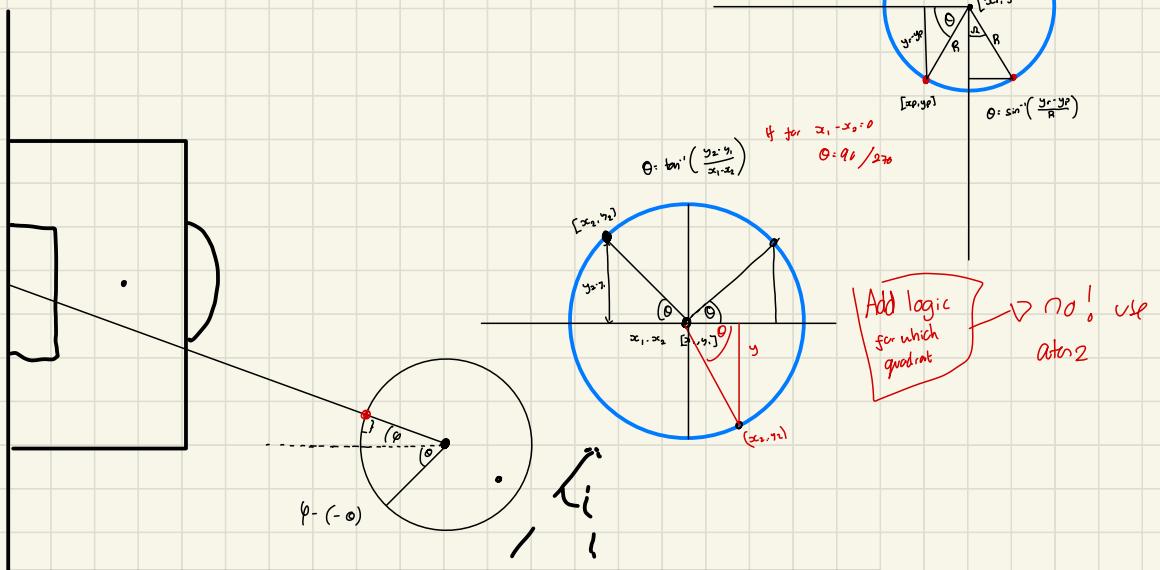
transformation: $(+5250, +3400)$

Coord
grid.size

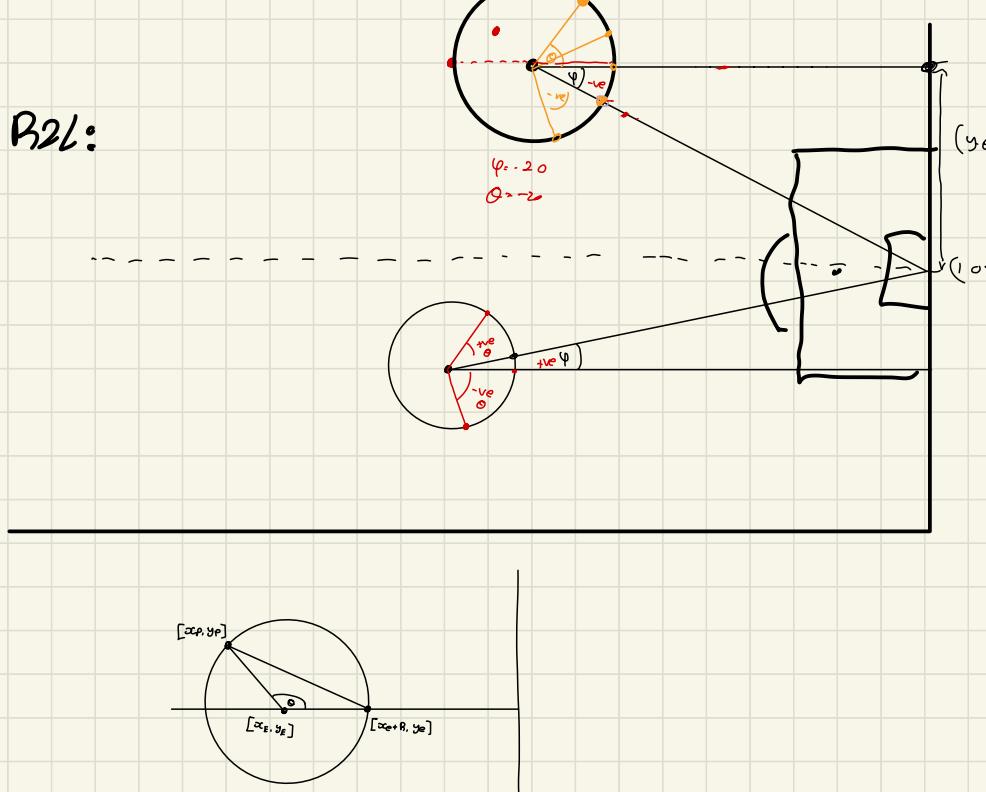




if off center:



For R2L:



$$(x_e + R - x_p)^2 + (y_e - y_p)^2 = R^2 + R^2 - 2R^2 \cos \theta$$

$$(x_e + R - x_p)^2 + (y_e - y_p)^2 = 2R^2(1 - \cos \theta)$$

$$= R^2 + (x_e - x_p)^2 + (y_e - y_p)^2 - 2R \sqrt{(x_e - x_p)^2 + (y_e - y_p)^2} \cos \theta$$

$$x_e^2 + R^2 + x_p^2 + 2x_e R - 2x_e R - 2x_p R = R^2 + x_e^2 + x_p^2 - 2x_e x_p - 2R \sqrt{(x_e - x_p)^2 + (y_e - y_p)^2} \cos \theta$$

$$\frac{(x_p - x_e)}{\sqrt{(x_e - x_p)^2 + (y_e - y_p)^2}} = \cos \theta$$

X

atan 2

Note Sequences need to be same length \rightarrow padding

\hookrightarrow LSTM, GAN \rightarrow Also consider activation functions.

- Consider how the initial hidden state is defined

(45, 32, 2105x68)

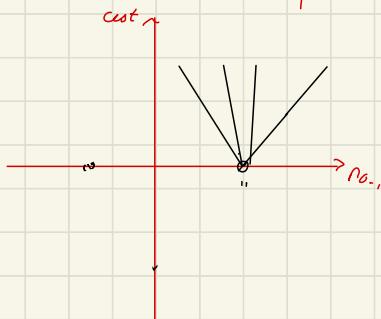
[45, 32, (105, 68)]

- 128 hidden dimensions

- 1 layer

Sparse matrix Cost

$$\text{Cost} = \left| 1000 n_{\text{p.}} - 11000 \right|$$



Cost for when # players != 11:

$$C = 100,000 \exp(0.0083 \text{ mod})$$

So when $\text{mod} = 12, 10$ the cost is $\approx 100,000$

(still leaves range for next Cost function to behave exponentially)

↳ next Cost function will always be below the minimum of this cost

Cost maximum is less than $100,000$ so always have knowledge of increasing

• Cost for when # of players = 11

$$\exp(1376 \text{ P}) + 100 = 110,000 \quad \exp(0.0083 \text{ mod}) + \text{angle}$$

$$P = 0.0084$$

maximum mod: when all plays on opposite side of pitch (mod = 1376)

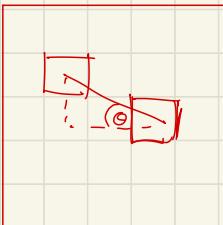
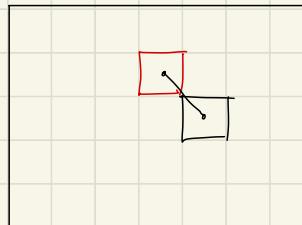
want this to be just below 110,000

$$\text{mod} + \frac{1}{1 + e^{-\text{mod}}}$$

$$[100,000]$$

$$\approx 200$$

$$\text{Cost} = 10000(\text{mod} - 11) + (e^{0.0083(\text{mod} - 11)} - 1)$$



$$e^{0.0000083 \text{ mod}} = 270$$

$$0.0000083 \text{ mod} \approx \ln(270)$$

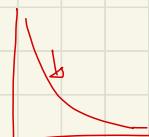


$$e^{\sqrt{68^2 + 105^2}}$$

$$11 \times \sqrt{68^2 + 105^2}$$

$$e^{p \times 1376} = 100,000$$

$$[0.0083]$$



model \rightarrow 12

Cost-A \rightarrow 1,000

model \rightarrow 11

Cost-B \rightarrow 999

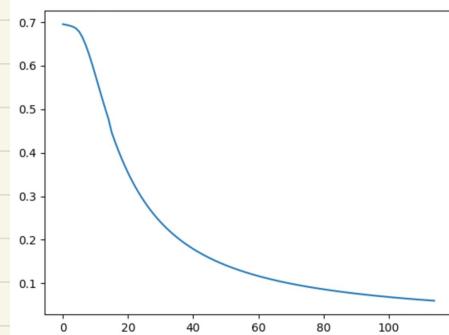
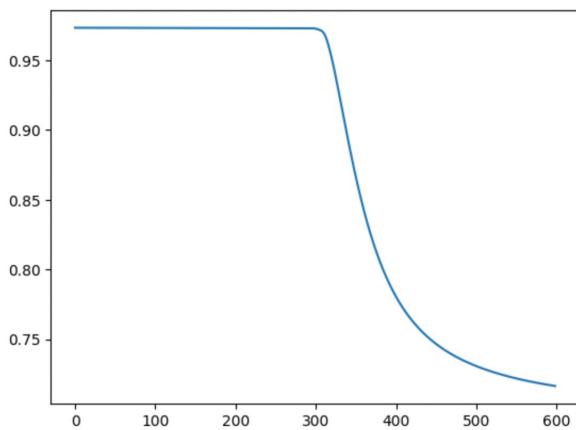
Sparse matrix
Format

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

E — T

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 2 \\ 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 0.1 \\ 0.2 \\ 0.8 \\ 0.2 \\ 0.3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0.1 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$



BCE Cost (sigmoid before) to try and force model to give 11 one's per frame.

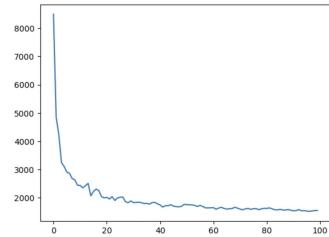
BCE, no sigmoid before.

Realised Can't use BCE because of occasional 2's in Sparse matrix!

Initially Cost not decreasing, tuned learning rate until Cost decreased as expected

- Tried taking the output of the model, normalising between 0-1 using sigmoid then rounding to nearest int (turning network output to binary)
- Then put this into a Cost function to force the model to predict 11 ones:

$$\text{Cost} = | \text{Pred} - \text{Expected} |$$



\Rightarrow Cost function over epochs.

However, realised this will try and predict the correct position (not what we want at this point)

- Current goal is to force the model to predict only 11 players per frame. (ie, all zeros apart from 11 ones, or occasional 2,3)
- So instead, sum each frame after the transform, and minimise $\text{Cost} = \text{sum} - 11$, should reach zero when successfully predicting 11 players.

However, at the moment, Cost is not reducing even after tuning learning rate. Possibly a problem with Sigmoid. Next, try a uniform transform.

