

Good to know

1. Where and why ReLU outperforms Sigmoid in hidden layers

1. Alleviation of the vanishing-gradient problem

- Sigmoid:

$$\sigma(x) = \frac{1}{1 + e^{-x}}, \quad \sigma'(x) = \sigma(x)(1 - \sigma(x)) \leq 0.25.$$

For large $|x|$, $\sigma'(x) \rightarrow 0$, so gradients vanish as they backpropagate through many layers.

- ReLU:

$$\text{ReLU}(x) = \max(0, x), \quad \text{ReLU}'(x) = \begin{cases} 1, & x > 0, \\ 0, & x \leq 0. \end{cases}$$

When active ($x > 0$), the gradient is constant 1, so it preserves gradient magnitude better in deep nets.

2. Sparse activations and computational efficiency

- ReLU sets all negative pre-activations to zero, producing sparsity in the network. Sparse activations reduce inter-neuron dependencies and can speed up both forward and backward passes.
- The piecewise-linear form of ReLU is cheaper to compute than the exponential in $\sigma(x)$.

3. Empirical performance in deep architectures

- Modern deep convolutional and feed-forward networks (e.g. AlexNet, ResNet) almost universally use ReLU (or its variants), because it leads to faster convergence and often better generalization.

2. When Logistic (Sigmoid) is preferable

1. Output layers for probabilistic interpretation

- For binary classification, the output neuron often needs to produce a probability in $[0, 1]$. The sigmoid's range matches this requirement directly, making it the natural choice for the final layer when used with cross-entropy loss.

2. Gating mechanisms in recurrent networks

- In architectures like LSTMs or GRUs, gates (forget, input, output) must produce values in $[0, 1]$ to modulate information flow. Logistic activations perfectly serve this gating role, whereas ReLU's unbounded positive output and zero gradient for negatives would break the soft gating behavior.

3. Shallow or specialized models

- In some shallow networks or autoencoders where depth is limited, vanishing gradients are less severe. Sigmoids can still work reasonably, especially when you need bounded, smooth outputs or want to exploit their probabilistic interpretation in energy-based models (e.g. Boltzmann machines).

Mathematical summary

Property	Sigmoid $\sigma(x)$	ReLU $\max(0, x)$	
Range	$(0, 1)$	$[0, \infty)$	
Derivative	$\sigma(x)(1 - \sigma(x)) \leq 0.25$	1 for $x > 0$, else 0	
Vanishing gradient	Severe for $ x \gg 1$	x	
Sparsity	No	Yes (zeros for $x \leq 0$)	
Computational cost	Higher (exponential)	Lower (simple comparison)	

By choosing ReLU in deep hidden layers, you maintain stronger gradients, encourage sparsity, and speed up training. Conversely, Sigmoid remains indispensable where bounded outputs or soft gating are required.

Batch Norm (1)

Implementing Batch Norm

Given some intermediate values in NN $\underbrace{z^{(1)}, \dots, z^{(n)}}_{z^{[l]}(:)}$

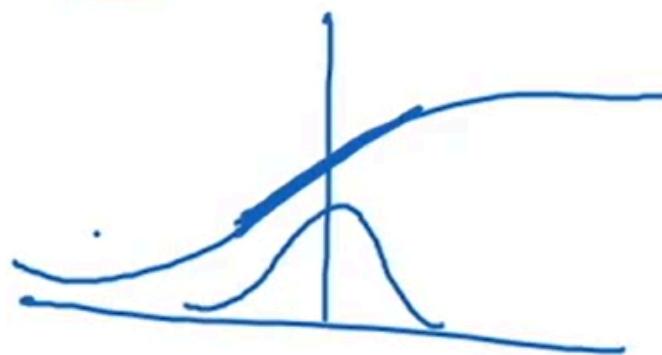
$$\mu = \frac{1}{m} \sum_i z^{(i)}$$
$$\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2$$
$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$
$$\hat{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

learnable parameters of model.

If gamma were $\text{sqrt}(\text{sigmasqr}d + \epsilon)$ and beta were mu, then no effect.

$X \leftarrow$

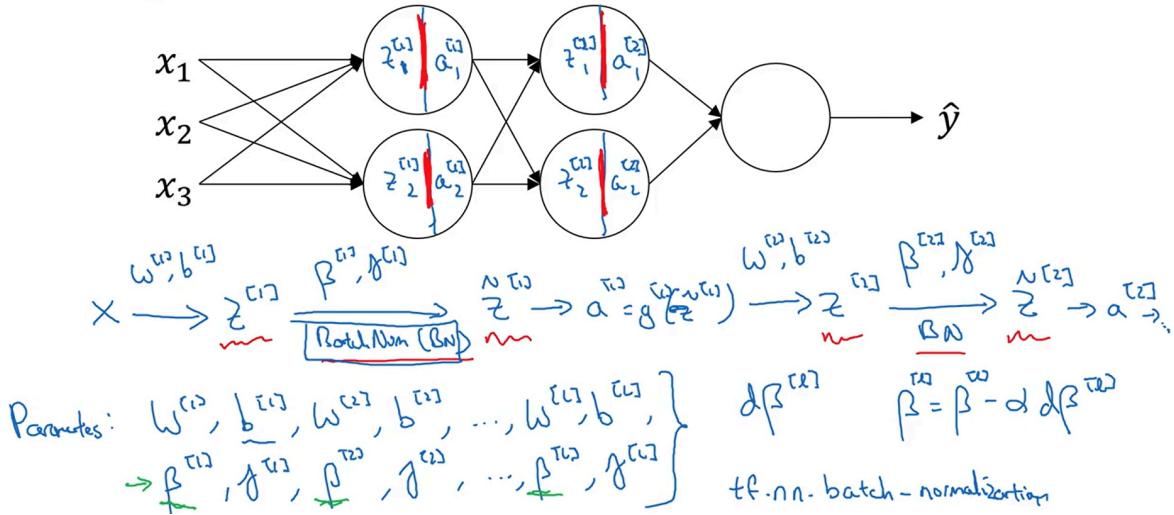
$\underline{z}^{(i)} \leftarrow$



But we need to be careful to avoid the above situation (where the values are clustered around 0) and we are unable to take advantage of the non-linearity of sigmoid. That's where gamma and beta come into play

Adding BatchNorm to a network

Adding Batch Norm to a network



If you use batch norm, just can ignore b

Implementing gradient descent

for $t = 1 \dots \text{numMiniBatches}$

Compute forward pass on X^{t+3} .

In each hidden layer, use BN to replace $\underline{z}^{(t)}$ with $\underline{\tilde{z}}^{(t)}$.

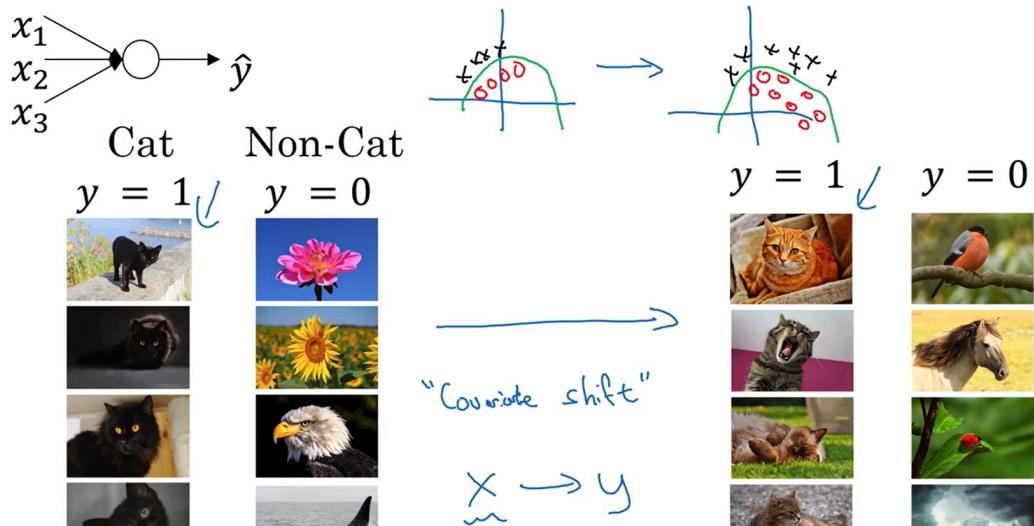
Use backprop to compute $\underline{dw}^{(t)}, \underline{db}^{(t)}, \underline{d\beta}^{(t)}, \underline{d\gamma}^{(t)}$

Update parameters $\left. \begin{array}{l} w^{(t)} := w^{(t)} - \alpha \underline{dw}^{(t)} \\ \beta^{(t)} := \beta^{(t)} - \alpha \underline{d\beta}^{(t)} \\ \gamma^{(t)} := \dots \end{array} \right\}$

Works w/ momentum, RMSprop, Adam.

Why does BatchNorm work?

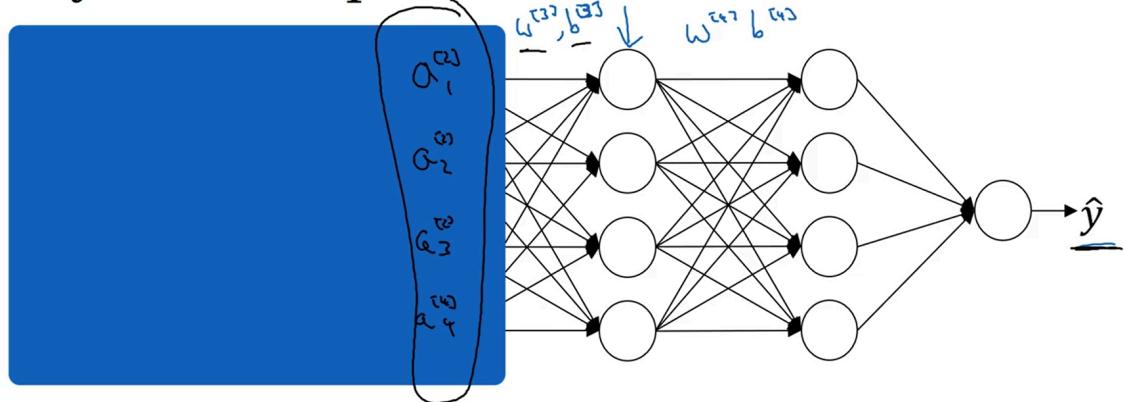
Learning on shifting input distribution



If you train to identify black cats and you test on coloured cats (:D), then your classifier may not work properly mah!

The change of the distribution from X (containing black cats) to X' (containing coloured cats is known as covariate shift)

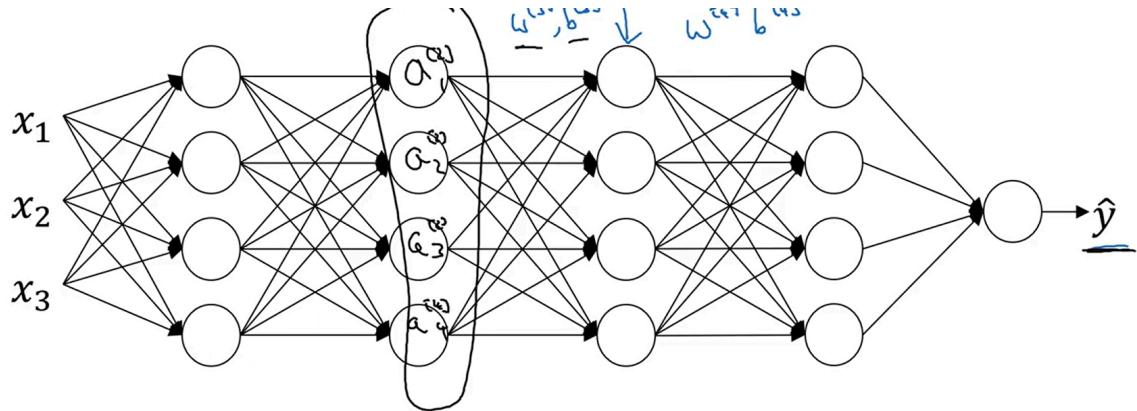
Why this is a problem with neural networks?



Ok bah?

The NN needs to try to work with the activations (can be interpreted as features). Job of the third layer needs to take these values and try to map to y

Let's uncover the blue mah?



From the view of the third layer, the inputs to it always keep changing! Come on lor! This is like that covariate shift bah!

So what batch norm does is that it reduces the amount that the distribution of hidden unit values shift around.

This is done by ensuring that the mean and variance doesn't change lor! Mean depends on gamma and beta for that BN layer

It limits the amount of how changing parameters in an earlier layer will affect the inputs to the current layer

i.e weakens coupling between weights parameters of learning almost making learning independent b/w layers.

Ok bah, got it? Sike! Here is another usage of BN! (You have more to think about lol)

Batch Norm as regularization

X

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch. $\hat{z}^{[l]}$ $\langle \cdot, \cdot \rangle$ $z^{[l]}$
- This adds some noise to the values $z^{[l]}$ within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations.
- This has a slight regularization effect.

Noise is added because the normalization is done w.r.t to the current minibatch and not w.r.t to the entire training set. (i.e estimate of the mean and st.dev is noisy lor!) This noise forces downstream hidden units not to be dependent on one hidden unit!

(So the larger mini-batch size you use, the less regularization you do! Ok mah? You seem confused lor!)

BN → Multiplicative noise + Additive noise

Dropout → Additive noise

Adam Optimizer

Adam = RMSProp + Momentum (Adaptive Moment Estimation)

- Start by initializing V_{dw} , S_{dw} and V_{db} , S_{db} to all zero
- Then do the Momentum calculations for V_{dw} , V_{db}
- Then do the RMS calculations for S_{dw} , S_{db}
- Then do bias corrections for all
- Then just combine Momentum update with RMS update 😊

$$V_{dw} = 0, S_{dw} = 0. \quad V_{db} = 0, S_{db} = 0$$

On iteration t :

Compute $\delta w, \delta b$ using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) \delta w, \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) \delta b \quad \leftarrow \text{"moment"} \beta_1$$
$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) \delta w^2, \quad S_{db} = \beta_2 S_{db} + (1-\beta_2) \delta b^2 \quad \leftarrow \text{"RMSprop"} \beta_2$$
$$V_{dw}^{corrected} = V_{dw} / (1 - \beta_1^t), \quad V_{db}^{corrected} = V_{db} / (1 - \beta_1^t)$$
$$S_{dw}^{corrected} = S_{dw} / (1 - \beta_2^t), \quad S_{db}^{corrected} = S_{db} / (1 - \beta_2^t)$$
$$w := w - \alpha \frac{V_{dw}^{corrected}}{\sqrt{S_{dw}^{corrected} + \epsilon}}$$
$$b := b - \alpha \frac{V_{db}^{corrected}}{\sqrt{S_{db}^{corrected} + \epsilon}}$$

Hyperparameters choice:

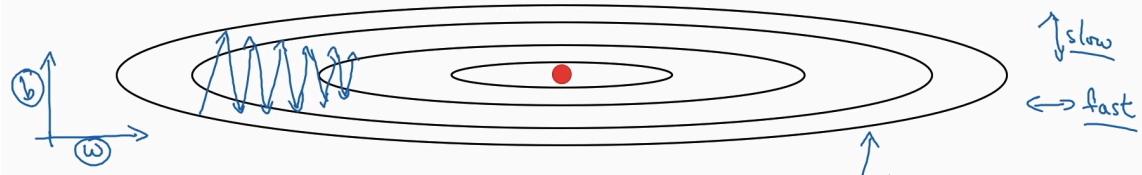
- α : needs to be tune
- β_1 : 0.9 (dw)
- β_2 : 0.999 ($d\omega^2$)
- ϵ : 10^{-8}

Only need to tune alpha usually. b1 and b2 were the default suggestions by the authors of the Adam optimizer paper

RMSProp

RMSPop → Root Mean Square Propagation

RMSprop



Want to:

Slow down learning in the b direction and increase learning in the W direction.

On iteration t :

Compute dW, db on current mini-batch
element-wise

$$S_{dW} = \beta S_{dW} + (1-\beta) \frac{dW^2}{\text{element-wise}}$$

$$S_{db} = \beta S_{db} + (1-\beta) \frac{db^2}{\text{element-wise}}$$

$$W := W - \alpha \frac{dW}{\sqrt{S_{dW}}} \quad b := b - \alpha \frac{db}{\sqrt{S_{db}}}$$

Do note that the squaring is done element wise!

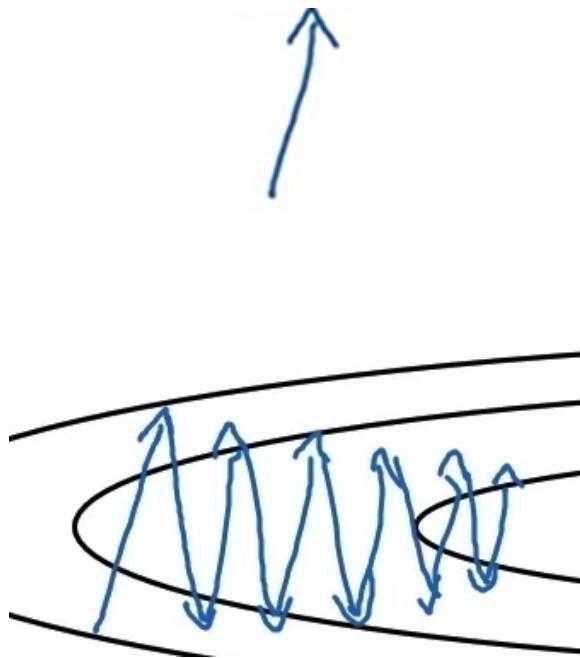
So what this is doing is keeping an EWA of the SQUARES of the derivatives.

Then look carefully at the update rule! We don't directly update by the EWA!

So, goal is to increase learning in the W direction and decrease/dampen learning in the b direction.

So what we hope is that S_{dW} is small so that when you divide by a smaller number, then the update to W is large. And we hope S_{db} is large so that when you divide by a larger number, then the update to b is small in order to slow down updates in the vertical direction.

And indeed, the dW 's are much smaller than the db 's!



Here, db is much larger compared to dW

Therefore oscillations in the vertical direction get damped out!

Think of this as adaptive learning rates!

Gradient Descent w/ Momentum

Almost always works faster than gradient descent.

Basic idea: compute an EWA of gradients and use that gradient to update weights instead.'

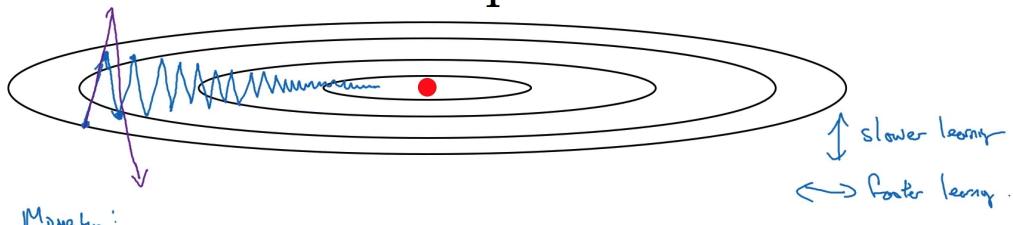
i.e smooth out gradient descent

To prevent oscillations, we need the vertical learning to be slower and we need the horizontal learning to be faster.

To do this, for every iteration

- Just compute the gradients dW and db on current minibatch
- Simply get the EWA $V_{dW} = \beta V_{dW} + (1-\beta) dW$
- And same for the EWA V_{db}
- Then do the updates $W = W - lr * V_{dW}$ and $b = b - alpha * V_{db}$

Gradient descent example



Momentum:

On iteration t :

Compute dW, db on current mini-batch.

$$V_{dW} = \beta V_{dW} + (1-\beta) dW \quad "V_0 = \beta V_{0\tau} + (1-\beta) \theta_t"$$
$$V_{db} = \beta V_{db} + (1-\beta) db$$

$$W = W - \alpha V_{dW}, \quad b = b - \alpha V_{db}$$

Eg: averaging over these gradients will cause vertical component to be almost zero and horizontal to be much more.



Now two hyperparameters: alpha (lr) and beta (most common in 0.9, average over last 10)

Implementation details

On iteration t :

Compute dW, db on the current mini-batch

$$v_{dW} = \beta v_{dW} + (1 - \beta) dW$$

$$v_{db} = \beta v_{db} + (1 - \beta) db$$

$$W = W - \alpha v_{dW}, \quad b = b - \alpha v_{db}$$

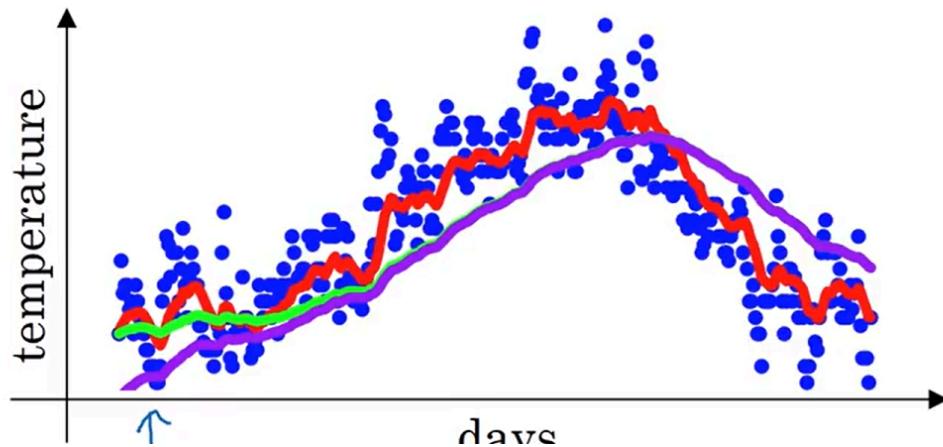
Hyperparameters: α, β $\beta = 0.9$

Bias correction in exponentially weighted averages

$\beta = 0.9 \rightarrow$ red line

$\beta = 0.98 \rightarrow$ green line (actually you don't get the green line, you get the purple line)

The purple line starts off low.

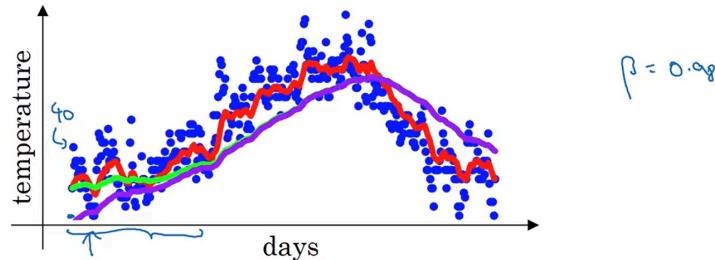


Initial EWA is bad, therefore use the following formula:

(essentially do the same EWA as before, but just divide each of those v_t 's by $(1-\beta^t)$)

And also note that $1-\beta^t$ is the sum of the coefficients of the summands, so it is a weighted average—it's at the scale of new temperatures.

Bias correction



$$\Rightarrow v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

$$v_0 = 0$$

$$v_1 = \underline{0.98 v_0 + 0.02 \theta_1}$$

$$v_2 = 0.98 v_1 + 0.02 \theta_2$$

$$= 0.98 \times 0.02 \times \theta_1 + 0.02 \theta_2$$

$$= 0.0196 \underline{\theta_1} + 0.02 \underline{\theta_2}$$

$$\left| \begin{array}{l} \frac{v_t}{1-\beta^t} \\ t=2: 1-\beta^t = 1-(0.98)^2 = 0.0396 \\ \frac{v_2}{0.0396} = \frac{0.0196 \underline{\theta_1} + 0.02 \underline{\theta_2}}{0.0396} \end{array} \right.$$

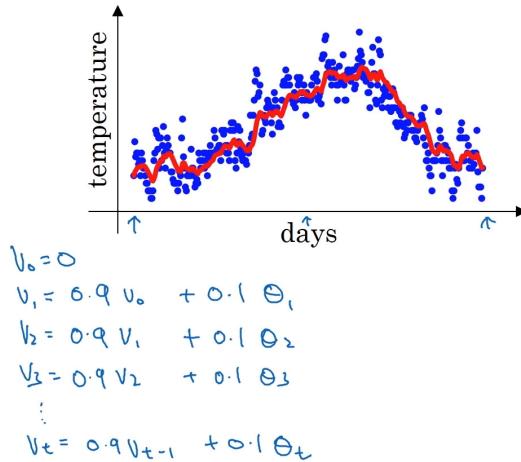
Andrew Ng

Now, as t becomes large $1-\beta^t$ will be 1 meaning bias correction is switched off.

Exponentially weighted averages

Temperature in London

$$\begin{aligned}\theta_1 &= 40^{\circ}\text{F} \quad 4^{\circ}\text{C} \\ \theta_2 &= 49^{\circ}\text{F} \quad 9^{\circ}\text{C} \\ \theta_3 &= 45^{\circ}\text{F} \quad \vdots \\ &\vdots \\ \theta_{180} &= 60^{\circ}\text{F} \quad 15^{\circ}\text{C} \\ \theta_{181} &= 56^{\circ}\text{F} \quad \vdots \\ &\vdots\end{aligned}$$



Andrew Ng

The red line is the exponentially weighted average

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

Here you can think of V_t as averaging over the past $1/(1-\beta)$ days' temperature.

V_t is approximately
averaging over
 $\approx \frac{1}{1-\beta}$ days'
temperature.

High beta \rightarrow (eg: green line), the plot is smoother because it keeps track of a longer range of days' temperatures. Adapts slowly to newer temperatures

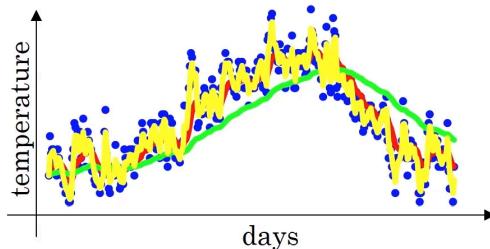
Very low beta \rightarrow (eg: red line), very noisy yellow line since you are very susceptible to new temperatures.

Exponentially weighted averages

$$\underline{V_t} = \underline{\beta} \underline{V_{t-1}} + \underline{(1-\beta) Q_t} \rightarrow$$

$\underline{\beta} = 0.9$: ≈ 10 days' temporal
 $\underline{\beta} = 0.98$: ≈ 50 days
 $\underline{\beta} = 0.5$: ≈ 2 days

Ve as approximately
average over
 $\rightarrow x \frac{1}{1-\beta}$ days
 temperature.



$$\frac{1}{1-0.98} = 50$$

Andrew Ng

Intuition for EWAs

Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

$$v_{100} = 0.9v_{99} + 0.1\theta_{100}$$

$$v_{99} = 0.9v_{98} + 0.1\theta_{99}$$

$$v_{98} = 0.9v_{97} + 0.1\theta_{98}$$

3

$$V_{100} = 0.1 \Theta_{100} + 0.9 \cancel{(0.1 \Theta_{99} + 0.9 \cancel{\Theta_{98}})} \quad \begin{matrix} \downarrow \\ 0.1 \Theta_{98} + 0.9 \Theta_{97} \end{matrix}$$

$$= 0.1 \Theta_{100} + 0.1 \times 6.9 \cdot \Theta_{99} + 0.1 (0.9)^2 \Theta_{98} + 0.1 (0.9)^3 \Theta_{97} + 0.1 (0.9)^4 \Theta_{96} + \dots$$

$$\underline{0.9} \approx 0.35 \approx \frac{1}{e} \quad \frac{(1-\varepsilon)^{\frac{1}{\varepsilon}}}{0.9} = \frac{1}{e} \quad 0.98 ?$$

$$\varepsilon = 0.02 \rightarrow 0.98^{\frac{1}{0.02}} \approx \frac{1}{e} \quad \text{Andrew Ng}$$

Explanation of the graphs on the top right

Top graph: just the temperature measurements over time.

Bottom one: exponentially decaying function of weightage factor for each of the elements on top.

So v_{100} = element wise product of those two graphs and sum them up

Kalman Filters

For the linear case, we have two models.

The linear dynamics model (how the state evolves over time)

Linear Dynamics Model

Dynamics model: State undergoes linear transformation plus Gaussian noise

$$\underline{\mathbf{x}}_t \sim N(D_t \mathbf{x}_{t-1}, \Sigma_{d_t})$$

The linear measurement mode

Linear Measurement Model

Observation model: Measurement is linearly transformed state plus Gaussian noise

$$\mathbf{y}_t \sim N(M_t \mathbf{x}_t, \Sigma_{m_t})$$

Constant Velocity (1D example)

Example: Constant velocity (1D)

State vector is position and velocity

$$x_t = \begin{bmatrix} p_t \\ v_t \end{bmatrix} \quad p_t = p_{t-1} + (\Delta t)v_{t-1} + \varepsilon$$
$$v_t = v_{t-1} + \xi$$

$$x_t = D_t x_{t-1} + \text{noise} = \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix} \begin{bmatrix} p_{t-1} \\ v_{t-1} \end{bmatrix} + \text{noise}$$

Example: Constant velocity (1D)

Measurement is position only

$$y_t = Mx_t + \text{noise} = [1 \quad 0] \begin{bmatrix} p_t \\ v_t \end{bmatrix} + \text{noise}$$

Constant Acceleration (1D example)

Example: Constant acceleration (1D)

State vector is position, velocity & acceleration

$$\begin{aligned} x_t = \begin{bmatrix} p_t \\ v_t \\ a_t \end{bmatrix} &\quad p_t = p_{t-1} + (\Delta t)v_{t-1} + \varepsilon \\ &\quad v_t = v_{t-1} + (\Delta t)a_{t-1} + \xi \\ x_t = D_t x_{t-1} + \text{noise} &= \begin{bmatrix} 1 & \Delta t & 0 \\ 0 & 1 & \Delta t \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} p_{t-1} \\ v_{t-1} \\ a_{t-1} \end{bmatrix} + \text{noise} \end{aligned}$$

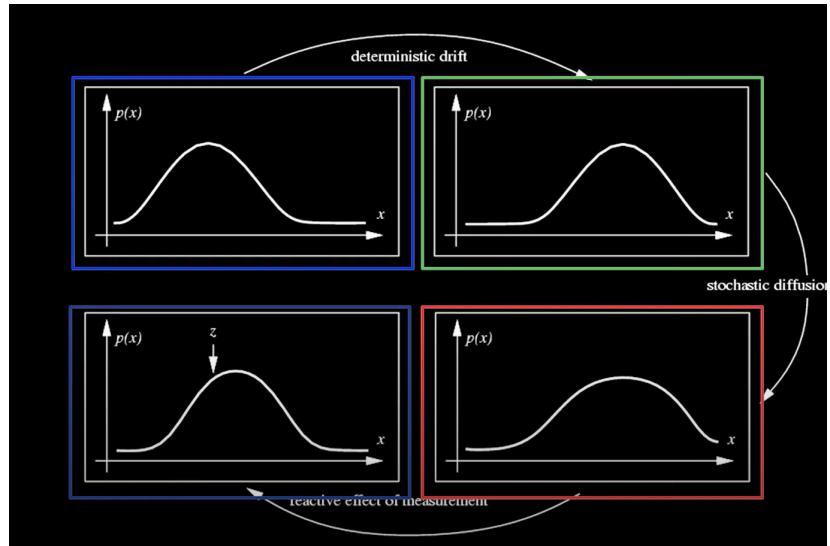
Example: Constant acceleration (1D)

Measurement is position only

$$y_t = Mx_t + \text{noise} = [1 \quad 0 \quad 0] \begin{bmatrix} p_t \\ v_t \\ a_t \end{bmatrix} + \text{noise}$$

Kalman Filter

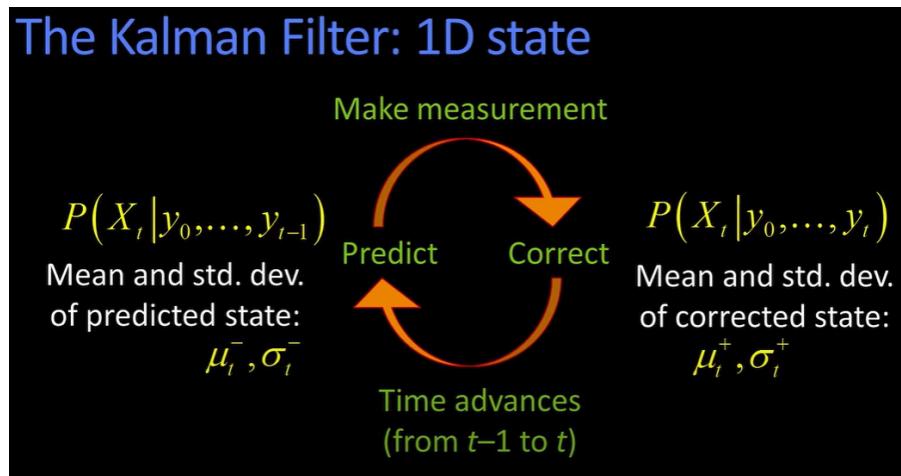
Method to track linear dynamical systems with GAUSSIAN NOISE. Both predicted and corrected states are both gaussian.



So you shift (scale mean) + expand (add noise) + measure + correct + shrink (because new information will only reduce uncertainty)

The first 3 images correspond to prediction.

Kalman Filter 1D State



Note that prediction uses mu minus and sigma minus (mean and covariance BEFORE I take the measurement, this is based on mu plus and sigma plus at t - 1)

Then after you make measurement and correct mean to mu plus and sigma plus at t

1D Kalman Filter Prediction

In this linear dynamics model the next state is gotten by multiplying the previous state by a factor of d and adding some noise modeled by a gaussian.

1D Kalman Filter: Prediction

Linear dynamics model defines predicted state evolution, with noise

$$X_t \sim N(dx_{t-1}, \sigma_d^2)$$

So, if that's the case, the distribution of the next predicted state is also gaussian:

The distribution for next predicted state is also a Gaussian

$$P(X_t | y_0, \dots, y_{t-1}) = N(\mu_t^-, (\sigma_t^-)^2)$$

Crucially note that they are mu and sigma minus because this prediction was made before measurement

And remember that we need to update mu and sigma

Update the mean:

$$\mu_t^- = d\mu_{t-1}^+$$

Update the variance:

$$(\sigma_t^-)^2 = \sigma_d^2 + (d\sigma_{t-1}^+)^2$$

1D Kalman Correction

In this linear dynamics model, the measurement is gotten by multiplying the previous state by a factor of m and adding some gaussian noise

1D Kalman Filter: Correction

Mapping of state to measurements:

$$Y_t \sim N(mx_t, \sigma_m^2)$$

Now, given our predicted state (that we got above), we need to correct it using the measurement!

Predicted state: $P(X_t | y_0, \dots, y_{t-1}) = N(\mu_t^-, (\sigma_t^-)^2)$

Want to estimate corrected distribution:

$$P(X_t | y_0, \dots, y_t) = \frac{P(y_t | X_t) P(X_t | y_0, \dots, y_{t-1})}{\int P(y_t | X_t) P(X_t | y_0, \dots, y_{t-1}) dX_t}$$

Now, the corrected distribution is defined to be a new Gaussian with new mean mu plus and sigma plus.

Kalman: With linear, Gaussian dynamics and measurements, the corrected distribution to be:

$$P(X_t | y_0, \dots, y_t) \equiv N(\mu_t^+, (\sigma_t^+)^2)$$

Formula for mu plus and sigma plus

Update the mean:

$$\mu_t^+ = \frac{\mu_t^- \sigma_m^2 + m y_t (\sigma_t^-)^2}{\sigma_m^2 + m^2 (\sigma_t^-)^2}$$

Update the variance:

$$\underline{(\sigma_t^+)^2} = \frac{\sigma_m^2 (\sigma_t^-)^2}{\sigma_m^2 + m^2 (\sigma_t^-)^2}$$

1D Kalman Intuition

1D Kalman Filter: Intuition

From:

$$\mu_t^+ = \frac{\mu_t^- \sigma_m^2 + m y_t (\sigma_t^-)^2}{\sigma_m^2 + m^2 (\sigma_t^-)^2}$$

What is this?

- *The weighted average of prediction and measurement based on variances!*

Weights Based on Variances: The crucial insight is how the weights are determined. They are inversely proportional to the uncertainty (variance) of each estimate:

- The weight for the prediction (μ_t^-) is proportional to the variance of the measurement (σ_m^2).
- The weight for the measurement-derived estimate (y_t/m) is proportional to the variance of the prediction ($(\sigma_t^-)^2$).

This means:

- If the measurement is very certain (small σ_m^2), the weight for the prediction is smaller, and the weight for the measurement is larger. The filter trusts the measurement more.
- If the prediction is very certain (small $(\sigma_t^-)^2$), the weight for the measurement is smaller, and the weight for the prediction is larger. The filter trusts the prediction more.

Special Cases

Prediction vs. correction

$$\mu_t^+ = \frac{\mu_t^- \sigma_m^2 + m y_t (\sigma_t^-)^2}{\sigma_m^2 + m^2 (\sigma_t^-)^2} \quad (\sigma_t^+)^2 = \frac{\sigma_m^2 (\sigma_t^-)^2}{\sigma_m^2 + m^2 (\sigma_t^-)^2}$$

What if there is no prediction uncertainty? $(\sigma_t^- = 0)$

$$\mu_t^+ = \mu_t^- \quad (\sigma_t^+)^2 = 0$$

The measurement is ignored!

What if there is no measurement uncertainty? $(\sigma_m = 0)$

$$\mu_t^+ = \frac{y_t}{m} \quad (\sigma_t^+)^2 = 0$$

Simplification of Mu+ and Sigma+

1D Kalman Filter: Intuition

$$\text{Also: } \mu_t^+ = \frac{\frac{\mu_t^- \sigma_m^2}{m^2} + \frac{y_t}{m} (\sigma_t^-)^2}{\frac{\sigma_m^2}{m^2} + (\sigma_t^-)^2}$$

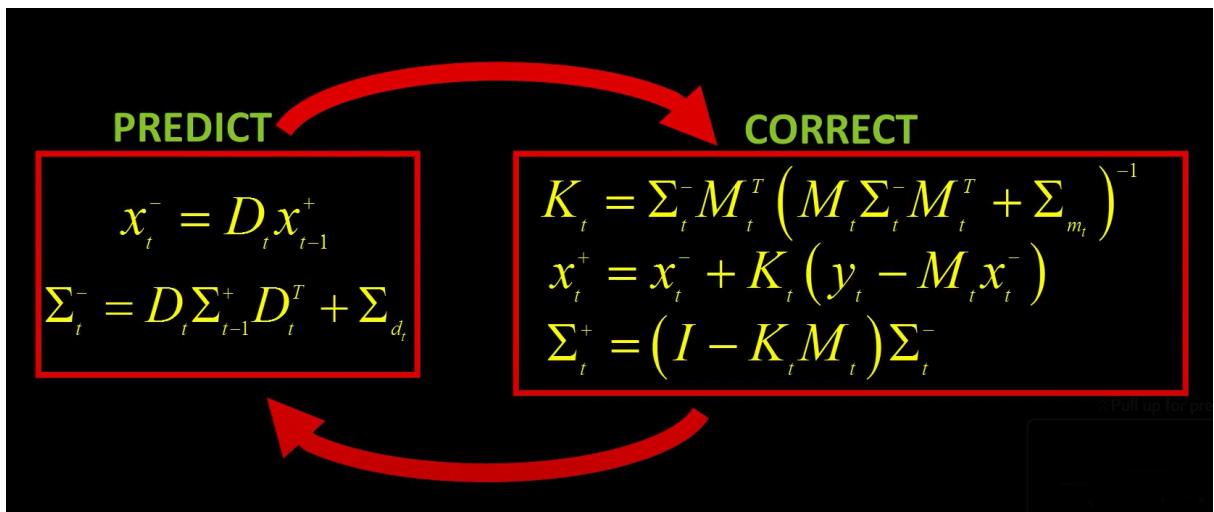
$$\mu_t^+ = \frac{a \mu_t^- + b \frac{y_t}{m}}{a + b} = \frac{(a + b) \mu_t^- + b \left(\frac{y_t}{m} - \mu_t^- \right)}{a + b}$$

$$\mu_t^+ = \frac{a\mu_t^- + b\frac{y_t}{m}}{a+b} = \frac{(a+b)\mu_t^- + b(\frac{y_t}{m} - \mu_t^-)}{a+b}$$

$$\mu_t^+ = \mu_t^- + \frac{b(\frac{y_t}{m} - \mu_t^-)}{a+b} = \mu_t^- + k(\frac{y_t}{m} - \mu_t^-)$$

μ_t^- is the prediction. The difference (red) is the residual. The residual is the difference between y_t/m (the measurements guess of x) and μ_t^- . Something like TD updates.

N-Dimensional Kalman Filters (the real deal)



Σ_t^- : This is the **predicted state covariance matrix** at time t . It represents the uncertainty in the predicted state x_t^- . The covariance matrix describes the variances of each state variable (diagonal elements) and the covariances between them (off-diagonal elements).

CORRECT

$$K_t = \Sigma_t^- M_t^T \left(M_t \Sigma_t^- M_t^T + \Sigma_{m_t} \right)^{-1}$$

$$x_t^+ = x_t^- + K_t \boxed{(y_t - M_t x_t^-)}$$

$$\Sigma_t^+ = (I - K_t M_t) \Sigma_t^-$$

Less weight on residual as a priori estimate error covariance approaches zero.

K_t is Kalman Gain The green box is the residual

Prof Amir's slide

ALGORITHM EXPLAINED

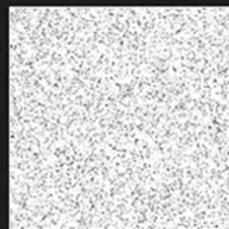
- The prediction step involves two main equations:
 - State Prediction: $\hat{x}_{k|k-1} = F \hat{x}_{k-1} + w_{k-1}$, where $\hat{x}_{k|k-1}$ is the predicted state at time k given the state at k-1, and w_k is the process noise, or a random vector drawn from a multivariate normal distribution with mean zero and covariance matrix Q, $w_k \sim \mathcal{N}(0, Q)$
 - The error covariance prediction: $P_{k+1} = FP_k F^T + Q$, where Q is the process noise covariance matrix.
- $Q = \begin{bmatrix} q_x & 0 & q_{xu} & 0 \\ 0 & q_y & 0 & q_{yu} \\ q_{xu} & 0 & q_u & 0 \\ 0 & q_{yu} & 0 & q_v \end{bmatrix}$
- Where q_x, q_y, q_u, q_v are uncertainty in positions and velocities, and q_{xu}, q_{yu} are correlation between position's and velocity's uncertainty.

M P A ... AMIR 2024 23

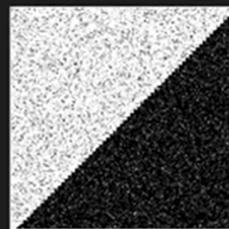
Harris Corner Detector

Corners

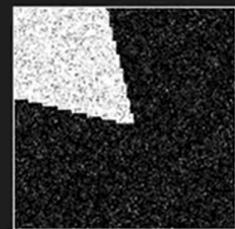
Corner: Point where Two Edges Meet. i.e., Rapid Changes of Image Intensity in **Two Directions** within a Small Region



"Flat" Region



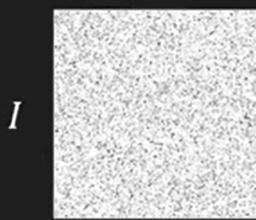
"Edge" Region



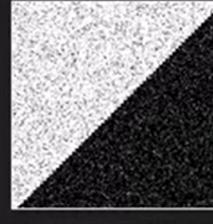
"Corner" Region

Image Gradients

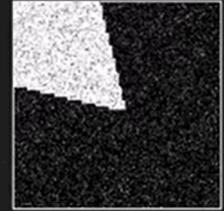
Flat Region



Edge Region

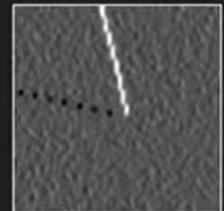
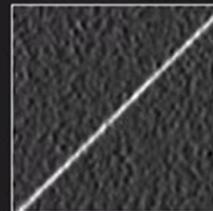
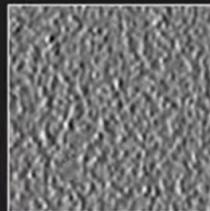


Corner Region

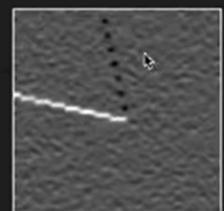
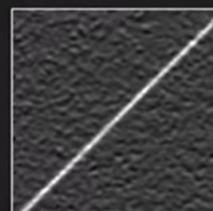
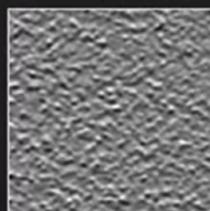


I

$$I_x = \frac{\partial I}{\partial x}$$



$$I_y = \frac{\partial I}{\partial y}$$



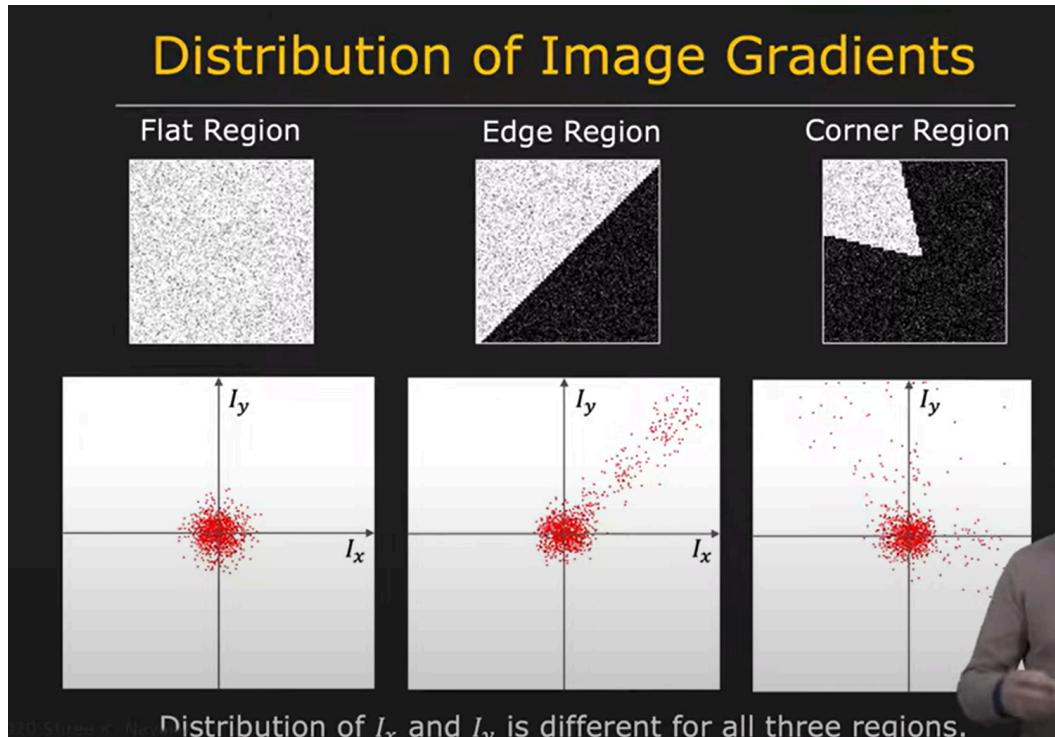
Bree K. Nayar

How come edge map can contain negative values?

If gradient of edge matches the gradient of the filter, then you will get some form of resonance and get a large positive number.

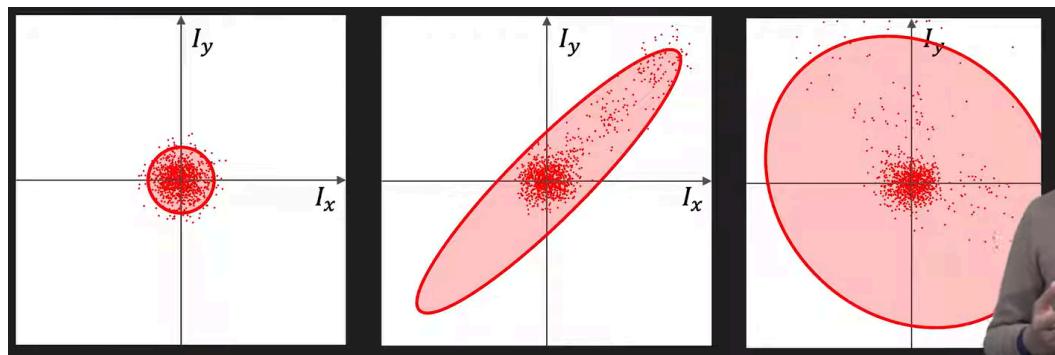
If the gradient of the edge is opposite to the gradient of the filter (for eg: the gradient of the edge is bright to dark and your filter is dark to bright), you will see on the edge map (in that direction, G_x , G_y or whatever), you will see large negative number at the position of that edge.

Images are normalized such that white means strong positive value, black means strong negative value and gray means near 0.

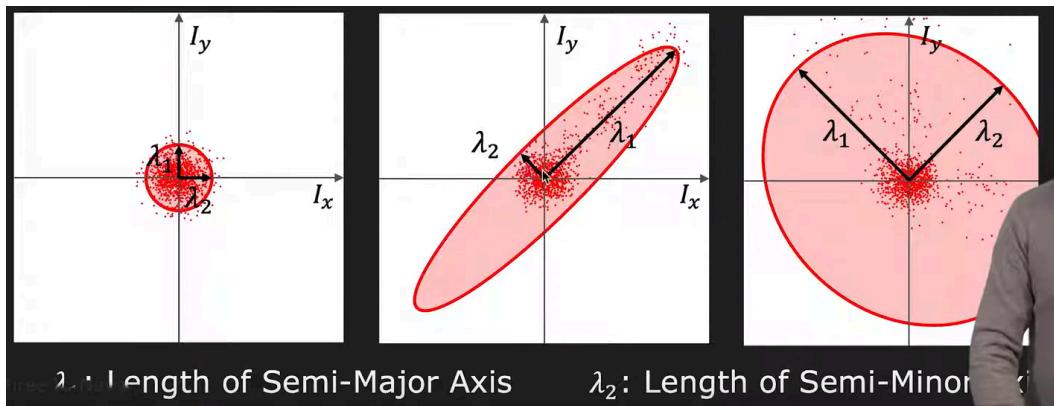


So the goal is to somehow quantify this structure and describe the structure of the distribution with a small number of parameters to classify the region.

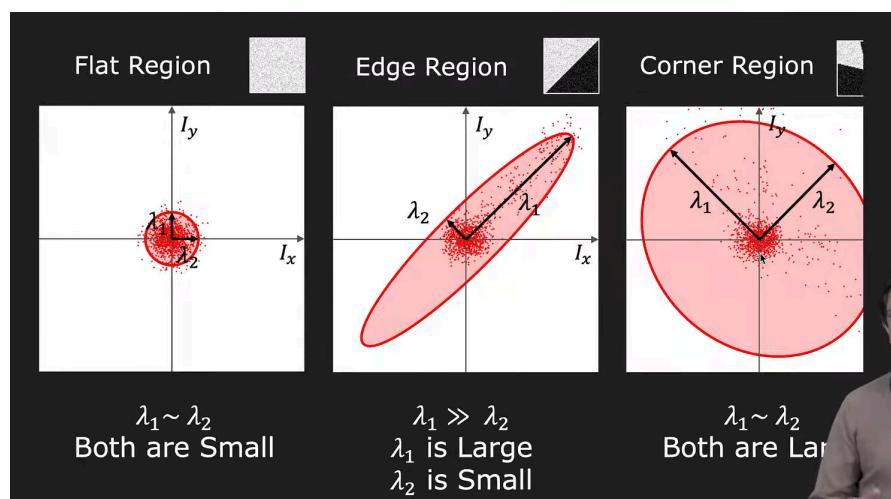
Do this by taking the distribution and fitting an ellipse around it.



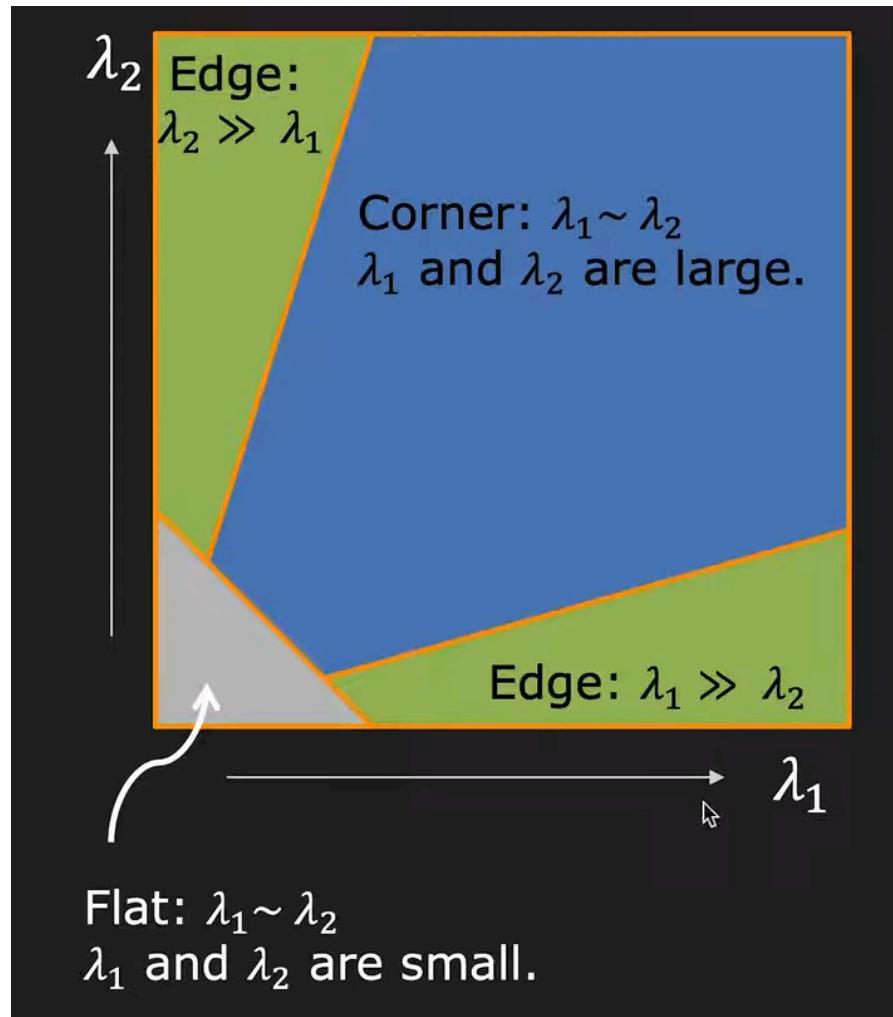
In doing so, you get a semi major (the longer) and a semi minor (the shorter) axis for the ellipse.



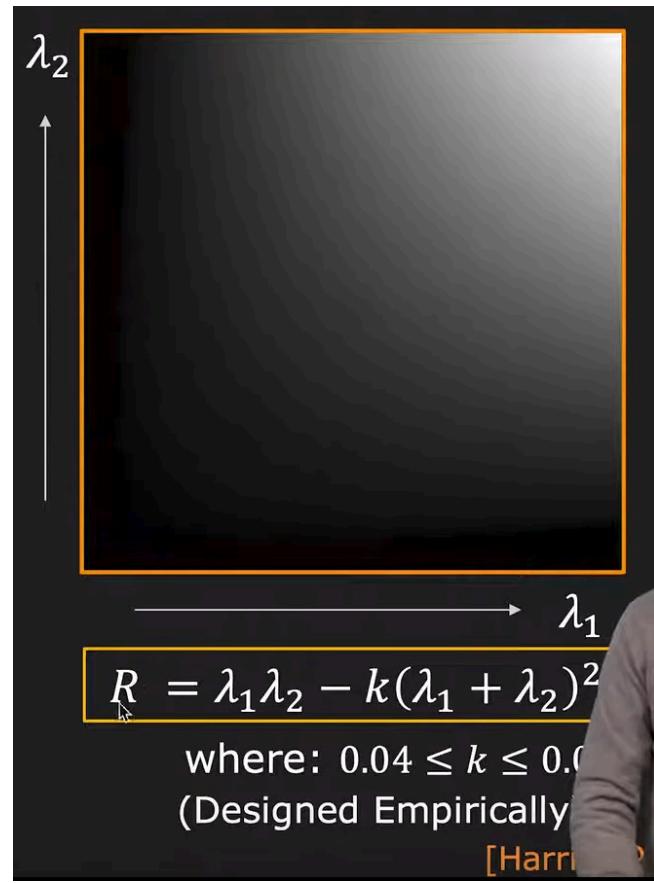
λ_1 is the length of the semi-major axis and λ_2 is the length of the semi-minor axis.



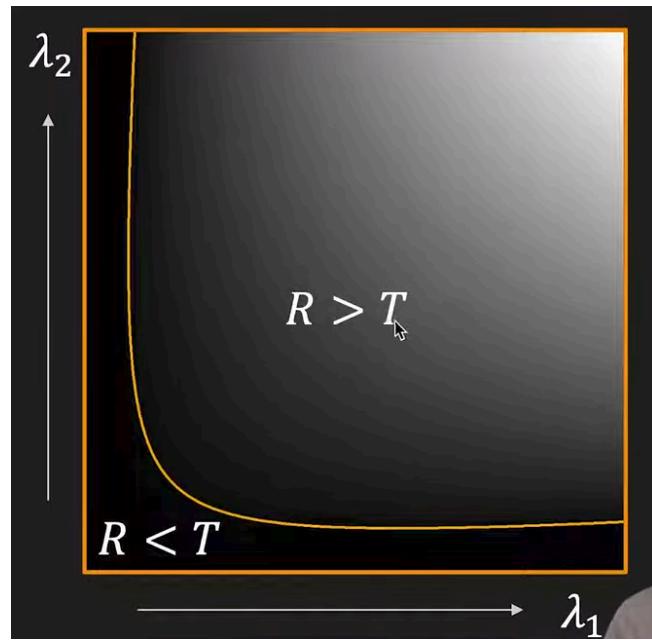
So in λ_1 , λ_2 space, we classify accordingly:



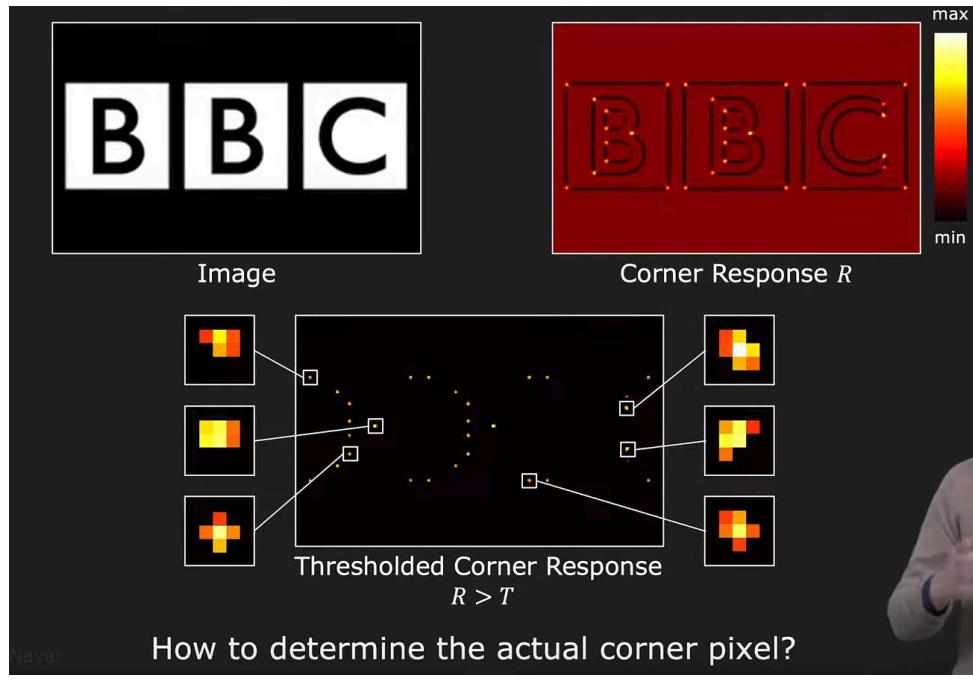
We want to do this simply, Harris came up with this expression where he maps the eigen values to a single real number. This is an empirically determined function ($0.04 \leq k \leq 0.06$ is also empirical). This is the Harris Corner response function



R is being plotted in this image where the brightness of the point is proportional to R. So if you threshold this space wrt R:



Eg:



The thing is, after you threshold, you get clusters of pixels near the corner signaled as the corner. But we want to single out the pixel that is the corner (this would be the peak), so we apply NMS

Non-Maximal Suppression

1. Slide a window of size k over the image.
2. At each position, if the pixel at the center is the maximum value within the window, label it as positive (retain it). Else label it as negative (suppress it).



Used for finding Local Extrema (Maxima/Minima)

HARRIS CORNER DETECTOR

Why Harris Corner Detector?

- **Rotation Invariance:** The Harris Corner Detector is invariant to image rotation, meaning the same corners can be detected even if the image is rotated.
- **Efficiency:** It's a computationally efficient way to detect corners, making it suitable for real-time applications.
- **Robust to Noise:** The algorithm can handle image noise reasonably well.

Limitations:

- **Scale Invariance:** The Harris detector is not scale-invariant, meaning it may not perform well if the image is scaled (i.e., zoomed in or out).
- **Fixed Parameters:** The choice of parameters such as the threshold and k is crucial and may need to be manually tuned for different images.

The biggest limitation is that it's not scale invariant.

Specifically:

- **Changing Gradient Magnitudes:** Scaling an image affects the magnitude of the intensity gradients. If an image is scaled down, the pixel differences across a certain distance decrease, leading to smaller gradient values. Since the Harris detector's response is based on the squares and products of these gradients, the corner response will be lower for a scaled-down image, potentially causing genuine corners to fall below the detection threshold.
- **Fixed Window Size:** The local window used to compute the structure tensor has a fixed size in terms of pixels. When the image scale changes, this fixed window corresponds to a different physical area in the original scene. A corner feature that perfectly fits within the window at one scale might occupy a much larger or smaller area at a different scale. If the corner feature becomes larger than the window, only a part of it might be analyzed, potentially leading to it being misclassified as an edge or a flat region. If the feature becomes too small, its influence on the gradients within the window might be diminished.

Lucas Kanade Point Tracking

Gemini 2.5 Flash

Okay, building on your knowledge of the Lucas-Kanade optical flow algorithm, understanding Lucas-Kanade point tracking is straightforward.

Essentially, **Lucas-Kanade point tracking is an application of the Lucas-Kanade optical flow algorithm specifically tailored to find the motion of a small number of predefined or detected feature points from one image frame to the next.**

Here's how to think about it:

1. **Lucas-Kanade Optical Flow (General):** This algorithm estimates the 2D motion vector (optical flow, (u,v)) for pixels within a local window, based on the assumptions of brightness constancy and spatial coherence. It solves a system of linear equations (using least squares) derived from the optical flow constraint equation $I_t u_x + I_t v_y + I_t = 0$ applied to all pixels in the window.
2. **Lucas-Kanade Point Tracking (Specific):**
 - **Goal:** Don't calculate the flow for every pixel or large regions. Instead, pick a few important points in the first frame (these are often "good features to track," like corners, detected using algorithms like Shi-Tomasi or Harris).
 - **Method:** For each point you want to track:
 - Define a small window (a patch of pixels) centered at *that specific point* in the first frame (I_t).
 - Apply the standard Lucas-Kanade optical flow calculation *just within that small window*, comparing the patch in frame I_t to the corresponding patch in the next frame (I_{t+1}).
 - The resulting optical flow vector (u,v) calculated for that specific window is the estimated movement of the feature point from frame t to frame $t+1$.
 - The new estimated position of the tracked point in frame $t+1$ is simply its old position plus the calculated flow vector: $(x_{old}+u, y_{old}+v)$.
 - **Iterative Refinement:** The standard Lucas-Kanade flow assumes small motion. For point tracking, especially if the motion is larger, the algorithm is typically made iterative. It estimates the flow, warps the patch from the second image towards the first based on the estimate, recalculates the residual flow, and repeats until the estimate converges or a maximum number of iterations is reached. This helps handle larger displacements more accurately.

In simple terms:

If Lucas-Kanade Optical Flow is about calculating a *motion field* (a vector for many pixels or regions), Lucas-Kanade Point Tracking is about applying the *same local calculation* only at the *locations of specific points of interest* to follow their individual trajectories.

So, knowing the LK optical flow algorithm means you already understand the core computation behind LK point tracking. The "point tracking" aspect simply defines *where* in the image you apply that computation (only around the feature points) and what your goal is (following those specific points over time).

Lucas Kanade

Adds additional constraint that for each pixel within a small local neighborhood W , assume that Motion Field and hence Optical Flow are constant

Assumption: For each pixel, assume Motion Field, and hence Optical Flow (u, v) , is constant within a small neighborhood W .



That is for all points $(k, l) \in W$:

$$I_x(k, l)u + I_y(k, l)v + I_t(k, l) = 0$$

© Shree K. Nayar

[Lucas 1981]

So, now you get a system of equations where the number of equations will be the size of the window.

Let the size of window W be $n \times n$

In matrix form:

$$\begin{bmatrix} I_x(1,1) & I_y(1,1) \\ I_x(k, l) & I_y(k, l) \\ \vdots & \vdots \\ I_x(n, n) & I_y(n, n) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} I_t(1,1) \\ I_t(k, l) \\ \vdots \\ I_t(n, n) \end{bmatrix}$$

So, why does this work? Why would this matrix on the left be invertible? All of these equations tend not to be linearly dependent on one another. Because if you have a patch that has an interesting texture, then your I_x 's and I_y 's would most likely be different from pixel to pixel.

$$\begin{bmatrix}
 I_x(1,1) & I_y(1,1) \\
 I_x(k, l) & I_y(k, l) \\
 \vdots & \vdots \\
 I_x(n, n) & I_y(n, n)
 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} I_t(1,1) \\ I_t(k, l) \\ \vdots \\ I_t(n, n) \end{bmatrix}$$

A \mathbf{u} B
 (Known) (Unknown) (Known)
 $n^2 \times 2$ 2×1 $n^2 \times 1$

n^2 Equations, 2 Unknowns: Find Least Squares Solution

Least Squares Solution

Solve linear system: $A\mathbf{u} = B$

$$A^T A \mathbf{u} = A^T B \quad (\text{Least-Squares using Pseudo-Inverse})$$

In matrix form:

$$\begin{bmatrix}
 \sum_w I_x I_x & \sum_w I_x I_y \\
 \sum_w I_x I_y & \sum_w I_y I_y
 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} -\sum_w I_x I_t \\ -\sum_w I_y I_t \end{bmatrix}$$

Indices (k, l)
not written
for simplicity

$A^T A$ \mathbf{u} $A^T B$
 (Known) (Unknown) (Known)
 2×2 2×1 2×1

$$\boxed{\mathbf{u} = (A^T A)^{-1} A^T B}$$

Fast and Easy to Solve

When would this not work again? What if $A^T @ A$ is not invertible?

$$\boxed{A\mathbf{u} = B} \qquad \boxed{A^T A \mathbf{u} = A^T B}$$

- $A^T A$ must be **invertible**. That is $\det(A^T A) \neq 0$

More than wanting invertibility, we want $A.T @ A$ to be "well-conditioned"

GPT: **Well-conditioned** means that the solution is *stable* with respect to noise or small changes in the input. A poorly conditioned (but still invertible) matrix can produce huge variations in the output for tiny variations in the input.

Also GPT:

Well-conditionedness is important because it ensures **stability** in the solution. In the context of Lucas-Kanade, you are solving a system of linear equations derived from image gradients. If that system's coefficient matrix (often the structure tensor) is **ill-conditioned**, tiny amounts of noise or small changes in the gradients can lead to **large errors** in the computed optical flow. A well-conditioned matrix, on the other hand, keeps the solution more **robust** and **less sensitive** to inevitable noise or approximations in the image data.

- $A^T A$ must be **well-conditioned**.

If λ_1 and λ_2 are eigen values of $A^T A$, then

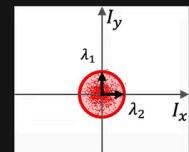
$$\lambda_1 > \epsilon \text{ and } \lambda_2 > \epsilon$$

$$\lambda_1 \geq \lambda_2 \text{ but not } \lambda_1 \gg \lambda_2$$

For $A^T @ A$ to be well conditioned, the eigen values must be significant enough. But just that one eigen value cannot be SIGNIFICANTLY larger than the other one.

How non well-conditioned matrices manifest IRL

Smooth Regions (Bad)



$\lambda_1 \sim \lambda_2$
Both are Small

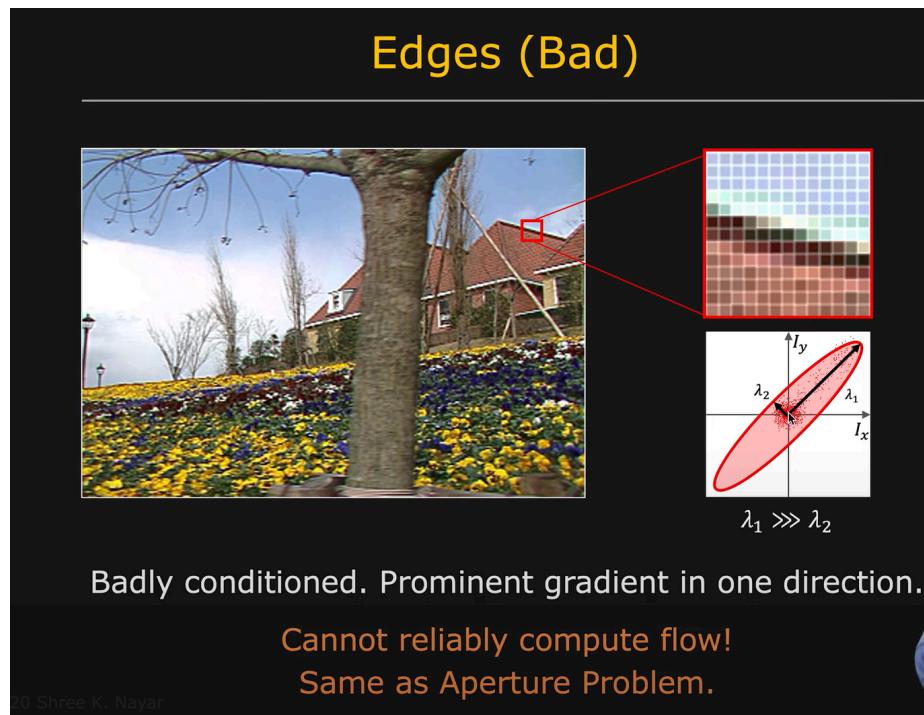
Equations for all pixels in window are more or less the same

Cannot reliably compute flow!

Note that $A^T @ T$ is the covariance matrix and its eigen vectors will show the direction where the data varies the most. If you have a smooth region then the I_x and I_y values are very small (these are the spatial gradients) then they will form a cloud round the origin. To find the eigen vectors of $A^T @ T$, fit an ellipse to these points. Then the semi major and major axis are the eigenvectors (of the covariance matrix $A^T @ A$). The covariance matrix also has very small eigen values.

So optical flow cannot be computed reliably.

Another example of a bad region:



Edges are also bad because you have strong gradients (of pairs (I_x, I_y)) in one direction, so one of the eigen values is much larger than the other one.

How do well conditioned matrices manifest IRL?

Textured Regions (Good)



Well conditioned. Large and diverse gradient magnitudes.

Can reliably compute optical flow.

Good changes in brightness in both directions, so pairs of spatial gradient pairs (I_x, I_y) are well scattered. So that the eigen values are large but not large w.r.t each other. Therefore the matrix $A.T @ T$ is well conditioned so that the optical flow computed is reliable.

Math Primer

Taylor Series Expansion

Expand a function as an infinite sum of its derivatives

$$f(x + \delta x) = f(x) + \frac{\partial f}{\partial x} \delta x + \frac{\partial^2 f}{\partial x^2} \frac{\delta x^2}{2!} + \cdots + \frac{\partial^n f}{\partial x^n} \frac{\delta x^n}{n!}$$

So, now you can see that if δx is really small, then the higher order terms are almost zero (i.e if δx is small to begin with, the δx^2 is even smaller)

If δx is small:

$$f(x + \delta x) = f(x) + \underset{\delta x}{\frac{\partial f}{\partial x}} \delta x + \boxed{O(\delta x^2)} \rightarrow \text{Almost Zero}$$

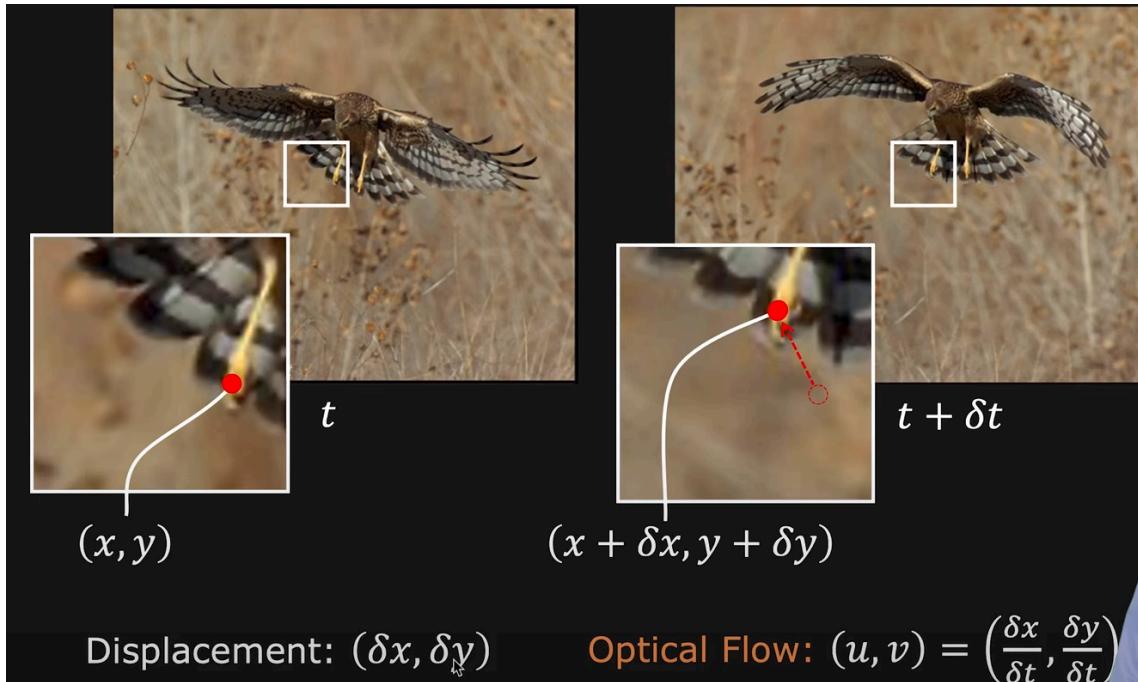
So you model $f(x + \delta x)$ by $f(x) + df/dx * \delta x$. This is a linear approximation because the approximation is linear in δx . AKA First order Taylor approximation.

For a function of three variables with small $\delta x, \delta y, \delta t$:

$$f(x + \delta x, y + \delta y, t + \delta t) \approx f(x, y, t) + \frac{\partial f}{\partial x} \delta x + \frac{\partial f}{\partial y} \delta y + \frac{\partial f}{\partial t} \delta t$$

Optical Flow Constraint Equation

Consider a small window (that's the same in both images). Focus on a point in this window. Say we are focused on the foot. Then at time $t + \delta t$, this point has moved to a new location which is $x + \delta x, y + \delta y$.



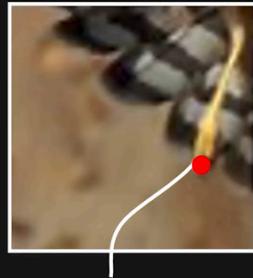
So, the displacement is (dx, dy)

So speed in x direction is dx/dt and speed in y direction is dy/dt . And this is the optical flow.

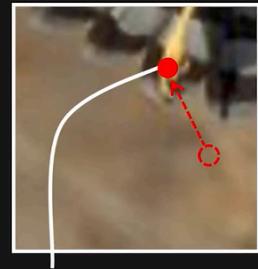
This is what we want to measure, (u, v)

Assumptions

Assumption #1: That brightness of image points remain the same, at least between consecutive images taken in quick succession. This is a reasonable assumption since dt is small.



$$I(x, y, t)$$



$$I(x + \delta x, y + \delta y, t + \delta t)$$

Assumption #1:

Brightness of image point remains constant over time

$$I(x + \delta x, y + \delta y, t + \delta t) = I(x, y, t)$$

Assumption #3: That the displacements of $(\delta x, \delta y)$, and δt are small.

Assumption #2:

Displacement $(\delta x, \delta y)$ and time step δt are small

This has to be true for us to derive a constraint equation.

What this allows us to do is to come up with an approximation for $I(x + \delta x, y + \delta y, t + \delta t)$. This approximation is based on the Taylor series expansion

Math Primer

$$I(x + \delta x, y + \delta y, t + \delta t) = I(x, y, t) + \frac{\partial I}{\partial x} \delta x + \frac{\partial I}{\partial y} \delta y + \frac{\partial I}{\partial t} \delta t$$

Alternatively,

$$I(x + \delta x, y + \delta y, t + \delta t) = I(x, y, t) + I_x \delta x + I_y \delta y + I_t \delta t$$

Results of assumptions

The two assumptions give rise to:

Optical Flow Constraint Equation

$$I(x + \delta x, y + \delta y, t + \delta t) = I(x, y, t) \quad \text{--- (1)}$$

$$I(x + \delta x, y + \delta y, t + \delta t) = I(x, y, t) + I_x \delta x + I_y \delta y + I_t \delta t \quad \dots \quad (2)$$

Then simply subtract:

Subtract (1) from (2): $I_x \delta x + I_y \delta y + I_t \delta t = 0$

Then divide by Δt and take limit to 0

Divide by δt and take limit as $\delta t \rightarrow 0$: $I_x \frac{\partial x}{\partial t} + I_y \frac{\partial y}{\partial t} + I_t = 0$

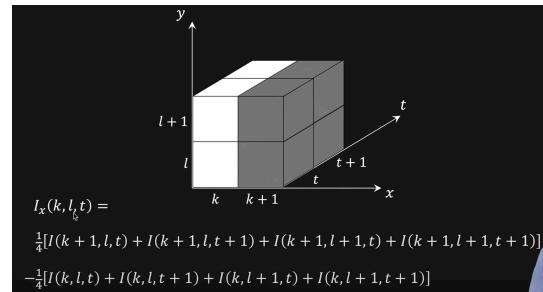
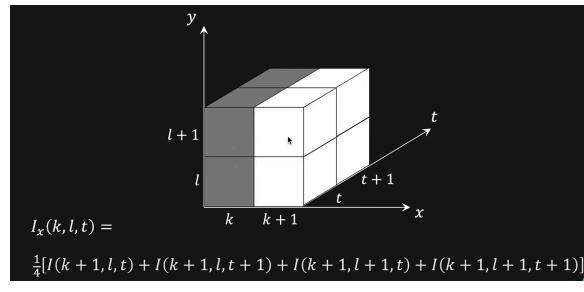
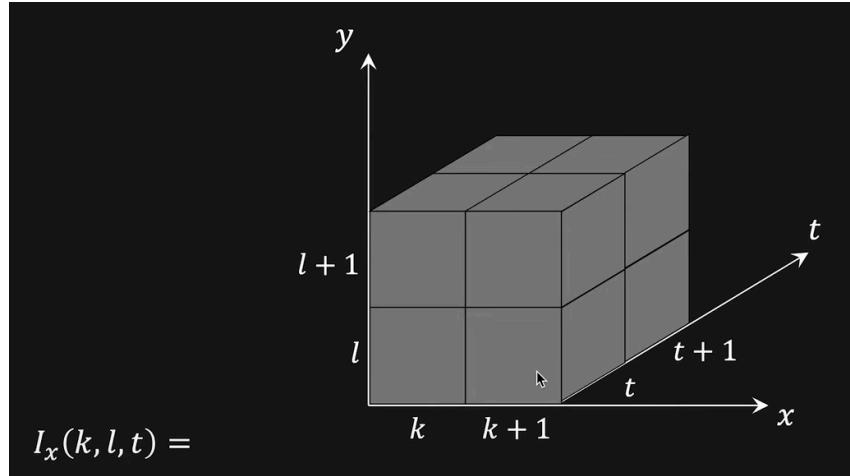
Then deltas turn into derivatives.

Then $\Delta x/\Delta t$ becomes $dx/dt = u$, $\Delta y/\Delta t$ becomes dy/dt .

Constraint Equation: $I_x u + I_y v + I_t = 0$ | (u, v) : Optical Flow

Then! Given two equations taken in quick succession (at t and $t + \Delta t$), we can find these derivatives I_x , I_y and I_t . Done using finite differences using edge detection.

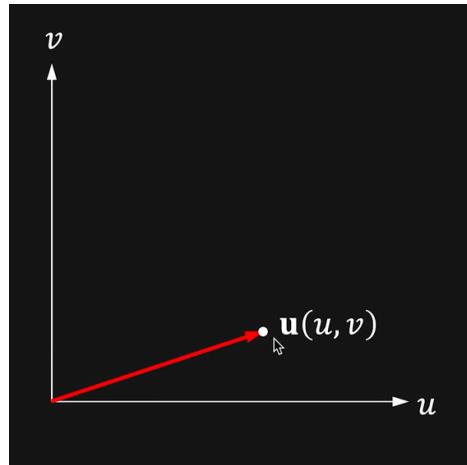
Computing I_x , I_y and I_t



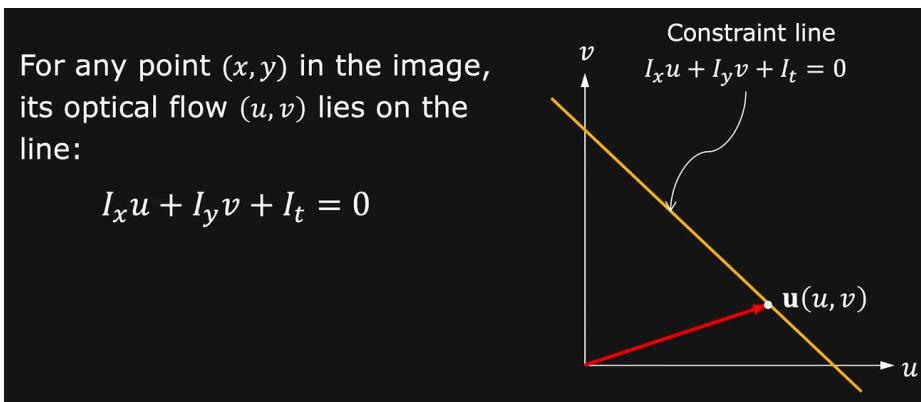
This will give you the derivative of the BRIGHTNESS in the x direction. Similarly you can do so for the y direction and t direction.

Geometrical Interpretation of Optical Flow (to show that u and v are really indeterminable, by showing under-constrainedness)

Let the true flow for a particular point in the image be (u, v)



(Unfortunately) the only thing that we know about (u, v) is that it must lie on the constraint line.



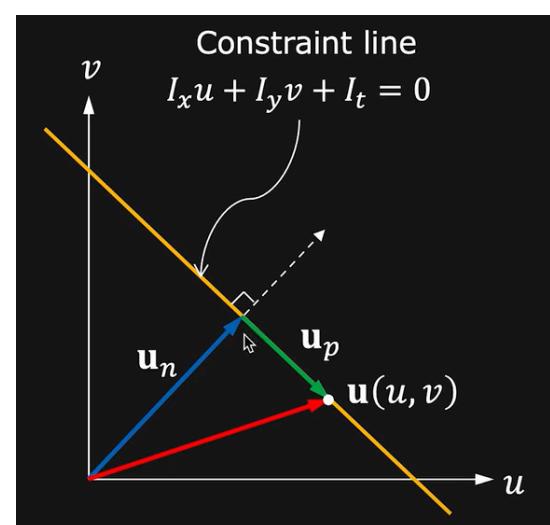
We know that (u, v) lies on the line. But don't exactly know where. This is what makes optical flow an under-constrained problem.

First, note that the vector can be split into two components.

$$\mathbf{u} = \mathbf{u}_n + \mathbf{u}_p$$

\mathbf{u}_n : Normal Flow

\mathbf{u}_p : Parallel Flow



U_n is the normal flow and U_p is the parallel flow

U_n is obviously easy to calculate.

Direction of Normal Flow:

Unit vector perpendicular to the constraint line:

$$\hat{\mathbf{u}}_n = \frac{(I_x, I_y)}{\sqrt{I_x^2 + I_y^2}}$$

$$\mathbf{u}_n = \frac{|I_t|}{(I_x^2 + I_y^2)} (I_x, I_y)$$

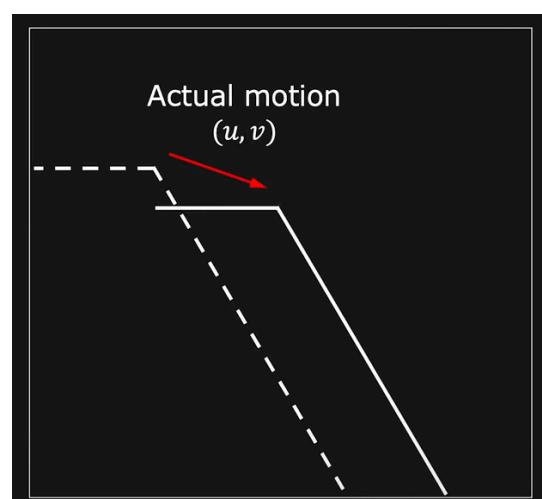
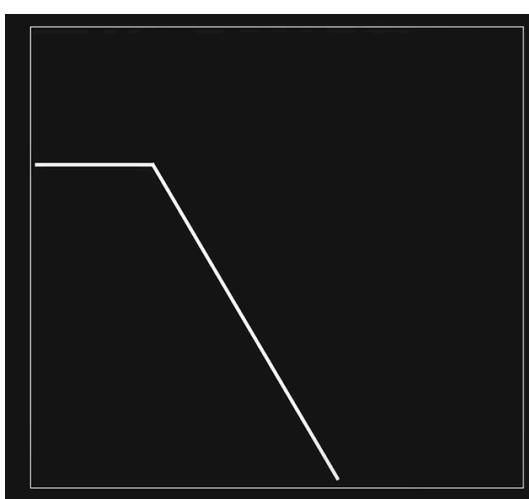
Magnitude of Normal Flow:

Distance of origin from the constraint line:

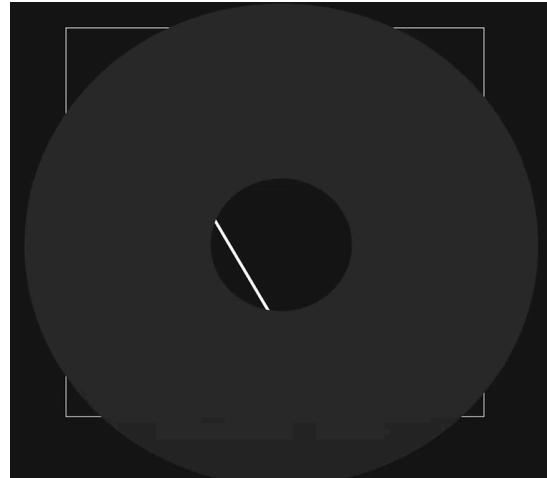
$$|\mathbf{u}_n| = \frac{|I_t|}{\sqrt{I_x^2 + I_y^2}}$$

But we cannot compute U_p !

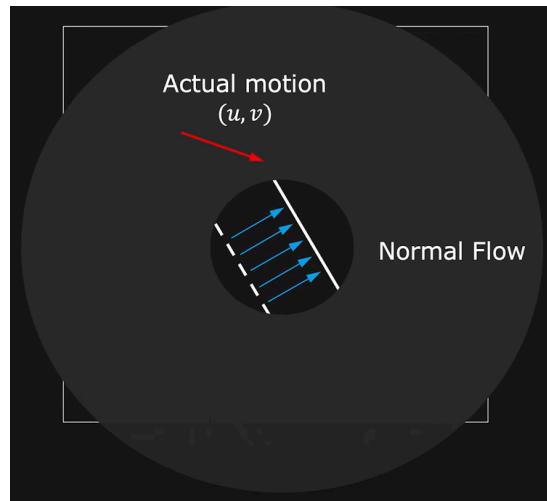
Example of how the problem of inability to calculate U_p manifest in a real scenario?



Our image is not one object, we have potentially different flows for every local region. So it necessitates to look at small local patches, called our aperture.



But if you look at the same motion through this aperture, you end up seeing ONLY the normal flow:



Unable to measure the parallel flow, and thus unable to measure the actual flow.

So locally, we are only able to determine the normal flow.

Optical Flow is under constrained!

Constraint Equation: $I_x u + I_y v + I_t = 0$

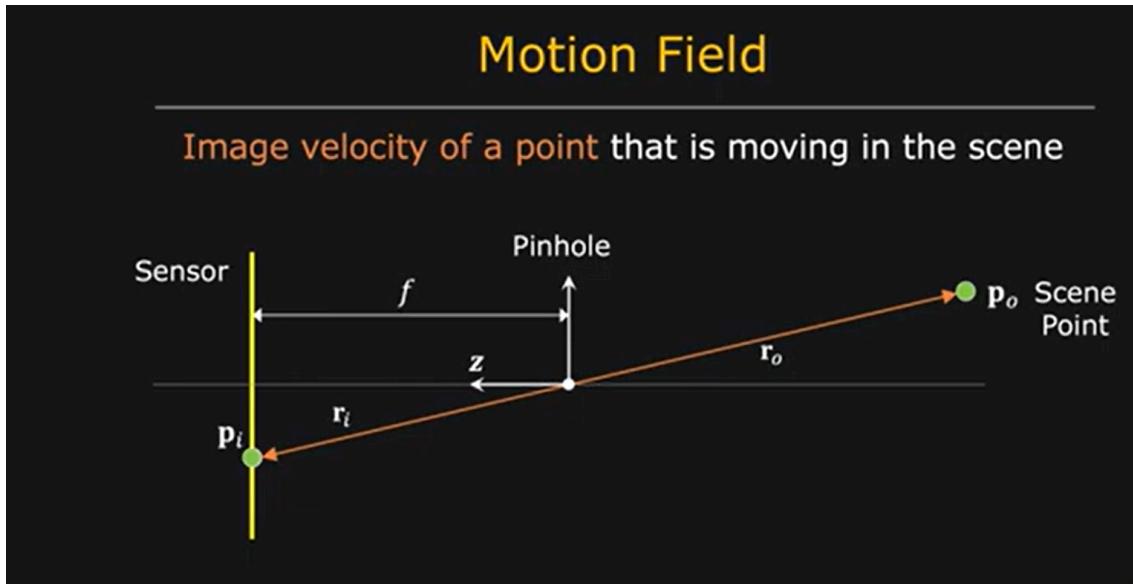
2 unknowns, 1 equation.

2 unknowns are u and v .

To solve this, we introduce additional constraints.

Motion Fields and Optical Flow

Motion Fields



p_0 is a point in the 3d scene and it's location is given by the vector r_0 . p_0 projects onto the point p_i via perspective projection and it's location on the image plane is described by the vector r_i .

If p_0 has some velocity in the 3d world, say v_0 , then its displacement is given by $v_0 * dt$. So its new location is $r_0 + \Delta r_0$.

Also by definition, scene velocity v_0 is the rate of change of r_0 , so

$$\text{Scene Point Velocity: } \mathbf{v}_o = \frac{d\mathbf{r}_o}{dt}$$

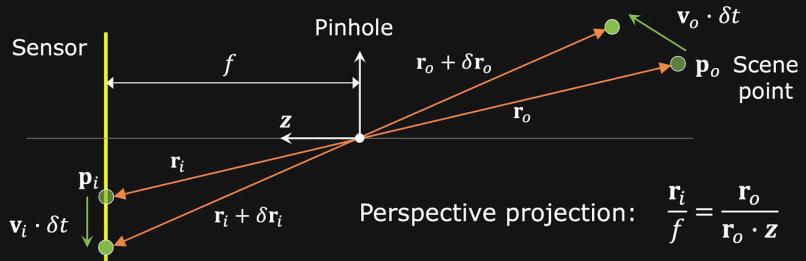
What we are interested in is the Image point velocity:

$$\text{Image Point Velocity: } \mathbf{v}_i = \frac{d\mathbf{r}_i}{dt} \quad (\text{Motion Field})$$

We know v_0 but we don't know v_i . So try to relate these to each other using perspective projection. Then simply take derivatives (using quotient rule) and then use cross product and get the answer.

Motion Field

Image velocity of a point that is moving in the scene



$$\text{Image Point Velocity: } \mathbf{v}_i = \frac{d\mathbf{r}_i}{dt} = f \frac{(\mathbf{r}_o \cdot \mathbf{z})\mathbf{v}_0 - (\mathbf{v}_o \cdot \mathbf{z})\mathbf{r}_0}{(\mathbf{r}_o \cdot \mathbf{z})^2}$$

(Motion Field)

$$\boxed{\mathbf{v}_i = f \frac{(\mathbf{r}_o \times \mathbf{v}_0) \times \mathbf{z}}{(\mathbf{r}_o \cdot \mathbf{z})^2}}$$

[Horn 1981]



© 2020 Shree K. Nayar

This \mathbf{v}_i is the motion field corresponding to the moving point in the 3D scene.

This is what we hope to measure, but there is no guarantee that we can measure this.

Optical Flow

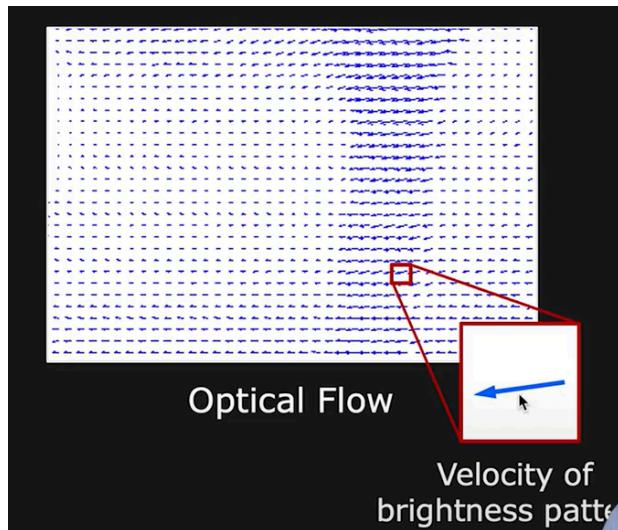
Instead, what we hope to measure is the motion of brightness patterns in the image.



(Here, if you take each point in the first frame and try to figure out where the point landed in the second image, that would be the motion field)

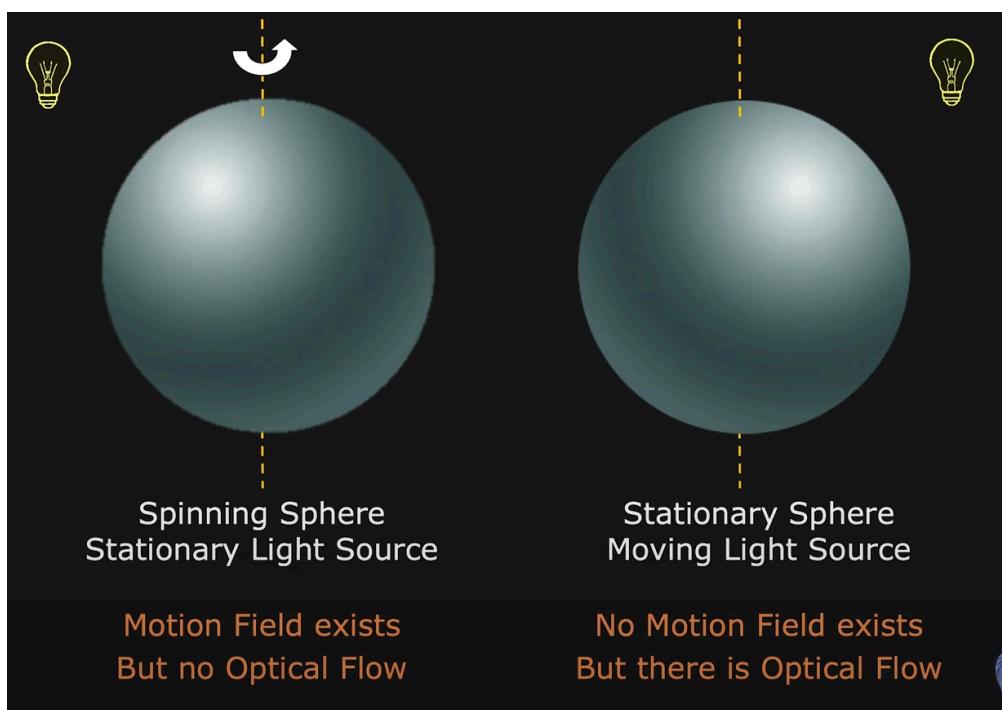
But what we can do is take a brightness pattern in one image and figure out where the brightness pattern is in the second image.

And if you have a successful algorithm that can do that, we can get a result that looks like this:



At each pixel, you have a vector that tells you what the optical flow at the point is. So, the motion of brightness patterns is the optical flow. And the optical flow vector's length tells you how fast its moving and direction tells you direction it's moving in the IMAGE PLANE.

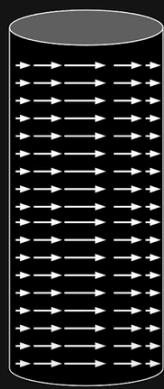
Ideally Optical Flow = Motion Field. Unfortunately, not most of the time.



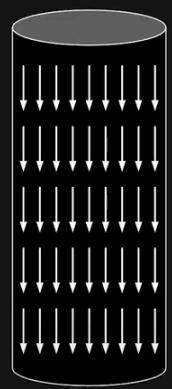
Another case where motion field and optical flow exists, but not the same



Barber Pole
Illusion



Motion Field



Optical Flow

Extra Hadamard Lecture

$A = hah$, a is the image, h is the sequency ordered hadamard matrix.

So if you just do $H @ a$, you've only done the first pass (transforms along the *horizontal* direction but leaves each column untouched) and haven't yet projected onto the vertical-basis patterns—hence it won't match the full 2D transform. After you have post-multiplied by h , this transforms both *rows* and *columns*, giving the true 2D Hadamard spectrum you see in your 4×4 coefficient matrix.

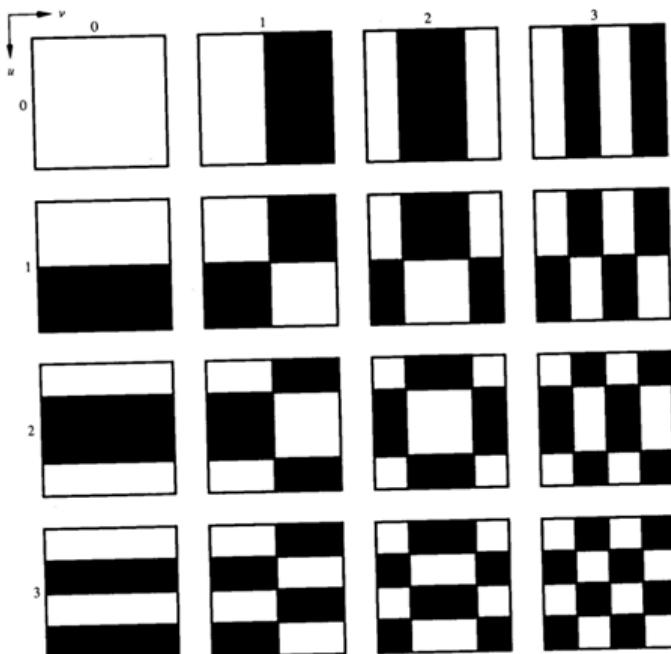
natural hadamard can do compression and coding

sequency ordered can do those plus filtering in hadamard domain

need bit reverse etc to convert natural hadamard to sequency ordered hadamard

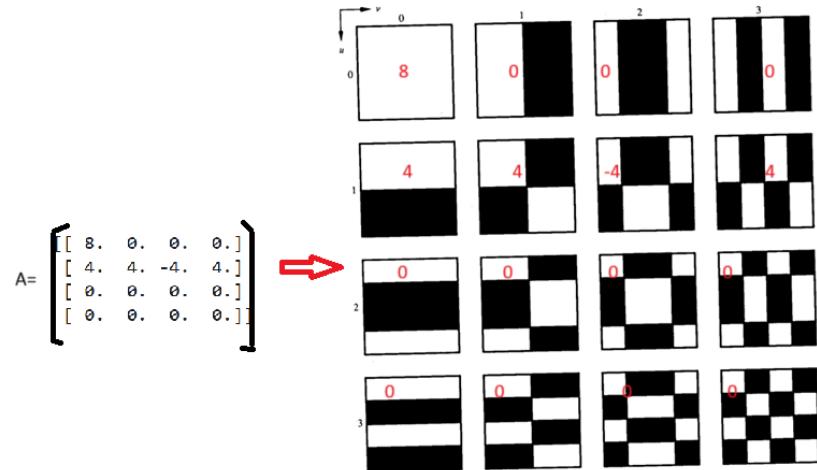
Any 4×4 image can be represented as a weighted sum of the following 16 4×4 square waves in the SPATIAL DOMAIN. These are the Walsh basis patterns.

(white is 1, black is 0)



And what sequency hadamard transform actually does it calculating the magnitude of these weightages in order (low \rightarrow high spatial frequency basis)

So A (4×4 matrix) = hah where a is the spatial image has these coefficients



Notice the u and v axis in the A matrix.

u represents the vertical frequency and v represents horizontal frequency.

Eg $A[0,0]$ represents the wavelet with no horizontal and vertical frequency

$A[0,1]$ represents the wavelet with low horizontal and no vertical frequency

$A[0,2]$ represents the wavelet with higher horizontal and no vertical frequency

$A[1,0]$ → no horizontal frequency but low vertical frequency

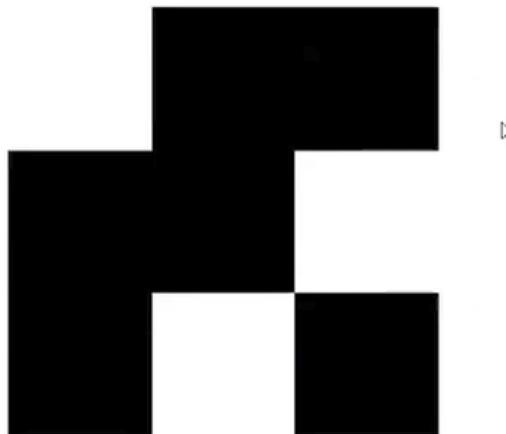
$A[2,0]$ → no horizontal but higher vertical frequency

$A[3,0]$ → no horizontal but more higher (lol) vertical frequency

Other $A[i, j]$ has a combination of horizontal and vertical frequency

$A[3,3]$ → highest horizontal and highest vertical frequency

Now, why is $A[0][0] = \text{sum of all the pixels?}$ Look at hah and the sequency ordered hadamard matrix's first row and first column.



the (0,0) entry, $A[0][0]$, is calculated using the first row of h on both sides. That is:

$$A[0][0] = (\text{first row of } h) \cdot a \cdot (\text{first column of } h)$$

So then since u and v represents frequency then you can filter out A in whatever way you like to preserve either high 2d frequencies by keeping coefficients of high frequency square waves intact or by removing them and preserving only the low frequency square waves by keeping their coefficients intact

How to convert Natural hadamard matrices to sequency ordered hadamard matrices?

(now the slide deck begins)

(by the way, take a look at the recursive formulation of the hadamard matrices)

Here, you only need to rearrange the rows, and the columns will be rearranged automatically.

1) Bit-by-bit GRAY 2 Bin algorithm

Let your Gray code be an n -bit string $g_{n-1}g_{n-2}\dots g_1g_0g_{\{n-1\}}g_{\{n-2\}}\dots g_1g_0g_{n-1}g_{n-2}\dots g_1g_0$. You build the binary bits $b_{n-1}\dots b_0b_{\{n-1\}}\dots b_0b_{n-1}\dots b_0$ as follows:

1. **MSB copies directly:** $b_{n-1}=g_{n-1}$.

$$b_{n-1} = g_{n-1}, b_{\{n-1\}} \leftarrow g_{\{n-1\}},$$

2. **Each lower bit is the XOR of the Gray bit with the previously computed binary bit:** $b_i = b_{i+1} \oplus g_i$ for $i=n-2, n-3, \dots, 0$.

$$\begin{aligned} b_i &= b_{i+1} \oplus g_i \text{ for } i=n-2, n-3, \dots, 0. \\ &\quad b_i \leftarrow b_i \oplus g_i \quad \backslash \text{quad} \\ &\quad \text{\texttt{for }} i = n-2, n-3, \dots, 0. \end{aligned}$$

Example

Gray = 1101 (4-bit)

1. $b_3 = g_3 = 1$
2. $b_2 = b_3 \oplus g_2 = 1 \oplus 1 = 0$
3. $b_1 = b_2 \oplus g_1 = 0 \oplus 0 = 0$
4. $b_0 = b_1 \oplus g_0 = 0 \oplus 1 = 1$

So 1101 (Gray) \rightarrow 1001 (binary).

How to Obtain Sequence-Ordered H Matrices?

To convert that 8x8 Hadamard to SO, firstly, write the row indices in a 3-bit representation

Row	000	001	010	011	100	101	110	111
000	1	1	1	1	1	1	1	1
001	1	-1	1	-1	1	-1	1	-1
010	1	1	-1	-1	1	1	-1	-1
011	1	-1	-1	1	1	-1	-1	1
100	1	1	1	1	-1	-1	-1	-1
101	1	-1	1	-1	-1	1	-1	1
110	1	1	-1	-1	-1	-1	1	1
111	1	-1	-1	1	-1	1	1	-1

Then, apply bit-reverse and gray-to-bin functions to those indices. Now you have the new index of each row in the new SO Hadamard matrix

Row	000	001	010	011	100	101	110	111
000	000	100	010	011	100	001	001	101
001	100	000	011	111	110	101	110	011
010	010	000	111	100	001	101	110	111
011	011	111	100	100	001	011	011	000
100	100	111	001	001	110	101	010	000
101	101	101	101	101	110	110	110	101
110	110	110	110	110	011	011	011	011
111	111	111	111	111	101	101	101	101

New indices for the SO Hadamard matrix, e.g., row 1 of the NH matrix will be row 7 of SOH, and row 6 of NH matrix will be row 3 of SOH

How to Obtain Sequence-Ordered H Matrices?

NH, 8x8

1	1	1	1	1	1	1	1
1	-1	1	-1	1	-1	1	-1
1	1	-1	-1	1	1	-1	-1
1	-1	-1	1	1	-1	-1	1
1	1	1	1	-1	-1	-1	-1
1	-1	1	-1	-1	1	-1	1
1	1	-1	-1	-1	-1	1	1
1	-1	-1	1	-1	1	1	-1

SOH, 8x8

1	1	1	1	1	1	1	1
1	1	1	1	-1	-1	-1	-1
1	1	-1	-1	-1	-1	1	1
1	1	-1	-1	1	1	-1	-1
1	-1	-1	1	1	1	-1	1
1	-1	1	-1	-1	1	1	-1
1	-1	1	-1	1	-1	1	1
1	-1	1	1	-1	1	-1	1

Amirhassan Monajemi, NUS, SoC, 2024

4

Butterworth low pass

For one 1:

$$H(f) = \frac{1}{1 + \left(\frac{f}{f_c}\right)^{2n}}$$

This is a low pass filter

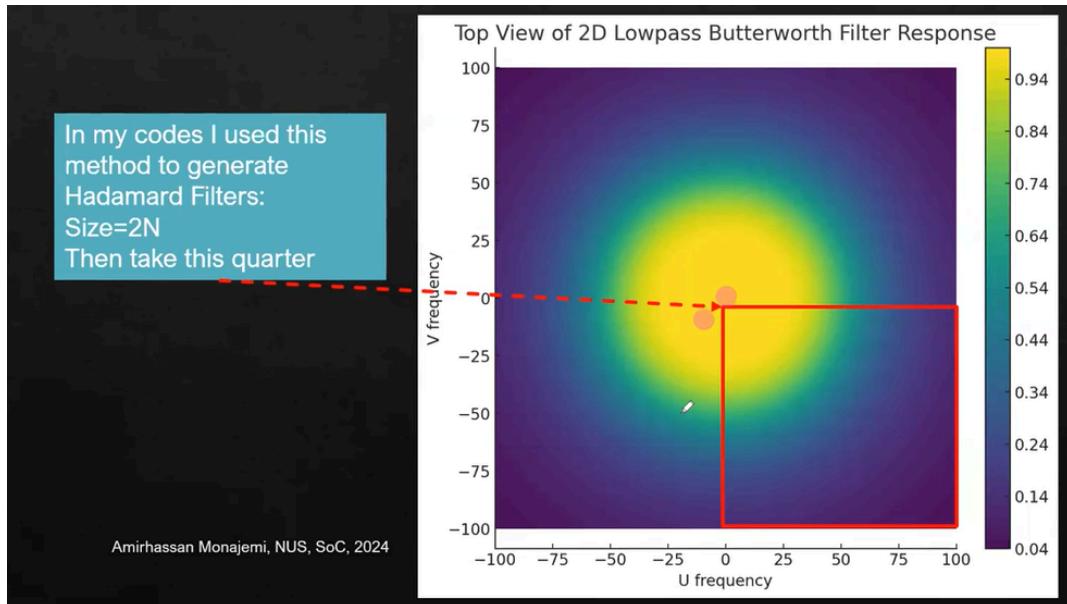

$$H(u, v) = \frac{1}{1 + \left(\frac{D(u, v)}{D_0}\right)^{2n}}$$

Here's what each term in this equation represents:

- $H(u, v)$ is the filter's transfer function at frequencies u and v , corresponding to the two spatial frequency dimensions in an image.
- $D(u, v)$ is the distance from the origin in the frequency domain to the point (u, v) , typically calculated as $D(u, v) = \sqrt{u^2 + v^2}$.
- D_0 is the cutoff frequency, similar to f_c in the 1D case.
- n is the order of the filter, controlling the steepness of the roll-off at the cutoff frequency.

Look at how if $D(u, v) < D_0$ then high value.

But as soon as $D(u, v) > D_0$, then quickly denominator explodes. Explosion is controlled by roll-off number n (order). Higher the order, higher the fall.



Butterworth highpass

Amithasan Monajemi

$\diamond H(f) = \frac{1}{1 + \left(\frac{f_c}{f}\right)^{2n}}$

$\diamond H(u,v) = \frac{1}{1 + \left(\frac{D_0}{D(u,v)}\right)^{2n}}$

highpass = 1 – lowpass

But note that for ideal low pass butterworth, there is no gradual decay of magnitude, but a SHARP cutoff.

Why divide by 256^2 ?

For the *unnormalized* Hadamard matrix H of size N :

1. $HH = NI$.
2. Forward 2D transform:

$$A = HaH.$$

3. Inverse then is

$$a = \frac{1}{N^2} HAH,$$

since

$$HAH = H(HaH)H = (HH)a(HH) = NIaNI = N^2a.$$

So in code you divide by N^2 . For a 256×256 image, $N = 256$, hence $\dots / 256^{**2}$.

Inverse via the matrix inverse

1. General inverse formula

If your forward was

$$A = HaH$$

then the inverse is

$$a = H^{-1}AH^{-1}.$$

2. Hadamard's special inverse

For the (unnormalized) Hadamard matrix $H \in \{+1, -1\}^{N \times N}$, one can show

$$HH = NI \implies H^{-1} = \frac{1}{N}H.$$

Moreover H is symmetric so $H^T = H$.

3. Putting it together

$$a = H^{-1}AH^{-1} = \left(\frac{1}{N}H\right)A\left(\frac{1}{N}H\right) = \frac{1}{N^2}HAH.$$

Hence your code's division by 256^{**2} is exactly the two factors of $1/N$ coming from each inverse.

A **normalized** Hadamard matrix is just the usual $\{\pm 1\}$ Hadamard matrix H rescaled so that its rows (and columns) become orthonormal rather than orthogonal. Concretely, if H is an $N \times N$ Hadamard (so that $H H^T = N I$), then you define

$$H_n = \frac{1}{\sqrt{N}} H.$$

Key properties of H_n

- Orthonormality:

$$H_n H_n^T = \frac{1}{N} H H^T = I.$$

- Symmetric and involutory:

$$H_n^T = H_n, \quad H_n^{-1} = H_n.$$

Why normalize?

1. Simplifies inversion.

With the un-normalized H you need

$$A = H a H \implies a = \frac{1}{N^2} H A H.$$

With H_n you get

$$A = H_n a H_n \implies a = H_n A H_n$$

with no extra $1/N^2$ factor.

2. Unit-energy basis.

Each row of H_n has Euclidean norm 1, so your transform becomes an orthonormal change of basis—very handy for energy-preserving signal analysis.

3. Interpretability.

The DC coefficient $A_{0,0}$ now gives the *average* pixel value (rather than the raw sum), because projecting onto an all-ones unit vector divides by \sqrt{N} .

What are possible applications of representing an image using the first few walsh basis images?

- **Image Compression**

- **Low-complexity codecs**

The Walsh–Hadamard transform (WHT) uses only additions/subtractions, so it's extremely fast in hardware. Truncating to the first KKK low-sequence coefficients yields a coarse but compact representation—ideal for very low-power or real-time systems (e.g. embedded vision).

- **Progressive coding**

Send only the DC term and a handful of first-sequence basis images to give a quick "preview" of the frame; refine gradually by adding higher-sequence terms as bandwidth allows.

- **Denoising and Smoothing**

- Noise in images tends to manifest in high-sequence (rapid sign-changes) coefficients. By zeroing out all but the first few Walsh coefficients, you perform a form of low-pass filtering that suppresses speckle or salt-and-pepper noise while keeping the large-scale structure intact.

- **Watermarking & Steganography**

- Embedding a watermark in mid- or low-sequence Walsh coefficients makes it robust to JPEG-style compression (which also retains low frequencies) but less perceptible to the human eye. You can add or modify a small offset in a few chosen basis coefficients to carry hidden data.

- **Feature Extraction for Recognition/Classification**

- Projecting image patches onto the first few Walsh functions gives you a low-dimensional feature vector that captures bulk texture/brightness patterns. These “Walsh features” have been used successfully for face recognition, fingerprint matching, and texture classification because they’re rotation- and scale-invariant to a degree, and very fast to compute.

- **Template Matching & Fast Correlation**

- Because convolution in the Walsh domain reduces to element-wise products, you can rapidly correlate a template against an image by transforming both to sequency order, multiplying only the first few coefficients (for a coarse match), and inversely transforming. This accelerates object detection or motion tracking in video.

- **Compressed Sensing & Sparse Sampling**

- The Walsh basis is incoherent with many natural-image bases. If you know an image is sparse in the Walsh domain, you can randomly sample a small set of pixel measurements and reconstruct it by solving a sparse-recovery problem—keeping only the largest few Walsh coefficients.

- **Optical & Holographic Processing**

- In optical computing, lenses can implement the WHT natively. Truncating to low-sequencey patterns lets you perform rapid, analog smoothing, correlation, or pattern recognition directly in hardware without digital post-processing.

- **Progressive Rendering in Graphics**

- Rather than rasterizing every pixel at full resolution, a renderer can project frames onto low-sequencey Walsh patterns first to quickly display a blurred version, then refine it. This improves interactivity in remote or VR streaming scenarios.