

```

import tensorflow as tf
hello_const = tf.constant("hello")
with tf.Session() as sess:
    output = sess.run(hello_const)
    print(output)

```

- TF data is not stored as integers, floats, or strings. The values are encapsulated in an object called **Tensors**
- hello_const \rightsquigarrow 0-dimensional string tensor.
 \rightsquigarrow a constant tensor
- $A = \text{tf. constant}(1234)$ \rightsquigarrow 0-dim int32 tensor
 $B = \text{tf. constant}([123, 456, 789])$ \rightsquigarrow 1D int32 tensor
 $C = \text{tf. constant}([[[123, 456, 789], [2, 3, 4]]])$ \rightsquigarrow 2D \times \times

② ~~Variables~~

- **Computational Graph**
- **tf.Session()** an environment for running computational graphs

• **tf.placeholder()**, **feed_dict**

```

x = tf.placeholder(tf.string)
with tf.Session() as sess:
    output = sess.run(x, feed_dict={x: "hello"})

```

• **tf.add()**, **tf.subtract()**, **tf.multiply()**

• **tf.cast()** \rightsquigarrow no implicit casting!



- Classification \rightarrow core building block task in many ML tasks
- \downarrow
- reinforcement learning
- detection
- ranking
- regression

• **Logistic Classifier (Linear Classifier)**

$$Wx + b = \hat{y}$$

$$\text{Softmax function: } S(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

y_i, y_j : scores

Softmax function takes the scores (a.k.a. logits) and turns them to probabilities.

Tensorflow uses $xW + b$ instead of $Wx + b$

~~• tf.Variable(), n = tf.Variable(5)~~

~~• tf.global_variables_initializer()~~

• **tf.global_variables_initializer()**

• **tf.truncated_normal()** \rightsquigarrow generate random numbers from Normal dist
(n-features, n-tables)

• **tf.zeros()**

• **tf.matmul()**

• **tf.nn.softmax()**

Cross Entropy

- Using one-hot encoding of input we can measure how well we are doing by comparing the two vectors (onehot label, and the output of softmax (probabilities))

$s(y)$ ← probabilities

0.7
0.2
0.1

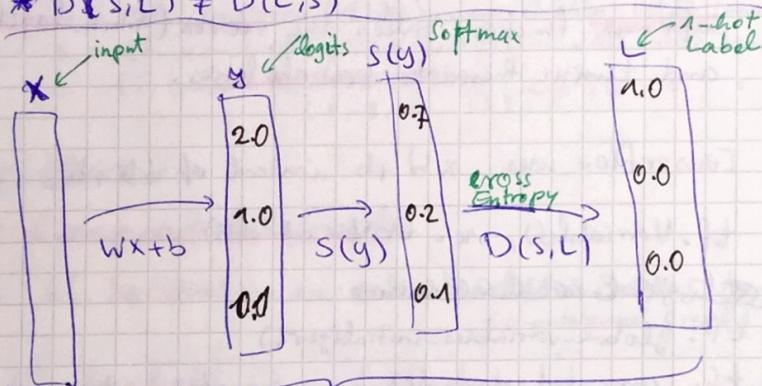
$$D(s, L) = - \sum_i L_i \log(s_i)$$

distance

L ← label
(one-hot)

1.0
0.0
0.0

* $D(s, L) \neq D(L, s)$



Multinomial Logistic Classification

$$D(S(WX+b), L)$$

- classification goal ($D\downarrow$ for correct classification and $D\uparrow$ for wrong classification)



• Loss = Average cross-entropy distance

$$L = \frac{1}{N} \sum_i D(S(WX_i + b), L_i)$$

Two Questions

- How do we feed image pixels to the classifier
- Where do we initialize the optimization.

Numerical Stability

- be worried when working with very large and very small numbers
- specially when adding them together.

$$a = 1000000000$$

for i in range(1000000):

$$a = a + 1e-6$$

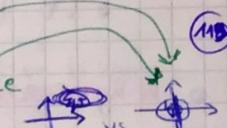
print(a - 1000000000) # prints 0.953674316406 instead of 1.0

- We want the values involved in the calculation of this big loss function to never get too big or too small

- Guiding principle:

$$\mu(x_i) = 0 \text{ zero mean}$$

$$\sigma(x_i) = \sigma(x_j) \text{ equal variance}$$



- Badly conditioned problem vs. well conditioned problem

PINK - 128

128 and ~~not make~~ from every channel

- For images subtract 128 from every channel
 - doesn't change the content of the image
 - But makes it much more suitable for numerical calculation. (optimization).

Weight Initialization

- draw weights randomly from $N(0, \sigma)$
- σ large \Rightarrow output of the first initial point of optimization high order of magnitude
 \Rightarrow prob. dist. by Softmax too peaky (opinionated)
- σ small \Rightarrow output small
 \Rightarrow prob. dist. of Softmax uncertain

\Rightarrow choose a small: σ

Training, Validation, Test set

- measuring performance

- Once you know how to measure performance on a problem, you have already solved half of it

Overfitting

- when the classifier memorizes the training set
 - it fails to generalize.

Why separate validation and test set?

- Because every time you tune model parameters as a result of validation set, you actually implicitly leaking information about your validation set to the classifier!



Rule of 30:

a change that affects 30 examples in your validation set, is usually statistically significant and can be trusted.

e.g. ~~dataset~~

Validation set size 3000

- improvement from 80% to 80.1%
are noise, because they are affecting 3 or 15 validation examples
- improvement from 80% to 81% is significant because it affects 30 validation examples.
 \Rightarrow that's why most people hold back a validation set of > 30,000 examples
 \Rightarrow changes $> 0.1\%$ in accuracy are significant and can be trusted

- If don't have enough data \rightarrow do cross-validation
- But cross-validation can be a slow process.
 \Rightarrow get more data is better.

Training Logistic regression using gradient descent is great

- we directly optimizing the error measure (great idea)
- That's why in practice a lot of ML research is about designing the right loss function to optimize
- But the biggest problem is scaling

④ Problem of scaling GD

- rule of thumb
!
- if computing the loss $\sum_i D_i$ takes n flops,
then computing the gradient of L takes about $3 \times$ more flops.
 - The L is huge, it depends on every single data point in data set
 - But we want to be able to train on lots of data.
 - also GD is iterative \Rightarrow we have to go through the entire DS many times 10s or 100s of times.

\Rightarrow Let's cheat:

- instead of computing the Loss we compute an estimate of it! (a terrible estimate in fact :))

⑤ Stochastic GD (SGD)

- compute the loss for a very small random subset of data. Between 1 to 100 samples each time
- it is very important that it is random!
- iterate many many times more.



- Vastly more effective than simple GD
 - scales well with both data and model

⑥ Momentum

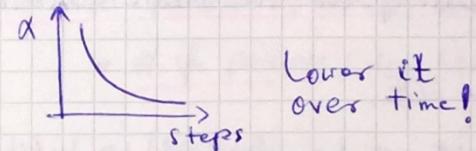
- Helping SGD :
- | | |
|------------------------|-------------------|
| mean = 0 | } inputs |
| equal variance | |
| initial | } initial weights |
| equal variance (small) | |

Momentum:

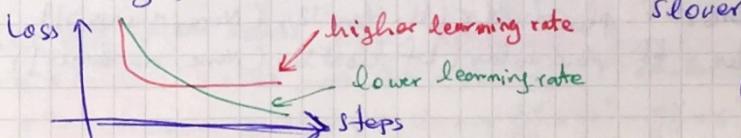
- take advantage of the accumulated knowledge from previous step about direction to go
 - e.g. keep a running avg. of the gradient
- $$M \leftarrow 0.9 M + \Delta L$$

⑦ Learning rate decay

- SGD: take smaller noisier steps towards goal
- How small? It is a whole area of research itself
- One thing however is always beneficial: make that step smaller and smaller as you train!
- e.g. on exponential decay, or reduce it every time loss reaches a plateau



- Learning faster is not better than learning slower



- SGD ≈ Black Magic (with so many hyperparam to tune with)
 - If ~~love to~~ remember just one thing:

When things don't work, Keep calm and Lower your Learning rate
 - ADAGRAD: a modification of SGD that implicitly does momentum and LR decay
⇒ Less sensitive to hyperparameters.
 - Mini-Batching
 - train on a subset of data instead of all
 - train model even if the memory is not enough to contain the entire DS.
 - quite useful with SGD:
 - shuffle the data at the start of each epoch, then create minibatches
 - Setting batches in TF
 - the data might not be dividable by batch size
 - ⇒ use tf. placeholder
- features = tf.placeholder(tf.float32, [None, n_inputs])
 labels = tf.placeholder(tf.float32, [None, n_classes])
-
- TF Lab Notebook:
- Epoch: A single forward and backward pass over the whole DS.
- Lesson 12: Deep Neural Networks
- Min-Max Scaling $x' = a + \frac{(x - x_{\min})(b-a)}{x_{\max} - x_{\min}}$
bring the data to a range of $[a, b]$
 - Linear Model Complexity

N inputs } $(N+1)K$ parameters
 K outputs
 - Linear model limitations
 - # of params & too few, we want many more
 - Linearity:
 $y = x_1 + x_2$ ✓
 $y = x_1 * x_2$ ✗
 - But, linear models have advantages too:
 - Stable: $y = WX \rightarrow \Delta y \propto \|W\| \Delta x$ bounded (125)
 - Derivatives are constant: $y = WX \rightarrow \frac{dy}{dx} = W^T, \frac{dy}{dw} = x^T$

- ⇒ We want to keep out parameters inside big linear functions.

$$Y = W_1 W_2 W_3 X = W X$$

But we would also like our entire model to be non-linear.

⇒ add non-linearities

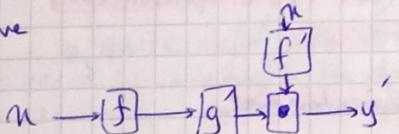
ReLU (Rectified Linear Unit)

- Lazy engineer's favorite non-linear function?
- $f(x) = \max(0, x)$
- `tf.nn.relu()`

Chain Rule Computational Graph

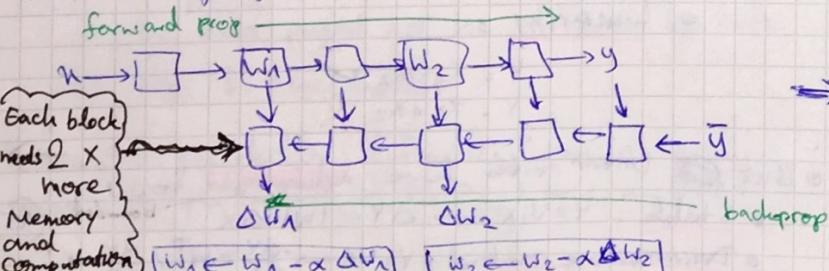
Function $n \rightarrow [f] \rightarrow [g] \rightarrow y$

Derivative



- Efficient data pipeline
- Lots of data reuse
- Simple data pipeline

Back-PROP



tf. train. Saver

Regularization

why DL now?

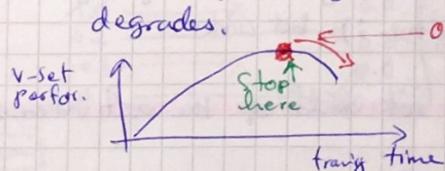
- become large enough DSs available
- because we know better ways to regularize data.

General issue with Numerical Optimization:

- the skinny jeans problem
- Training a network that is just the right size for your network is very hard to train.
- ⇒ we most often train networks that are way too big for our data, and then try our best to prevent them from overfitting

Early Termination

- Check performance of validation set and stop training as soon as performance degrades.



Regularization

Prevent overfitting (127)

④ Regularization :

- applying artificial constraints on the network that implicitly reduce the number of free params while not making it more difficult to optimize
- stretch pants in skinny jeans analogy :)

⑤ L2 Regularization :

- add another term to Loss functions which penalizes large weights.

$$L' = L + \beta \frac{1}{2} \|w\|_2^2$$

- adding L_2 norm of weights to loss, multiplied by a small constant.

- + it is very simple and does not need to change the network
- + its derivative is also very simple: w

⑥ Dropout :

- developed recently, works amazingly well
- looks insane :), but works
- randomly set half of the activations to zero
 \Rightarrow randomly destroy half of the information that flows through the network. Sounds crazy
- \Rightarrow Network can never rely on any given activation to be present \Rightarrow it is forced to learn a redundant representation for everything!



- Forcing the network to learn a redundant representation of everything
 - makes things more robust
 - prevents over fitting
 - it is as if the resulting network is taking a consensus over an ensemble of networks.

Dropout

- One of the most important techniques to emerge over the last few years.

! - If it doesn't work for you, you should probably be using a bigger network!

Evaluating a network trained with dropout:

- we obviously don't want the randomness when evaluating; we want something deterministic
- we will evaluate the consensus over these redundant models
- we average the activations

Training

$$\begin{array}{rcl} 1.0 & \xrightarrow{\quad X \quad} & 0.0 \\ 0.2 & \xrightarrow{\quad X \quad} & 0.2 \times 2 \\ -0.3 & \xrightarrow{\quad X \quad} & 0.0 \\ 0.5 & \xrightarrow{\quad X \quad} & 0.5 \times 2 \end{array} \left. \right\} y_t$$

Evaluation

$$\begin{array}{rcl} 1.0 & \longrightarrow & 1.0 \\ 0.2 & \longrightarrow & 0.2 \\ -0.3 & \longrightarrow & -0.3 \\ 0.5 & \longrightarrow & 0.5 \end{array} \left. \right\} y_e$$

$$y_e = E(y_t)$$

- \Rightarrow During training, scale the remaining activations by factor $\times 2$

• tf.nn.dropout()

- ↳ TF multiplies the units that are kept by $\frac{1}{\text{keep-prob}}$
- a good starting value for keep-prob 0.5
- during testing: keep-prob = 1.0

12.04.
20:15
2 h

Lesson 13: Convolutional Neural Networks

- If your data has some structure, and your network doesn't have to learn that structure from scratch, it performs better.

• Statistical Invariance

- e.g. translation invariance in image
- things that does not change over time and space

• Weight Sharing

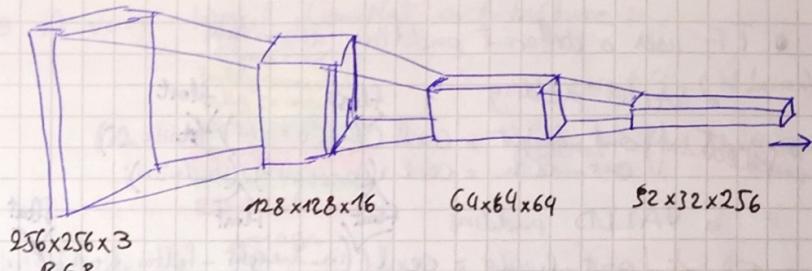
- When you know that two inputs can contain the same kind of information, then you share their weights, and train the weights jointly for those inputs

deep learning

• Weight sharing leads to CNNs, same for images and to embeddings and RNNs for text and general sequential data

• CONVNETS (Convolutional Networks)

- Share weights across space
- progressively squeeze the spatial dimension of the image while increasing the depth, which corresponds roughly to the semantic complexity of your representation



→ Classifier

- We have now a representation where all the spatial information have been squeezed out and only parameters map to content of the image (like Cat) remain.

- Patch = Kernel = filters

- Feature map = channels

- Stride: how much to shift the filter
stride of two \rightarrow input size is halved

- Padding:

↳ Valid padding: don't go past the edge

↳ zero padding: go off the edge and pad with zeros (same padding)
because size of output will be the same as input in case of stride = 1

- filter depth:

- different filters pick up different quality of a patch.

- Dimensionality (Size of model parameters)?

↓ ↓ ↓
↓ ↓ ↓
↓ ↓ ↓

Given: input $W \times H$, filter size F , stride S , Padding P , #filters = K

$\Rightarrow W_{\text{out}} = \lceil (W - F + 2P) / S \rceil + 1$ $H_{\text{out}} = \lceil (H - F + 2P) / S \rceil + 1$

$D_{\text{out}} = K$, output volume = $W_{\text{out}} \cdot H_{\text{out}} \cdot D_{\text{out}}$

- TF uses a different padding algorithm

◦ SAME padding float float
 $\text{out_height} = \text{ceil}(\text{in_height}) / (\text{stride_h})$
 $\text{out_width} = \text{ceil}(\text{in_width}) / (\text{stride_w})$

◦ VALID padding float float float
 $\text{out_height} = \text{ceil}((\text{in_height} - \text{filter_h} + 1) / (\text{stride_h}))$
 $\text{out_width} = \text{ceil}((\text{in_width} - \text{filter_w} + 1) / (\text{stride_w}))$

• Deep Visualization toolbox

• TF CNNs

- `tf.nn.conv2d()`
- `tf.nn.bias_add()`

• Pooling

- Stride of 2 is a very aggressive way to down sample an image. It removes a lot of information.
- What if instead of skipping one in every two convolutions, we still run with a very small stride, say e.g. 1, but then take all the convolutions in the neighbourhood and combine them somehow.

• Max Pooling

- does not add to the number of parameters
 \Rightarrow don't increase the risk of overfitting.

- often yields more accurate models.

- model becomes more expensive to compute

- more hyperparams: `pooling_size`
`pooling_stride`



- A very typical ConvNet architecture:

Classifier
 fully connected
 fully connected
 max pooling
 convolution
 Max pooling
 convolution
 image

LeNet '5 ('88 LeCun)

AlexNet ('12 Alex Krizhevsky)

• Average Pooling

- take avg instead of max
- like a blurred, low resolution view of the feature maps below

• `tf.nn.max_pool()`

- pooling decreases the size of the output

\Rightarrow prevents overfitting

by
 reducing the number of parameters in
 future layers.

• Pooling layers have fallen out of favor recently:

- recent DNNs are so big and complex, that we are more concerned with underfitting
- Dropout is a much better regularizer

- Pooling results in loss of information. E.g. in max pooling on n numbers, we discard information contained in n-1 of them.

• 1×1 convolution

- why would we ever want a 1×1 convolution?
 They are not really looking for a patch in image
 but just a pixel.

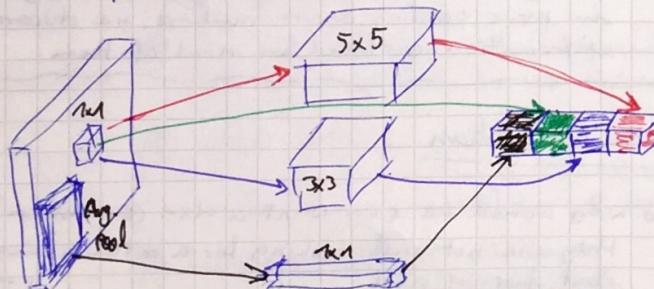
- ④ A normal convolution is basically a small classifier for a patch of the image, but it is only a linear classifier
- ⑤ by adding a 1×1 conv in the middle, we suddenly have a neural network instead of a linear classifier
- ⑥ it is a very inexpensive way to make models deeper and have more parameters without completely changing the structure
- ⑦ they are also very cheap. it is actually only a matrix multiplication with relatively few parameters

⑧ Avg pooling, 1×1 convs

⇒ a very successful strategy for building ConvNets that are better and smaller than convnets that simply use a pyramid of convolutions.

⑨ Inception Modules

- ⑩ instead of making a choice at each layer between pooling, convolution, 1×1 , 3×3 , 5×5 , let's have a combination of all of them
 - ⑪ a composition of Avg Pooling followed by 1×1 conv
 - a 1×1 conv
 - or 1×1 followed by a 3×3
 - a 1×1 followed by a 5×5
- and concatenate the output of each of them



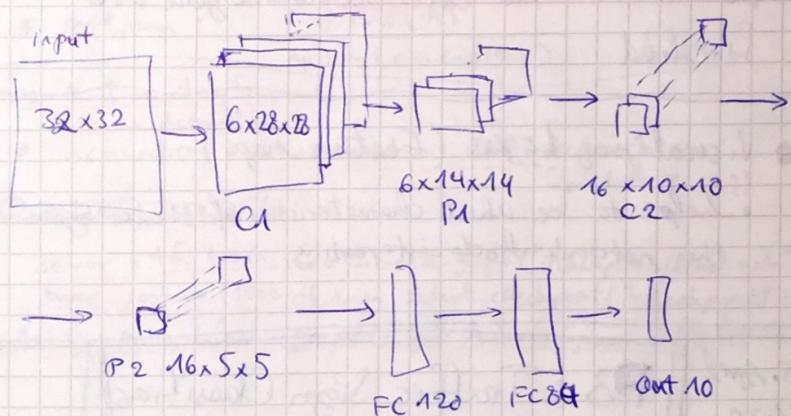
- ⑫ It looks complicated, but we can choose those parameters in such a way that the total number of params in network is very small, but the performance is much better than having individual convolutions only.

⑬ LeNet

- ⑭ always a good starting point

⑮ tf.contrib.layers.flatten

⑯ How many parameters LeNet?



⑰ Number of Params

$$\begin{aligned}
 & [6 \times (5 \times 5 \times 1) + 6] + [16 \times (5 \times 5 \times 6) + 16] + [(400 \times 120) + 120] \\
 & + [(120 \times 84) + 84] + [(84 \times 10) + 10] = 61706
 \end{aligned}$$

⑱ TensorBoard for visualizing networks

15. April

17-19

(2h, actually only 30 min)

Lesson 14: LeNet for Traffic Signs

- Always shuffle the data, otherwise the ordering in the data can have a major impact on the training
- training pipeline / evaluation pipeline
- sklearn. train-test-split()
 - 20% is a good rule of thumb for validation set
- Don't run the test set until you are finished

Visualizing Layers (Feature maps)

- helps to see which characteristics of the image the network finds interesting.

15. April

21:30

26. April

P3: Traffic Sign Classifier

- dataset distribution
- normalization (pixel - 128) / 128 (zero mean)
- shuffle
- Learning rate: $0.1 \times 0.01 \times 0.001 \checkmark$
- increasing batch size.

- Epochs, perceptron oscillation and stagnation on validation set
 - early stopping. (Dropout)

restore model

↳ don't init ~~all~~ variables again

construction code (network layers,)
variables

saver = tf.train.Saver()
with tf.Session() as sess:

saver.restore("./checkpoint/---")
sess.run(logits, feed_dict={---?})

softmax :

$$\frac{\exp(n)}{\exp(n)}$$

sess = tf.InteractiveSession()
sess.as_default()

activation = tf.activation.eval(session=sess,
~~feed_dict={---?}~~)

~~sess = tf.Session()~~

saver = tf.train.import_meta_graph("./checkpoints/lenet.meta")
saver.restore(sess, tf.train.latest_checkpoint("./checkpoints"))

const = tf.get_default_graph().get_tensor_by_name("C1:0")

26. April
18:30

15
26. April
17. 20:30
(2)

Lesson 16: Keras

2h

otto, Voyage

Andrew Gray
& Uber ATG

Behavioral Cloning

End-to-End Learning:

- Network learns the steering wheel angle and speed using only the input from sensors.

Why DL over Robotic approach?

- Both approaches are used.
- DL has the potential to fundamentally change the way SDCs are programmed.
- Robotic approach involves programming a lot of detailed knowledge about sensors and rules and control into the car.
- With DL on the other hand, we feed all the data we have and let the network decide what is important and what is not.
 - + with DL we have a feedback loop:
The more we drive, the more data we collect, and the better becomes the network.

Keras:

- a front-end to TF that makes creating networks much faster and cleaner.
- Sometimes if you need very fine-grained control you need to create the layers directly using TF.

• keras.model.Sequential()

- a wrapper for NN model
- fit(), evaluate(), compile()

• two ways to build Keras models:

- sequential
- functional

• Sequential(), Flatten(), Activation(), Dense(), Conv2D(), MaxPooling2D(), Dropout!

↑
dropout prob
not
keep-prob

• sklearn.preprocessing.LabelBinarizer()

• model.fit(x-normalized, y-one-hot, epochs, validation_split)

• model.compile("adam", "categorical_crossentropy", ["accuracy"])

• Lambda layers for preprocessing

• Cropping2D()

26. April

20:30

27. April

22:00

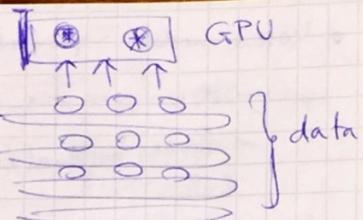
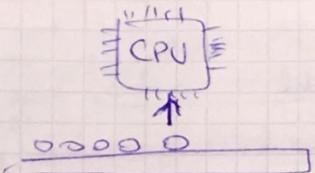
L17: Transfer Learning

(VP. DL)

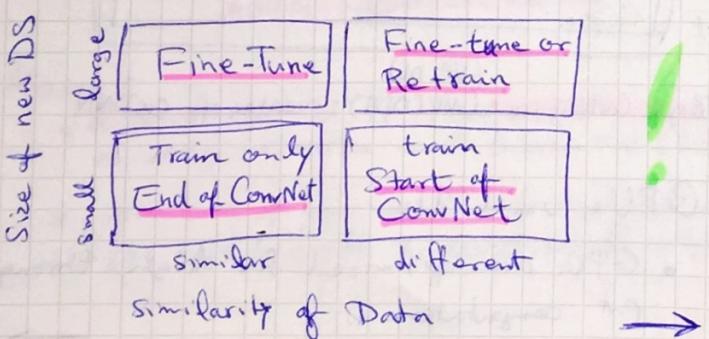
• Bryan Catanzaro (NVIDIA): creator of cuDNN

• GPUs vs. CPUs

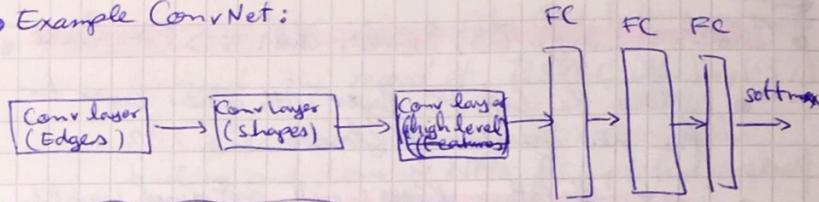
- GPUs are optimized for high throughput computation



- CPUs are optimized for latency: running a single thread of instructions as quickly as possible
- GPUs are optimized for throughput: running as many simultaneous operations as possible.
- Rule of thumb
 - Training is 5X faster on a GPU.
- Check graph on "CPU vs. GPU GigaFLOPs."
- Transfer Learning

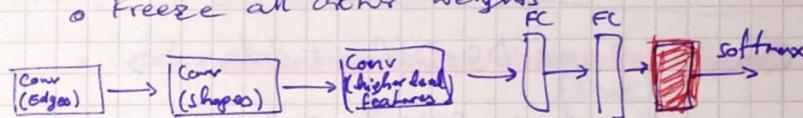


• Example ConvNet:



• Case 1: Small DS, Similar Data:

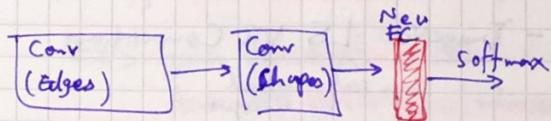
- Cut-off the last FC Layer and retrain it for your problem (with new number of classes)
- Freeze all other weights



- both DSs have similar higher level features
⇒ keep all the feature detector layers.

• Case 2: Small DS, Different data

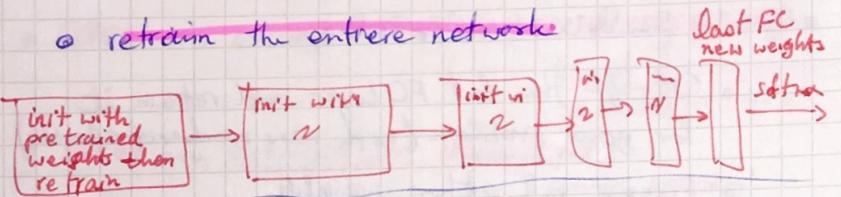
- Keep only few of the Conv layers at the beginning (Feature detectors for general features, like edges and shapes)
- the new DS has different higher level features than the original DS
⇒ The conv layer for higher level features is not useful.
- Add a new FC layer for the new number of classes and train its weights.



→ Finet

Case 3: Large DS, Similar data \Rightarrow Fine-tune

- replace the last FC layer with a new layer for the new number of classes.
- init all other layers with pretrained weights
- retrain the entire network



Case 4: Large DS, Different data \Rightarrow

- either do as case 3, or
- or retrain entire net with random weights.

Rise of DL

- huge labeled DSs because of Internet
- greater computational power.

Keras Applications: import a pre-trained network:

Xception, VGG16, VGG19, ResNet, ResNetV2, InceptionV3, InceptionResnetV2, MobileNet, MobileNetV2, DenseNet, NASNet

ImageNet: A large database of hand-labeled Images.

- ImageNet LS VR Competition

Large scale visual
Recognizer

- AlexNet (2012) parallel GPUs
 - first to use ReLU
 - first to use Dropout
- improved the record of ILSVRC (26.) by (15%) error error
- still used today as a starting point for building and training NNs. (Simplified versions of it)

VGG

- visual geometry group @ oxford / ILSVRC error 7%
- simple and elegant arch.: great for transfer learning
 - a long sequence of 3×3 convolutions broken up by 2×2 pooling layers a trio of FC layers.

- used a lot as starting point; works really well. flexible

Advice, Wisdom

- DL is a hands-on empirical field
 - \rightarrow more experiments \rightarrow more ideas
 - practice is most of the times ahead of theory
 - \Rightarrow don't get too attached to theory and network architectures

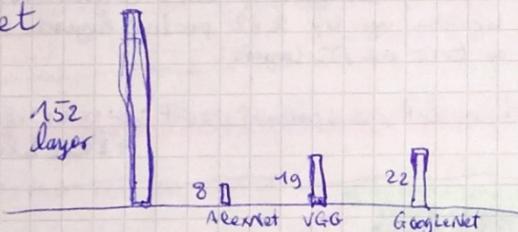
- To be a successful DL Practitioner you need to have a bias towards Action
 - \Rightarrow Experiment, try

- DL is actually one of the few places in CS that we actually do science \rightarrow experiment
- Get your hands dirty; Try to solve real problems
 - \Rightarrow gain practical skills develop your confidence

GoogleNet / Inception

- very fast
 - o shallowest as fast as AlexNet
 - great accuracy
 - great choice if you want to run your network in real-time
- Inception modules
⇒ total number of parameters is small.

ResNet



similar to VGG

- same structure, repeated again and again
- short cut connections (skip layers)
- less on ImageNet 3%!
 ⇒ better than human

Freezing layers

⊕ memory usage, training speed

• `layer.trainable = False` / `model.layers`

`inception.summary()`
`model.pop()` / `include_top=False` / ^{top} `2 layers`

• Sequential vs. Model API

• CIFAR-10 Dataset

GlobalAveragePooling 2D()

callbacks

• ModelCheckpoint()

• EarlyStopping()

ImageDataGenerator()

o `model.fit_generator()`, `model.fit()`
with predict_generator() large DS
`Model.predict()`

- whole DS at once.
- use if whole DS fits into memory

28.April
05. May
(7 days)

P4: Behavioral Cloning

Lambda layers for preprocessing.

`model.add(Lambda(lambda, n: (x/255.0) - 0.5, input_shape=(160, 320, 3)))`

Cropping2D()

~~model.add(Cropping2D(cropping=((7, 25), (0, 0)))~~

• Simulator: 2 modi training mode
auto race mode

• press R to record

• drive on the center

• drive recovery if car drives off the center

• drive counter clockwise to generalize model

• Flipping images to augment data

• collect the data in the second track to help the model to generalize

• avoiding over fitting and under fitting

• knowing when to stop collecting data

• Data
 $\underbrace{\text{center image}, \text{left}, \text{right}}_{\text{features}}, \underbrace{\text{steering}, \text{throttle}, \text{brake}, \text{speed}}_{\text{labels}}$
 $[0, 1]$ $(0, 1)$ $(0, 1)$ $(0, 1)$

- Merge RGB, 160×320
 - * may need to crop top and bottom of the image (sky, trees, hills, car's hood). distract the model
- left and right images help the network learn how to recover when off the center
- regression model
 - * loss mse
- Save model. h5

```

input CSV
lines = []
with open("../data/driving-log.csv") as csv_file:
    reader = csv.reader(csv_file)
    for line in reader:
        :
images = []
meas = []
y_train = np.array(meas)
model = Sequential()
model.add(Flatten(input_shape=(160, 320, 3)))
model.add(Dense(1))
model.compile(loss="mse", optimizer="adam")
model.fit(x_train, y_train, validation_split=0.2,
          shuffle=True)
model.save("model.h5")
    
```

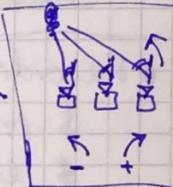
- drive-by model. h5 / underfitting
- improve model:
 - * increase epochs
 - * add more convolutions

overfitting

- * use dropout and pooling
- * collect more data
- * fewer conv and FC



- Lambda layer: Preprocessing
 - * normalization: pixel / 255
 - * mean center: pixel-normalized ≈ 0.5
- Lambda (lambda x: (x / 255) - 0.5)
- LeNet
- Track and sense the training data is biased towards steering to the left
 - * np.random.choice(image)
 - * Data augmentation: flip the images
 - * Steering = -steering
- \Rightarrow Balanced data set
- Use left and right images for learning how to recover from off-center
 - correction = 0.2
 - steer-left = steer-center + correction
 - steer-right = steer-center - correction
 - images.extend(left \rightarrow right)
 - steerings.extend(steer-left, steer-right)
- Cropping images
 - * Cropping 2D () layer
- NVIDIA network
- Collecting more data
 - * driving counter-clockwise \Rightarrow balanced data \Rightarrow better
 - * driving other track \Rightarrow model generalizes
- How much data?
 - * loss \rightarrow MSE
 - * loss high on training and validation? \Rightarrow underfitting
 - * loss high on training but low on validation? \Rightarrow overfitting
 - * loss low on training but high on validation? \Rightarrow high loss
- two or three laps center driving
- one lap recovery driving from center, one lap driving smoothly on curves



• Outputting training and validation loss metrics.

- verbose option for fit() and fit_generator()

1 2 → loss on training and valid set after each epoch

- model history object : contains the training & validation loss for each epoch

-
- Generators ↗ Coroutines ↘ ↙
 - large data : don't store the whole DS at once in memory ;
 - pull chunks of data
 - yield keyword in python
-

• Recording video

```
$ python drive.py model.h5 run1  
$ ls run1  
$ python video.py run1 --fps 48
```

05. May
22-00:30

0:30m

Lesson 20: Sensors

- Michael Maile : Director of sensor fusion team at Mercedes Benz.

- Andrei Vatavu

- Dominik Nauß → Luminar (lidar)
↳ Ulm uni.

- Kalman Filter : Extended Kalman Filter
 - track a pedestrian relative to the car with help of radar and lidar data

- Sebastian Thrun : AI for Robotics course

- C++ to fuse ~~Radar~~ & lidar data

- Unscented Kalman Filter:
 - non-linear tracking

-
- Prototype SDC of Mercedes Benz (choobaka:))

- 15 sensors :

- a stereo camera rig (2 cameras)
- a long range lens camera for detection of traffic lights (1 camera)
- radars (behind the front bumper)
- lidar

- Radar uses Doppler effect to measure speed directly. → Important for Sen. Fus. → (fig)

← Doppler effect very important for
Sens. Fns. because it gives us the velocity as
an independent measured parameter

⇒ Fusion Algorithms converge much faster

- Radar also used for localization by generating the radar map of the environment.

- Because lidar signals waves bounce off hard surfaces they can provide measurements to objects without direct line of sight

e.g. see underneath other vehicles
on spot objects and buildings that would be obscured otherwise.

- Radar:

- most effective by rain and fog
- wide FoV about 150°
- a long range 200+ meters

	Cam.	Lidar	Radar
Resolution	○	○	○ ← low res. spec. vertical
Noise	○	○	○
Velocity	○	○	○
All-weather	○	○	○
size	○	○	○

- Low Res. of Radar ⇒ reflection from static objects can cause problems. [Radar Clutter]

RADAR: Radio Detection and Ranging

LIDAR: Laser Detection and Ranging [Velodyne]
For MB.
SDC

- Lidar:

- Infrared laser beam (900 nm wave length)
 - some use longer WLs ⇒ better in rain and fog
- much higher spec. res. than Radar
 - larger numbers of scan layers in ver. direct
 - high density of points per layer.
- cannot measure velocity directly
 - have to rely on position change between two or more scans.
- more affected by weather and dirt on sensor
 - ⇒ requires keeping them clean.
- bulkier

06. May, 21:00
08. May, 21:00

Lesson 21: Kalman Filters

- Stanford automotive research center.
- Jammin: SDC of Stanford, child of Stanley, first SDC that is now in National Museum of computing history.
- Lidar and sensor data to track other cars to avoid colliding with them
 - not only know where the other cars are (localization)
 - but also how fast they are moving!

• Avoid collision for future : tracking

• Kalman filters

• Estimate continuous state of a system

Kalman Filter

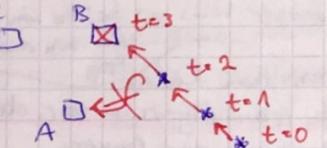
Continuous
Unimodal

Monte Carlo

Discrete
Multi-modal

Particle Filters

Continuous
Multi-modal



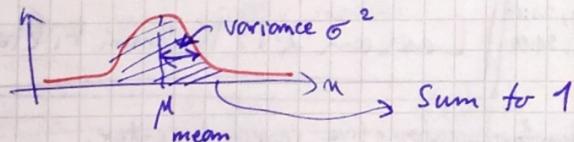
• Assuming no drastic change of velocity (A),

Kalman filters

Noise → • takes observations

→ • estimates future locations and velocities

• In Kalman Filters the dist is given by a Gaussian



1-D Gaussian (μ, σ^2)

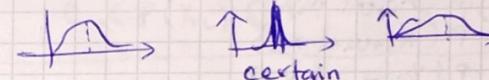
• Task in Kalman Filters:

to maintain a μ and σ^2 that are best estimate of the location of the object

$$f(m) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp^{-\frac{1}{2} \frac{(m-\mu)^2}{\sigma^2}}$$

• Unimodal → one peak

σ^2 : certainty. The larger the covariance, the more uncertain we are.



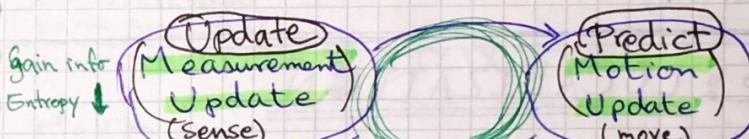
• Gaussian as Belief in Kalman filters.

• $m = \mu \Rightarrow$ Gaussian value is max: $\frac{1}{\sqrt{2\pi\sigma^2}}$

Kalman Filters: [recap localization montecarlo (histogram) filters]

• represent the distribution in Gaussians.

• iterate on two main cycles:



Loose Info
Entropy ↑

- requires product
- uses Bayes Rule

World = ['green', 'red', 'green', 'green', 'green']

Z = 'red' # measurement

p Hit = 0.6 } # sensor may not read correctly (sensor prior)
p Miss = 0.2 }

def sense(p, Z):

q = [] ← measurement prior

for i in range(len(p)):

hit = (Z == world[i])

q.append(p[i] * (hit * pHit + (1 - hit) * pMiss))

return q ← new belief (posterior)

normalization:

S = sum(q)

for i in range(len(q)):

q[i] = q[i] / S

(not normalized:
the probabilities do not add to one. Needs normalization!)

Measurement
Update:

- take measurement and prior
- update the belief.
- output posterior

Motion Update (move() function):

$$\begin{aligned} p_{\text{Exact}} &= 0.8 \quad \# P(X_{i+0} | X_i) \\ p_{\text{Overshoot}} &= 0.1 \quad \# P(X_{i+1} | X_i) \\ p_{\text{Undershoot}} &= 0.1 \quad \# P(X_{i+0-1} | X_i) \end{aligned}$$

def move(p , v):
 prior belief of where
 number of steps
 the robot is

```
a = []
for i in range(len(p)):
    q-exact = p[(i-v)%len(p)] * pExact
    q-overshoot = p[(i-v-1)%len(p)] * pOvershoot
    q-undershoot = p[(i-v+1)%len(p)] * pUndershoot
    q.append(q-exact + q-overshoot + q-undershoot)
return q
```

← posterior

Recap : LOCALIZATION

- The localization maintains a function over all possible places a robot might be, where each cell has an associated probability value

~~BELIEF~~ = PROBABILITY

SENSE = PRODUCT ← measurement update:
followed by
normalization

take those probabilities
and multiply them up or
down based on the exact
measurement.

MOVE = CONVOLUTION (= ADDITION)

For each possible location
after the motion, we
reversed engineered the
situation, and guessed where the robot might have
come from, and then collected (added) the
corresponding probabilities

Formal Definition of Localization

$$0 \leq P(X_i) \leq 1$$

$$\sum P(X_i) = 1 \quad \text{for } X_i \in \text{World} \quad (\text{a grid cell})$$

Measurements

[Bayes Rule]

X : grid cell Z : measurement

- the measurement update seeks to calculate a belief over my location after seeing the measurement.

$$P(X_i | Z) = \frac{P(Z | X_i) \cdot P(X_i)}{P(Z)}$$

Prior
multiplied by
 $P(Z | X_i)$: the chance
of seeing a red or
green cell for
every possible location

probability of seeing a measurement
devoid of any location information

our normalization
term gives this

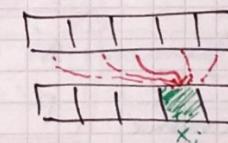
$$P(Z) = \sum_i P(Z | X_i) \cdot P(X_i)$$

Motion

Total Probability

$$P(X_i^t) = \sum_j P(X_j^{t-1}) \cdot P(X_i | X_j)$$

↑ time
grid cell



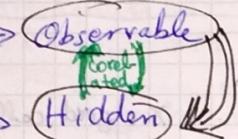
$$P(A) = \sum_B P(A | B) \cdot P(B)$$

Total probability
Convolution
(weighted sum)

check the Kalman Filter
notes from Intro to SDC course.

- Kalman Filters allow inferring a hidden variable, like velocity, from measurements, like (x, y) position

KF States



allows inference as a result of H's correlation

- For design of a KF we need two things:

- State transition function

F

$$\begin{pmatrix} \dot{x} \\ \dot{v} \end{pmatrix} \leftarrow \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ v \end{pmatrix}$$

$$\begin{array}{l} \dot{x} = 1 \\ \dot{v} = 0 \end{array} \quad \begin{array}{l} x' \leftarrow x + v \\ v' \leftarrow v \end{array}$$

- Measurement function

H

$$z \leftarrow \begin{pmatrix} 1 & 0 \end{pmatrix} \begin{pmatrix} x \\ v \end{pmatrix}$$

\hat{x} : estimate

F: state transition matrix

u: motion vector

P: process uncertainty covariance (error cov. matrix)

z: measurement

H: measurement function

R: measurement noise

K: Kalman Gain

(error)
we compare the measur.
with our prediction

projecting sys. uncertainty
into the measurement space

$$K = P \cdot H^T \cdot S^{-1}$$

$$\begin{aligned} x' &= x + (Ky) \\ P' &= (I - K \cdot H) \cdot P \end{aligned}$$

Predict:

$$\begin{aligned} \hat{x}' &= F\hat{x} + Bu \\ \hat{P}' &= FPFT + Q \end{aligned}$$

Process noise cov. matrix

Update

$$y = z - H\hat{x}$$

$$S = R + HP^T H^T$$

$$K = P \cdot H^T \cdot S^{-1}$$

$$\hat{x} = \hat{x}' + Ky$$

$$P = (I - KH)P'$$

$$y = z - H\hat{x}$$

May 09.
22:00
May 14.
21:00

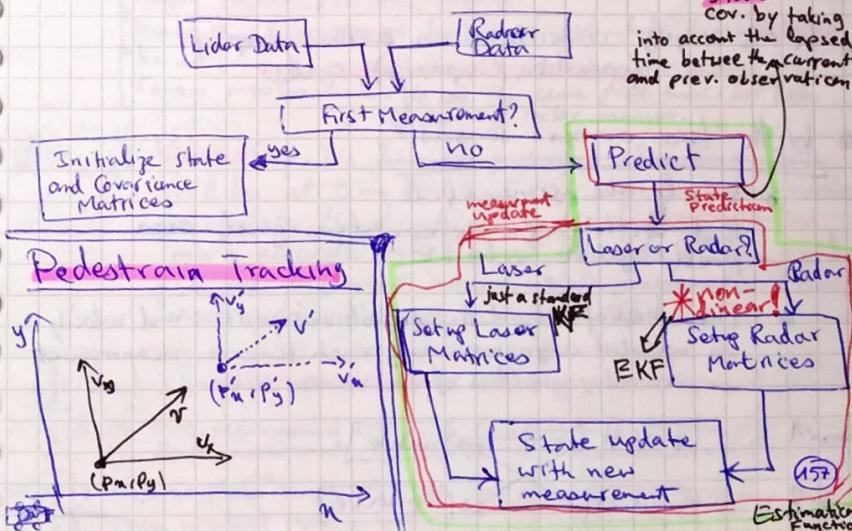
4 hours

Lesson 24: Extended Kalman Filters

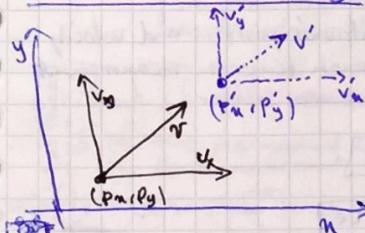
- **Extended KF (EKF):** A nonlinear version of Kalman Filter which linearizes about an estimate of the current mean and covariance.

- Structure of Sensor Fusion using KFs

Predict pedestrian state and its cov. by taking into account the elapsed time between the current and prev. observations



Pedestrian Tracking

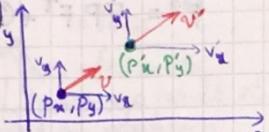


157

Estimation Function

• Overview: Sensor Fusion with KF

- We want to track a pedestrian moving in front of our AV (autonomous vehicle)



- The KF uses a two-step Estimation problem

use information we have so far to predict the state of the pedestrian until the next measurement arrives.

**State prediction
1. (Predict)**

**measurement update
2. (Update)**

we new obs. to correct our belief about state of pedestrian

- we always start with what we already know about the pedestrian

- we use that information to predict the state of the pedestrian up to the time when the next measurement arrives. This is called prediction step.

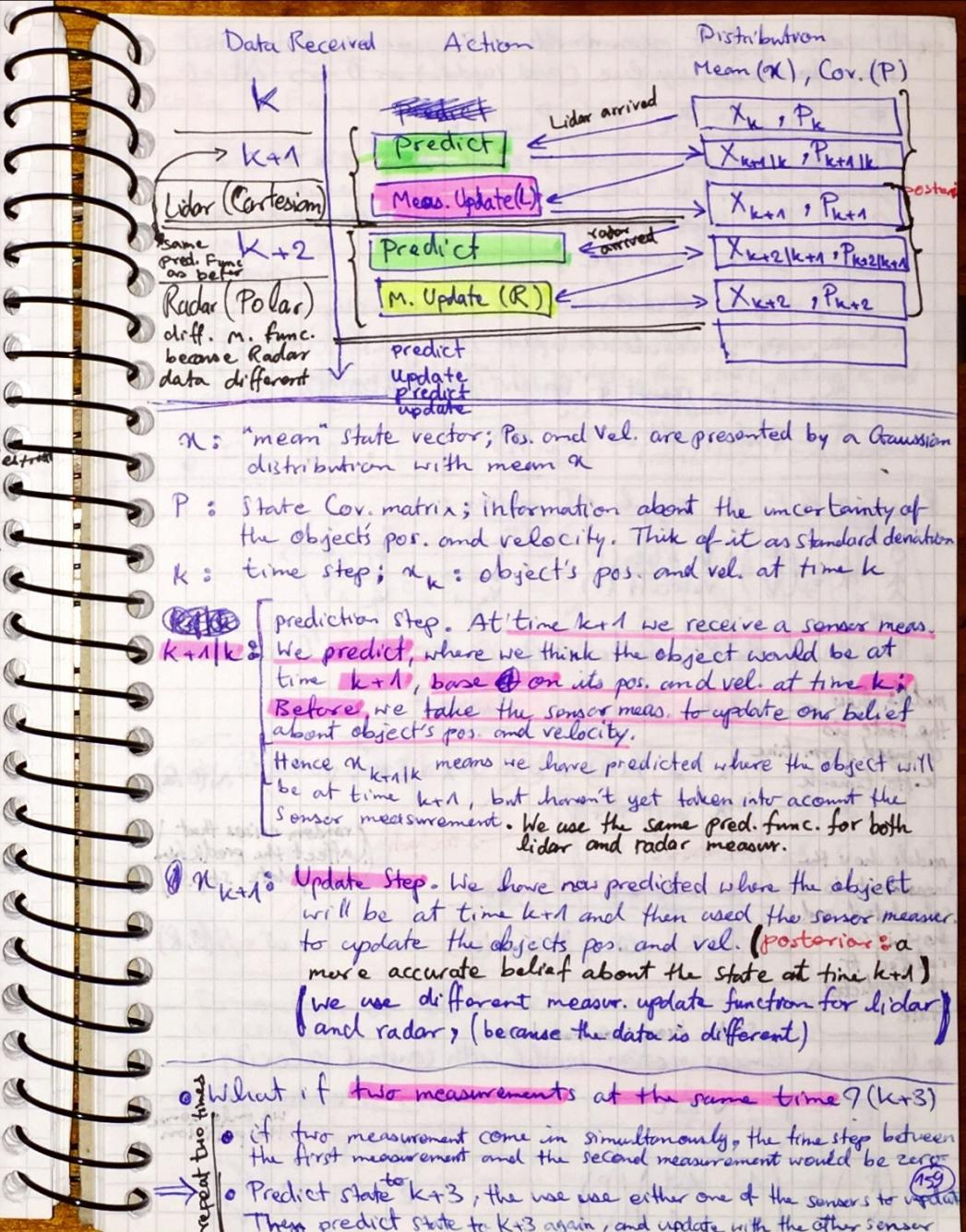
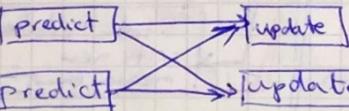
Prediction Step

- Update Step**: we use the new observation to correct our belief about the state of the pedestrian.

- KF**: simply consists of an endless loop of (predict / update) cycles.

- With two or more sensors?

- Simply the same cycle. But with each sensor having its own pred / update scheme.
- The belief about the pedestrian's position and velocity is updated asynchronously each time the measurement is received regardless of the source sensor.



← our 2D motion model is: $\dot{v}_x, \ddot{v}_x, \dot{v}_y, \ddot{v}_y \rightarrow \frac{\Delta v}{\Delta t} = \frac{v_{t+1} - v_t}{\Delta t}$

$$\begin{aligned} P'_x &= P_x + v_x \Delta t + \frac{a_x \Delta t^2}{2} \\ P'_y &= P_y + v_y \Delta t + \frac{a_y \Delta t^2}{2} \\ V'_x &= V_x + a_x \Delta t \\ V'_y &= V_y + a_y \Delta t \end{aligned}$$

deterministic part

Stochastic: acceleration is considered in noise component.

$$N \in \mathcal{N}(0, Q), \quad a_x \sim \mathcal{N}(0, \sigma_{ax}^2), \quad a_y \sim \mathcal{N}(0, \sigma_{ay}^2)$$

$$V = \begin{pmatrix} a_x \Delta t^2 \\ a_y \Delta t^2 \\ a_x \Delta t \\ a_y \Delta t \end{pmatrix} = \begin{pmatrix} \frac{\Delta t^2}{2} & 0 \\ 0 & \frac{\Delta t^2}{2} \\ \Delta t & 0 \\ 0 & \Delta t \end{pmatrix} \begin{pmatrix} a_x \\ a_y \end{pmatrix} = G a$$

G : does not contain random variables.
 a : contains random acc. components.

process Cov. matrix:

$$Q = E[V V^T] = E[G a a^T G^T]$$

because G does not contain random variables (it is constant)

$$= G E[a a^T] G^T = G \begin{pmatrix} \sigma_{ax}^2 & \sigma_{axy} \\ \sigma_{ayx} & \sigma_{ay}^2 \end{pmatrix} G^T$$

$$= G Q_v G^T \quad \text{expect. of } a \text{ and cov.}$$

$$\Rightarrow Q_v = \begin{pmatrix} \sigma_{ax}^2 & \sigma_{axy} \\ \sigma_{ayx} & \sigma_{ay}^2 \end{pmatrix} : \text{three statistical moments}$$

expectation of a_y (cov.)

$\sigma_{axy} := 0$ because a_x and a_y are assumed to be uncorrelated noise processes

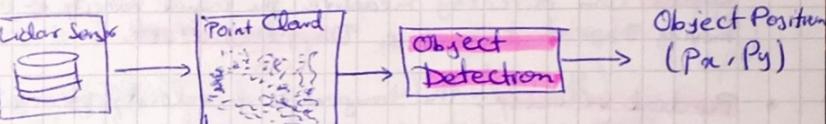
$$\Rightarrow Q_v = \begin{pmatrix} \sigma_{ax}^2 & 0 \\ 0 & \sigma_{ay}^2 \end{pmatrix} \Rightarrow Q = G Q_v G^T = \begin{pmatrix} 4 \times 4 \end{pmatrix}$$

(see prev. page)

Now Measurement Model for 2D

measur. vectors Z
 measur. matrix H
 cov. matrix R

- The laser sensor output is a point cloud. But here we assume we already analyzed the point cloud to compute the 2D location of the pedestrian.



$$Z = \begin{pmatrix} P_x \\ P_y \end{pmatrix}, \quad \text{state } a = \begin{pmatrix} P_x \\ P_y \\ V_x \\ V_y \end{pmatrix} \quad Z = H \cdot a + w$$

- H : measurement matrix projects our belief about the object's current state into the measurement space of the sensor. For lidar, this means, we discard the velocity information, since the lidar only measures the position. The state vector a contains information about (P_x, P_y, V_x, V_y) . The measurement vector Z contains only information about (P_x, P_y) . Multiplying them allows us to compare a , our belief, with z , the sensor measurement.

$$\Rightarrow H = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}$$

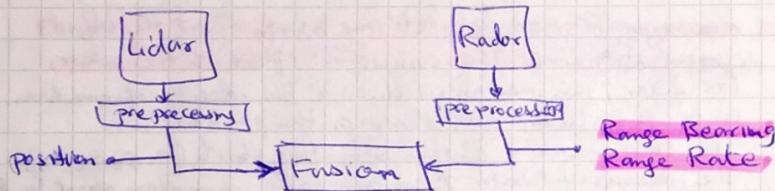
- R : measurement cov. matrix. Uncertainty of sensor about its measurements. Square matrix, with each side the same as the number of measurement parameters.

$$R = E[W W^T] = \begin{pmatrix} \sigma_{px}^2 & 0 \\ 0 & \sigma_{py}^2 \end{pmatrix}$$

- usually the parameters for the random noise measurement matrix will be provided by the sensor manufacturer.
- off-diagonal 0 s in R indicate that the noise processes are uncorrelated.

Radar Measurements:

- LIDAR → preprocessing → Post-tran (but no velocity directly)
- RADAR → preprocessing → radial velocity
- Radar can directly measure the radial velocity of a moving object using the **Doppler Effect**.
- Radial velocity: the component of velocity moving towards or away from the sensor.
- But Radar has a lower ~~radial~~ spatial resolution, than Lidar. ⇒ Combine both sensors together.



- State transition function will be the same as in the lidar case

$$\mathbf{x}' = f(\mathbf{x}) + \mathbf{v}$$

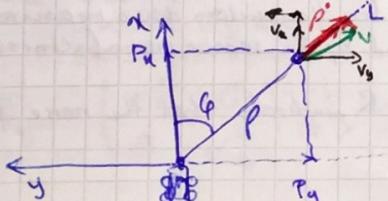
$$\mathbf{x} = \begin{pmatrix} p_x \\ p_y \\ v_x \\ v_y \end{pmatrix}$$

- measurement function is different

$$\mathbf{z} = h(\mathbf{x}') + \mathbf{w}$$

same state vector
Same linear motion model and process noise

- The Radar sees the world differently



n-axis always points in the vehicle's direction of movement

y-axis points to the left

- instead of a 2D pose, the Radar can directly measure
- $\mathbf{z} = \begin{cases} \text{Range}(\rho) & : \text{radial distance from origin}, \\ \text{Bearing}(\theta) & : \text{angle between } \rho \text{ and } n\text{-axis} \\ \dot{\rho} & : \text{range rate } (\dot{\rho}) \text{ (radial velocity)} \end{cases}$

$$\mathbf{z} = h(\mathbf{x}') + \frac{\mathbf{w}}{\text{noise}}$$

(assuming)
Three measurement components are not cross-correlated

$$\Rightarrow \mathbf{R} = \begin{pmatrix} \sigma_p^2 & 0 & 0 \\ 0 & \sigma_\theta^2 & 0 \\ 0 & 0 & \sigma_{\dot{\rho}}^2 \end{pmatrix}$$

Radar measure. cor. matrix

$$\begin{pmatrix} \rho \\ \theta \\ \dot{\rho} \end{pmatrix} \xleftarrow{h(\mathbf{x}')} \begin{pmatrix} p_x \\ p_y \\ v_x \\ v_y \end{pmatrix}$$

\mathbf{z} : measurement vector
 \mathbf{x}' : state vector
maps the predicted state \mathbf{x}' into the measurement space.

$$h(\mathbf{x}') = \begin{pmatrix} \sqrt{p_x'^2 + p_y'^2} \\ \arctan(p_y'/p_x') \\ p_x' v_x + p_y' v_y / \sqrt{p_x'^2 + p_y'^2} \end{pmatrix}$$

$h(\mathbf{x}')$ is not linear

Derivation of $\dot{\rho}(t)$:

$$\rho(t) = \sqrt{p_x(t)^2 + p_y(t)^2}$$

$$\dot{\rho} = \frac{\partial \rho(t)}{\partial t} = \dots = \frac{p_x v_x + p_y v_y}{\sqrt{p_x^2 + p_y^2}}$$

$$\frac{\partial p_x}{\partial t} = v_x, \frac{\partial p_y}{\partial t} = v_y$$

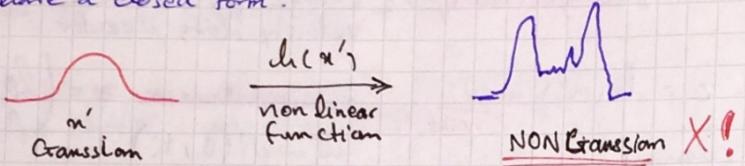
Alternative:

$\dot{\rho}$ can be seen as the scalar projection of velocity vector \vec{v} onto the $\vec{\rho}$.
 $\vec{\rho} = (p_x, p_y) \quad \vec{v} = (v_x, v_y)$
 $\dot{\rho} = \frac{\vec{v} \cdot \vec{\rho}}{\|\vec{\rho}\|} = \frac{(v_x, v_y) \cdot (p_x, p_y)}{\sqrt{p_x^2 + p_y^2}}$

- Say we have our predicted state \bar{x}' described by a Gaussian dist.

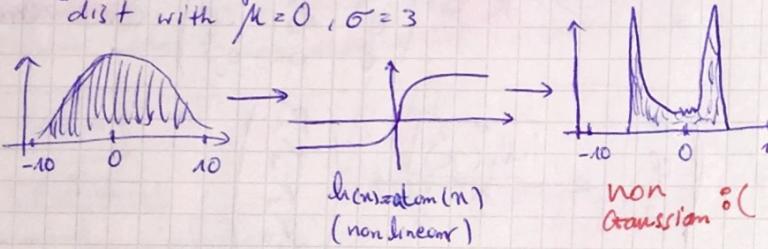
Extended Kalman Filters

- If we map this dist. by a non-linear function, the resulting dist is not a Gaussian. It doesn't also have a closed form.

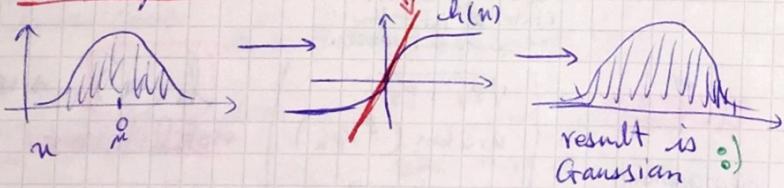


Example:

randomly sample 10 000 numbers from a Gaussian dist with $\mu=0, \sigma=3$



- Replace $h(n)$ with the linear approximation of h at μ .



- How to get the linear approximation? (How to linearize a non-linear function?)

• Taylor Expansion: $f(x) \approx f(\mu) + \frac{\partial f(\mu)}{\partial x} (x-\mu)$

Multivariate Taylor Series:

- Remember the measurement function $h(x')$ for Radar measurements was composed of three equations that showed how to map the state $x' = (P_x, P_y, V_x, V_y)^T$ to measurement space $z = (p, q, \dot{p})^T$

$$h(x') = \begin{pmatrix} \sqrt{P_x'^2 + P_y'^2} \\ \arctan(P_y'/P_x') \\ P_x' V_x + P_y' V_y \end{pmatrix}$$

Multivariate Taylor series at point a

$$T(a) = f(a) + (a-a)^T Df(a) + \frac{1}{2!} (a-a)^T D^2 f(a)$$

$Df(a)$ = Jacobian Matrix (first order derivative)

$D^2 f(a)$ = Hessian matrix (second order derivative)

Linear Approximation of the function

- only keep the Taylor expansion up to the Jacobian term $Df(a)$. We ignore the Hessian term and other higher order terms. Assuming that $(a-a)$ is small, $(a-a)^2$ or the multidimensional equivalent $(a-a)^T(a-a)$ will be even smaller; our extended KF assumes the higher order terms are negligible.

State vector $n = \begin{pmatrix} P_x \\ P_y \\ V_x \\ V_y \end{pmatrix}$

measurement function describes three components

$$z = \begin{pmatrix} p \\ q \\ \dot{p} \end{pmatrix}$$

← normalized
 $(-n, n)$

⇒ Jacobian of $h(n)$ is a 3×4 matrix

$$H_j = \begin{pmatrix} \frac{\partial p}{\partial P_x} & \frac{\partial p}{\partial P_y} & \frac{\partial p}{\partial V_x} & \frac{\partial p}{\partial V_y} \\ \frac{\partial q}{\partial P_x} & \frac{\partial q}{\partial P_y} & \frac{\partial q}{\partial V_x} & \frac{\partial q}{\partial V_y} \\ \frac{\partial \dot{p}}{\partial P_x} & \frac{\partial \dot{p}}{\partial P_y} & \frac{\partial \dot{p}}{\partial V_x} & \frac{\partial \dot{p}}{\partial V_y} \end{pmatrix}$$

$$H_g = \begin{bmatrix} \frac{P_x}{\sqrt{P_x^2 + P_y^2}} & \frac{P_y}{\sqrt{P_x^2 + P_y^2}} & 0 & 0 \\ -\frac{P_y}{P_x^2 + P_y^2} & \frac{P_x}{P_x^2 + P_y^2} & 0 & 0 \\ \frac{P_x(V_y P_y - V_y P_x)}{(P_x^2 + P_y^2)^{3/2}} & \frac{P_x(V_y P_x - V_x P_y)}{(P_x^2 + P_y^2)^{3/2}} & \frac{P_x}{\sqrt{P_x^2 + P_y^2}} & \frac{P_y}{\sqrt{P_x^2 + P_y^2}} \\ \end{bmatrix}$$

- Extended Kalman Filter: state transition
 - First linearize the nonlinear prediction (F) and measurement (H) function
 - Use the normal KF scheme using the linearized F and H (\approx Jacobians)

Kalman Filter

$$\mathbf{x}' = \mathbf{F} \mathbf{x} + \mathbf{u}$$

measurement
update

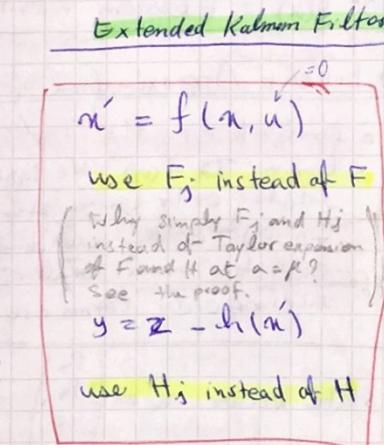
$$\begin{aligned} y &= z - H\alpha' \\ S &= HP'H^T + R \end{aligned}$$

$$K = P' H^T S^{-1}$$

$$\alpha = \alpha' + Ky$$

$$P = (I - KH)P'$$

- The mathematical proof is complex, but the equations for KF and EKF are very similar. The differences are:



- ←

 - the F matrix will be replaced by F_j when calculating P
 - the H matrix will be replaced by H_j when calculating S , K , and P
 - to calculate \hat{x} , the prediction function f , is used instead of F matrix
 - to calculate \hat{y} , the h function is used instead of H matrix.

- Evaluating performance of KF: RMSE metric

$$\sqrt{\frac{1}{n} \sum_{t=1}^n (\hat{x}_t^{\text{est}} - x_t^{\text{gttruth}})^2}$$

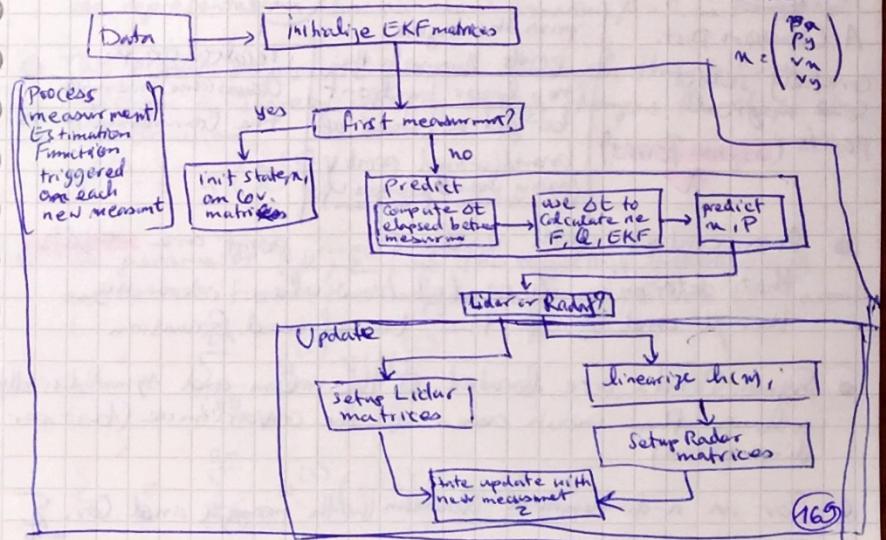
May 14, 21:0

Mym 18.20

Project 5: Extended Kalman Filter

- RViz : ROS Visualization tool.

- ## Sensor Fusion General Flow

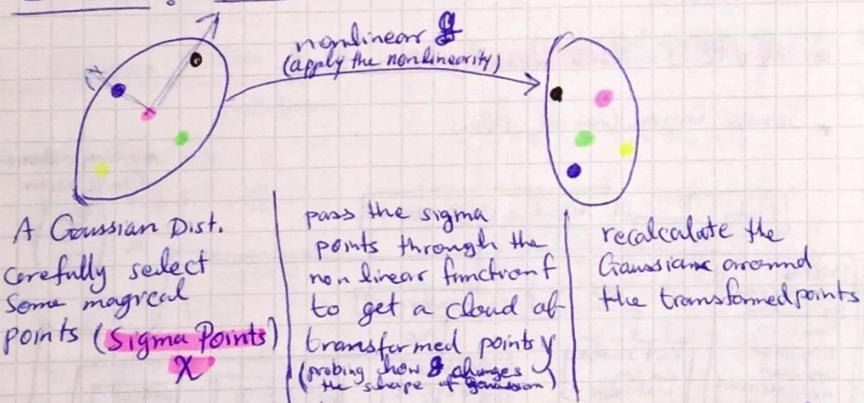


Extra: Unscented KF

May 18, 21:00
May 19, 18:00

- UKF does not use a Taylor series for linearizing $f(n)$ and $h(n)$
- It avoids to linearize around the mean as the EKF does.
- It is based on the assumption/intuition that it is easier to approximate a Gaussian distribution than it is to approximate an arbitrary non-linear function.
- It uses a stochastic linearization through the use of weighted statistical linear regression process.

Idea: Unscented Transform



- Associated with each sigma point are weights that determine its contribution when recovering the μ and Σ of the transformed Gaussian.
- Sigma Points are located at the mean and symmetrically along the main axes of the covariance (two per dimension)
- For an n -dimensional Gaussian with mean μ and Cov. Σ

← the $2n+1$ Sigma points are chosen according to the following rule:

$$X^{[0]} = \mu$$

$$X^{[i]} = \mu + (\sqrt{(n+\lambda)\Sigma})_i \quad \text{for } i=1, \dots, n$$

$$X^{[i+n]} = \mu - (\sqrt{(n+\lambda)\Sigma})_i \quad \text{for } i=n+1, \dots, 2n$$

where $\lambda = \alpha^2(n+\kappa)-n$ with α and κ scaling parameters that determine how far the sigma points are spread from the mean.

- Each sigma point $X^{[i]}$ has two weights associated with it: One weight $w_m^{[i]}$, is used when computing the mean, the other weight, $w_c^{[i]}$, is used when recovering the Cov. of the Gaussian.

$$w_m^{[0]} = \frac{\lambda}{n+\lambda}, \quad w_c^{[0]} = \frac{\lambda}{n+\lambda} + (1-\alpha^2 + \beta) \quad \sum w_i = 1$$

$$w_m^{[i]} = w_c^{[i]} = \frac{1}{2(n+\lambda)} \quad \text{for } i=1, \dots, 2n \quad M = \sum_i w_i \cdot X_i$$

$$\Sigma = \sum_i w_i (X_i - \mu)(X_i - \mu)^T$$

the param β can be chosen to encode additional (higher order) information about distribution underlying the Gaussian representation. For an exact Gaussian $\beta=2$ is the optimal

- The sigma points are passed through the nonlinear function g thereby probing how g changes the shape of the Gaussian

$$Y^{[i]} = g(X^{[i]})$$

- The parameter (μ', Σ') of the resulting Gaussian are extracted from the mapped sigma points $Y^{[i]}$ according to:

$$\mu' = \sum_{i=0}^{2n} w_m^{[i]} Y^{[i]}$$

$$\Sigma' = \sum_{i=0}^{2n} w_c^{[i]} (Y^{[i]} - \mu')(Y^{[i]} - \mu')^T$$

Algorithm Unscented KF ($\mu_{t-1}, \Sigma_{t-1}, u_t, z_t$)

prediction step $\gamma = \sqrt{n+1}$

$$\bar{\chi}_{t-1} = (\mu_{t-1} \quad \mu_{t-1} + \gamma \sqrt{\Sigma_{t-1}} \quad \mu_{t-1} - \gamma \sqrt{\Sigma_{t-1}})$$

$$\bar{x}_t^* = g(u_t, \bar{\chi}_{t-1})$$

$$\bar{\mu}_t = \sum_{i=0}^{2n} w_m^{(i)} \bar{x}_t^{*(i)}$$

$$\bar{\Sigma}_t = \sum_{i=0}^{2n} w_c^{(i)} (\bar{x}_t^{*(i)} - \bar{\mu}_t)(\bar{x}_t^{*(i)} - \bar{\mu}_t)^T + R_t$$

Measurement update Step

$$\bar{x}_t = (\bar{\mu}_t \quad \bar{\mu}_t + 8\sqrt{\bar{\Sigma}_t} \quad \bar{\mu}_t - 8\sqrt{\bar{\Sigma}_t})$$

$$\bar{z}_t = h(\bar{x}_t)$$

$$\hat{z}_t = \sum_{i=0}^{2n} w_m^{(i)} \bar{z}_t^{(i)}$$

$$S_t = \sum_{i=0}^{2n} w_c^{(i)} (\bar{z}_t^{(i)} - \hat{z}_t)(\bar{z}_t^{(i)} - \hat{z}_t)^T + Q_t$$

$$\bar{\Sigma}_{x,z} = \sum_{i=0}^{2n} w_c^{(i)} (\bar{x}_t^{(i)} - \bar{\mu}_t)(\bar{z}_t^{(i)} - \hat{z}_t)^T$$

cross-covariance between state and observation, use to compute the Kalman gain
corresponds to the term $P' H^T$ in EKF.

$$K_t = \bar{\Sigma}_{x,z}^{-1} S_t^{-1}$$

$$\mu_t = \bar{\mu}_t + K(z_t - \hat{z}_t)$$

$$\Sigma_t = \bar{\Sigma}_t - K_t S_t K_t^T$$

return μ_t, Σ_t

• UKF is almost as fast as EKF.

- For non-linear systems UKF is as good or better than EKF with difference being higher with the degree of nonlinearity. In many practical applications though the difference between UKF and EKF are negligible.
- UKF is derivative-free

19. May 18:00

Part 2

Localization | Path Planning
Control | System Integration

19. May 18:00

Localization

Lesson 1: Intro to Localization

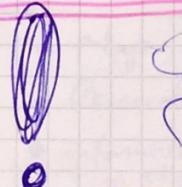
- Tiffany Huang: localization engineer (Mercedes-Benz) (MB RD North America)
- Maximilian Muffert: ...

- Localization: A Robot takes the information about the environment and

compares that to information it already knows about the world.

- Blind folded in the world \Rightarrow can be anywhere, no belief uncertainty as big as the whole world
- Open eye and see Eiffel tower reduce the uncertainty to few kilometers

Intritration

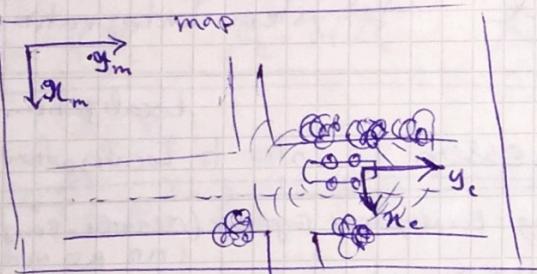


- Intuition: Reduce the uncertainty by measuring/sensing the world and comparing that measurement to the pre-known model of the world and hence building a belief.

- Answer the question:

where is the car in a given map with an accuracy of 10cm or less.

- GPS \rightarrow 3 - 50 meters (not usable)



- Find the transformation between local coordinate system of the car and the global coordinate system of the map by finding position of objects for whom we know their position in the map.

- 1D Bayesian Filter in C++
- Motion models
- 2D Particle Filter in C++

20. May 12:00
22. May 20:00

~~other 3h~~

Lesson 2: Markov Localization

- Localization: Multiplication (Bayes Rule sense/measurement) and Addition (motion (convolution))
- Whole underlying math behind the General localization problem

→ Derivation of the Bayes Localization Filters (Markov Localization) (one of the most common loc. frameworks)

→ 1D Realization of Markov Loc. Filters in C++

→ Motion & Observation Model.

- The project: Kidnapped Vehicle uses Particle Filters.
- Markov Localization is a generalized Filter for localization and all other loc. approaches are realizations of this approach.
- By Understanding the Markov Loc. Filter, we develop intuition and methods that help us solve any vehicle loc. tasks, including Particle Filters.

→ Think of Vehicle position generally as a prob. dist.

→ Each time we move \rightarrow dist. becomes wider
more uncertain

→ Each time we feed the filter with data (map data, measurement, control) \rightarrow the dist. becomes
more certain
more concentrated
more narrow

What do we have?

- Map with Landmark positions in global co.sy.
- Observation from the onboard sensors in the local co.sy.
- information about how the car moves between two time steps.