

# Announcements

**Matlab Grader homework**, emailed Thursday,  
1 and 2 (of less than 9) homeworks Due 21 April, Binary graded.

**Jupyter homework?:** translate matlab to Jupiter, TA Harshul h6gupta@eng.ucsd.edu or me  
I would like this to happen.

“GPU” homework. NOAA climate data in Jupyter on the datahub.ucsd.edu, released 17 April.

Projects: Any computer language. Access to Jupyterhub with GPU

**Podcast** might work eventually.

## Today:

- Stanford CNN
- Gaussian processes for concert hall
- Linear models for regression

Wednesday 10 April

Stanford CNN, Linear models for regression/classification (Bishop 3 and 4),

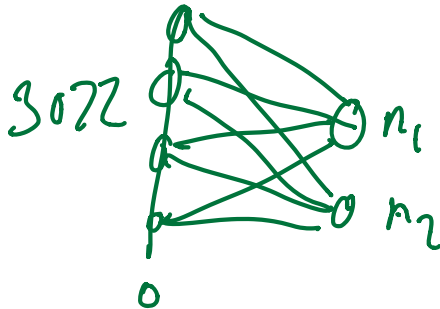
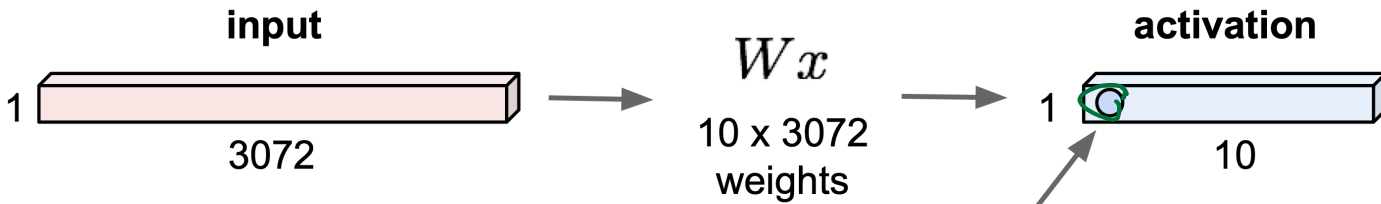
CNN

FNN

# Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1

Each neuron looks at the full input volume



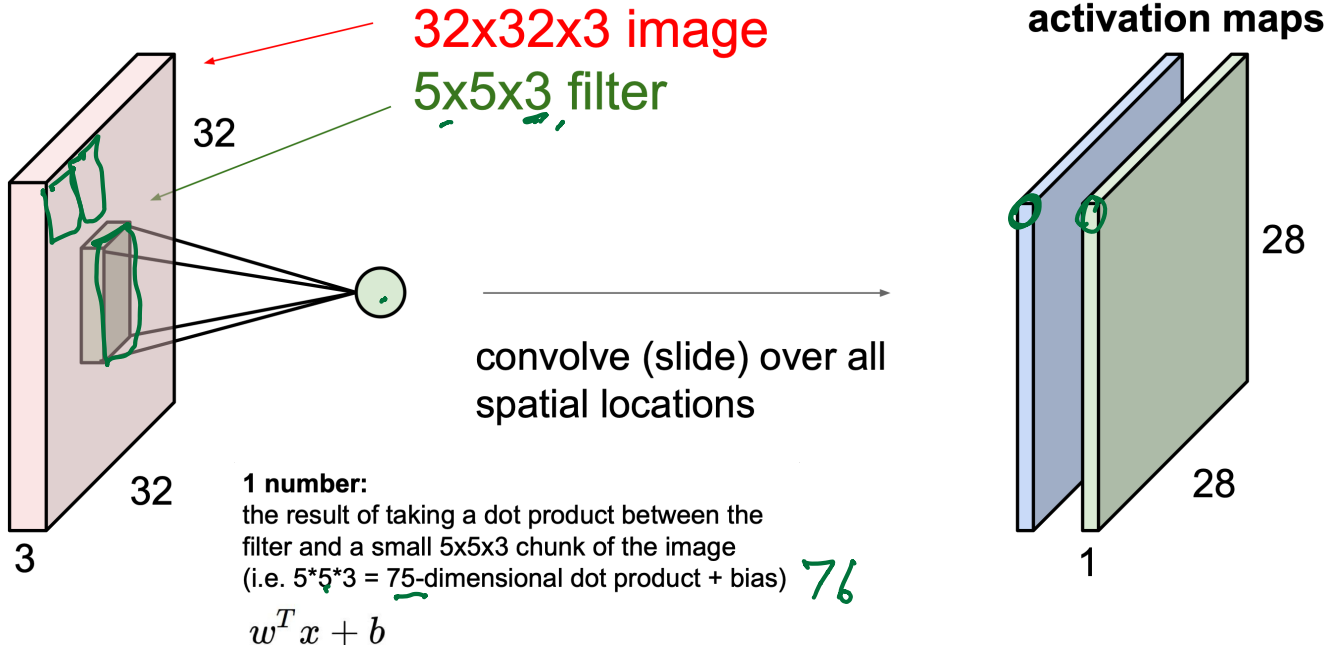
$$W^T x$$

**1 number:**  
the result of taking a dot product  
between a row of  $W$  and the input  
(a 3072-dimensional dot product)

# CNN

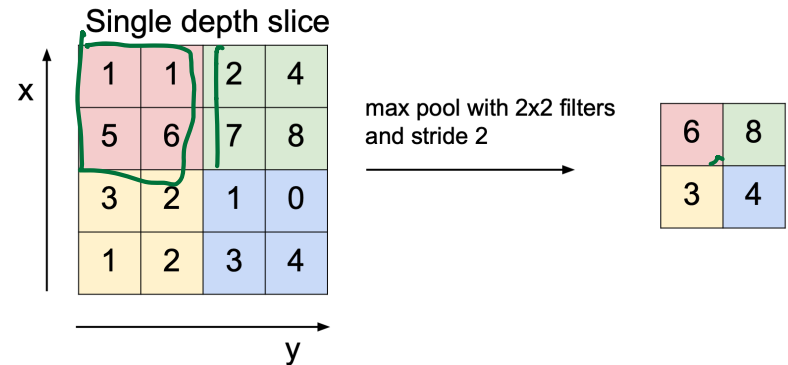
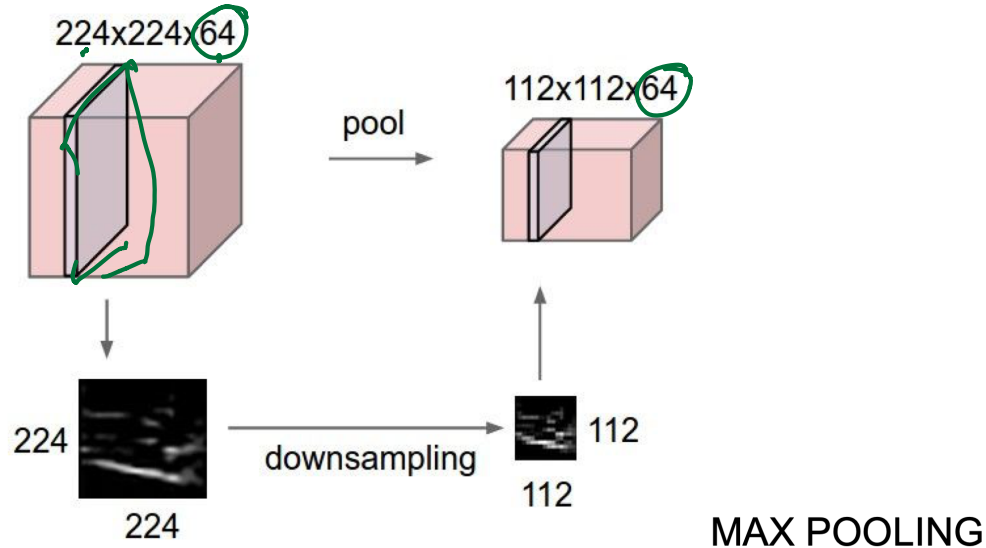
## Convolution Layer

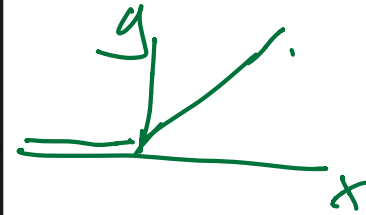
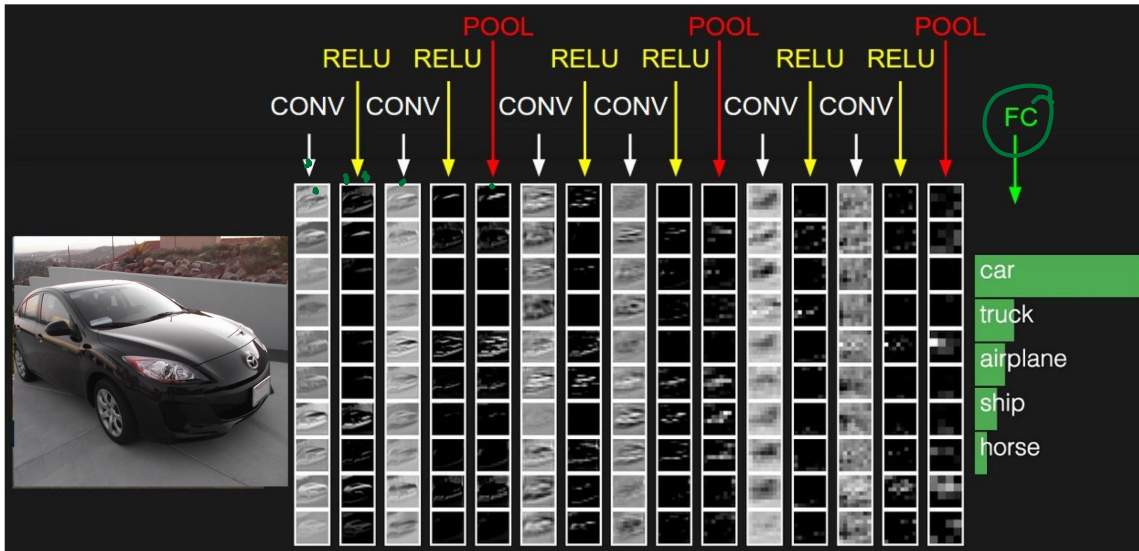
consider a second, **green** filter



# Pooling layer

- makes the representations smaller and more manageable
- operates over each activation map independently:





PA  
 $y = \max(c, x)$

## Summary

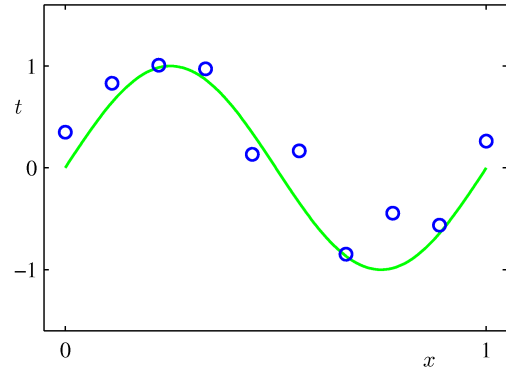
- ConvNets stack CONV, POOL, FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Typical architectures look like **[(CONV-RELU)\*N-POOL?]\*M-(FC-RELU)\*K, SOFTMAX** where N is usually up to ~5, M is large,  $0 \leq K \leq 2$ .
  - but recent advances such as ResNet/GoogLeNet challenge this paradigm

# Linear regression: Linear Basis Function Models (1)

Generally

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

- where  $\phi_j(x)$  are known as *basis functions*.
- Typically,  $\phi_0(x) = 1$ , so that  $w_0$  acts as a bias.
- Simplest case is linear basis functions:  $\phi_d(x) = x_d$ .

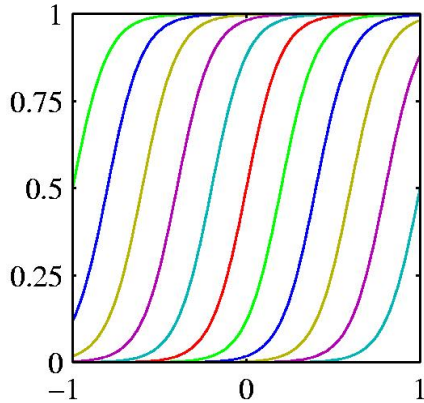


$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

<http://playground.tensorflow.org/>

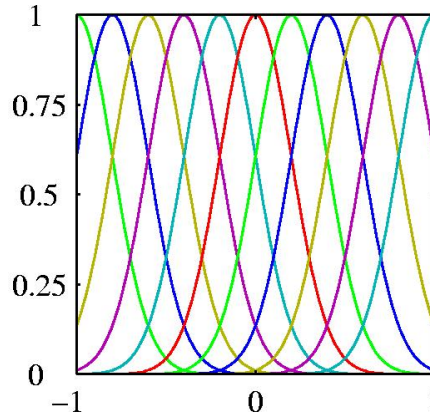
$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} \quad \phi = \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix}$$

## Some types of basis function in 1-D



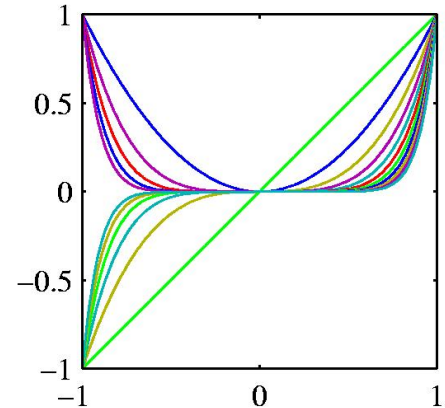
**Sigmoids**

$$\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right)$$
$$\sigma(a) = \frac{1}{1 + \exp(-a)}.$$



**Gaussians**

$$\phi_j(x) = \exp\left\{-\frac{(x - \mu_j)^2}{2s^2}\right\}$$



**Polynomials**

$$\phi_j(x) = x^j.$$

Sigmoid and Gaussian basis functions can also be used in multilayer neural networks, but neural networks **learn** the parameters of the basis functions. **This is more powerful but also harder and messier.**

Two types of linear model that are equivalent with respect to learning

$$y(\mathbf{x}, \mathbf{w}) = \overset{\text{bias}}{\downarrow} w_0 + w_1 x_1 + w_2 x_2 + \dots = \mathbf{w}^T \mathbf{x} \quad \text{↔}$$

$$\Rightarrow y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 \phi_1(\mathbf{x}) + w_2 \phi_2(\mathbf{x}) + \dots = \underline{\mathbf{w}^T \Phi(\mathbf{x})}$$

- The first and second model has the same number of adaptive coefficients as the number of basis functions +1.
- Once we have replaced the data by basis functions outputs, fitting the second model is exactly the same the first model.
  - No need to clutter math with basis functions



# Maximum Likelihood and Least Squares (1)

- Assume observations from a deterministic function with added Gaussian noise:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad \text{where} \quad p(\epsilon|\beta) = \mathcal{N}(\epsilon|0, \beta^{-1})$$

- or,

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1}).$$

- Given observed inputs,  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , and targets  $\mathbf{t} = [t_1, \dots, t_N]^T$ , we obtain the likelihood function

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n|\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}).$$

## Maximum Likelihood and Least Squares (2)

Taking the logarithm, we get

$$\begin{aligned}\ln p(\mathbf{t}|\mathbf{w}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta \underline{E_D(\mathbf{w})}\end{aligned}$$

Where the sum-of-squares error is

$$\begin{aligned}E_D(\mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 \\ \frac{\partial E}{\partial \mathbf{w}} &= \sum_{n=1}^N \phi_n^T \{t_n - \mathbf{w}^T \phi_n\} \\ &= \underline{\Theta}^T \underline{t} - \underline{\Theta}^T \underline{\Phi} \underline{w} = 0 \\ \underline{w}_{ML} &= (\underline{\Theta}^T \underline{\Phi})^{-1} \underline{\Theta}^T \underline{t}\end{aligned}$$

Handwritten diagram illustrating the matrix representation of the least squares problem:

$$\underline{\Phi} = \begin{bmatrix} \phi^T(\mathbf{x}_1) \\ \vdots \\ \phi^T(\mathbf{x}_N) \end{bmatrix} \quad \begin{matrix} \xrightarrow{M} \\ \downarrow \\ N \end{matrix}$$
$$\underline{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix} \quad N$$

# Maximum Likelihood and Least Squares (3)

Computing the gradient and setting it to zero yields

$$\nabla_{\mathbf{w}} \ln p(\mathbf{t}|\mathbf{w}, \beta) = \beta \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n)^T = \mathbf{0}.$$

Solving for  $\mathbf{w}$ ,

where

$$\mathbf{w}_{\text{ML}} = \left( \Phi^T \Phi \right)^{-1} \Phi^T \mathbf{t}$$

The Moore-Penrose pseudo-inverse,  $\Phi^\dagger$ .

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

*Handwritten annotations:* A green arrow labeled  $M$  points to the columns of the matrix. A green arrow labeled  $N$  points to the rows of the matrix.

# Maximum Likelihood and Least Squares (4)

Maximizing with respect to the bias,  $w_0$ , alone,

$$\begin{aligned} \underline{w_0} &= \bar{t} - \sum_{j=1}^{M-1} w_j \bar{\phi}_j \\ &= \frac{1}{N} \sum_{n=1}^N t_n - \sum_{j=1}^{M-1} w_j \frac{1}{N} \sum_{n=1}^N \phi_j(\mathbf{x}_n). \end{aligned}$$

We can also maximize with respect to  $\underline{\beta}$ , giving

$$\frac{\partial \mathcal{P}}{\partial \beta} = 0$$

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{n=1}^N \{t_n - \mathbf{w}_{\text{ML}}^T \phi(\mathbf{x}_n)\}^2$$

$\downarrow$   $M$

$M$   
 $T$

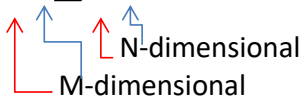
# Geometry of Least Squares

Consider

$$N \begin{bmatrix} \varphi_1 \\ \vdots \\ \varphi_M \end{bmatrix} = \begin{bmatrix} \varphi_1(x_1) \\ \vdots \\ \varphi_1(x_2) \end{bmatrix} / N$$

$$\mathbf{y} = \Phi \mathbf{w}_{ML} = [\varphi_1, \dots, \varphi_M] \mathbf{w}_{ML}$$

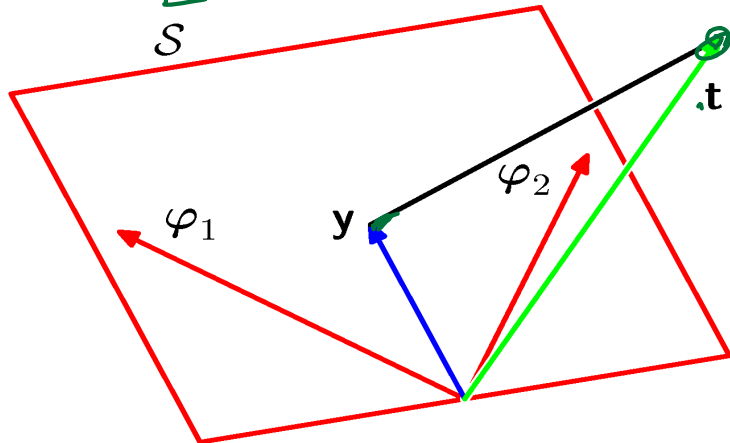
$$\mathbf{y} \in \mathcal{S} \subseteq \mathcal{T} \quad \mathbf{t} \in \mathcal{T}$$



$\mathcal{S}$  is spanned by

$$\varphi_1, \dots, \varphi_M$$

$\mathbf{w}_{ML}$  minimizes the distance between  $\mathbf{t}$  and its orthogonal projection on  $\mathcal{S}$ , i.e.  $\mathbf{y}$ .



$$\varphi_1 = \begin{bmatrix} \varphi_1(x_1) \\ \vdots \\ \varphi_1(x_2) \end{bmatrix} / N$$

## Least mean squares: An alternative approach for big datasets

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E_{n(\tau)}$$

$\mathbf{w}^{\tau+1}$  (weights after seeing training case tau+1) =  $\mathbf{w}^{\tau}$  (learning rate) -  $\eta$  (learning rate)  $\nabla E_{n(\tau)}$  (squared error derivatives w.r.t. the weights on the training case at time tau).

This is “**on-line**” learning. It is efficient if the dataset is redundant and simple to implement.

- It is called **stochastic gradient descent** if the training cases are picked randomly.
- Care must be taken with the learning rate to prevent divergent oscillations. Rate must decrease with tau to get a good fit.

$$\frac{\partial E}{\partial \mathbf{w}} = \sum_n \alpha_n (t_n - \mathbf{w}^T \mathbf{a}_n)$$

## Regularized least squares

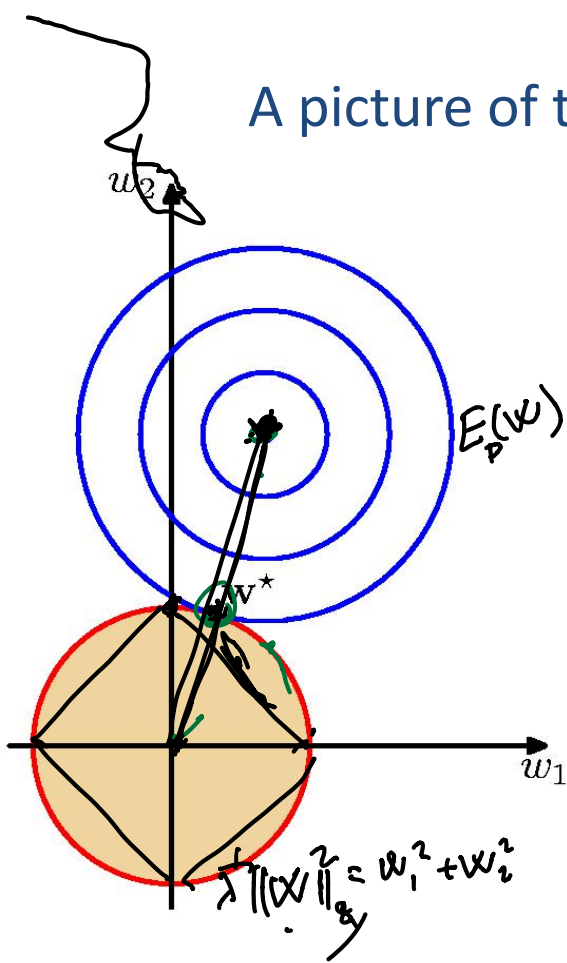
$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

The squared weights penalty is mathematically compatible with the squared error function, giving a closed form for the optimal weights:

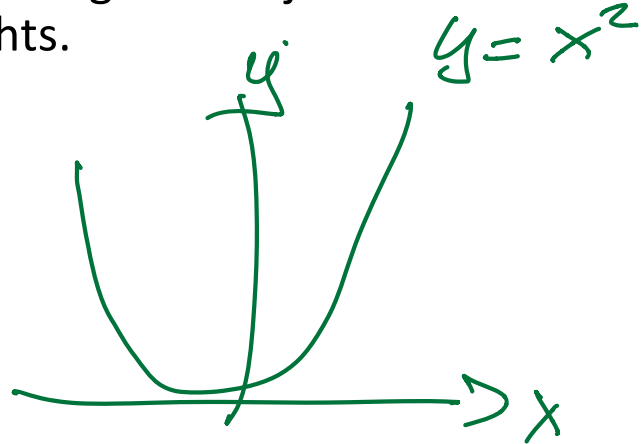
$$\underline{\mathbf{w}}^* = (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$$

↑  
identity matrix

## A picture of the effect of the regularizer



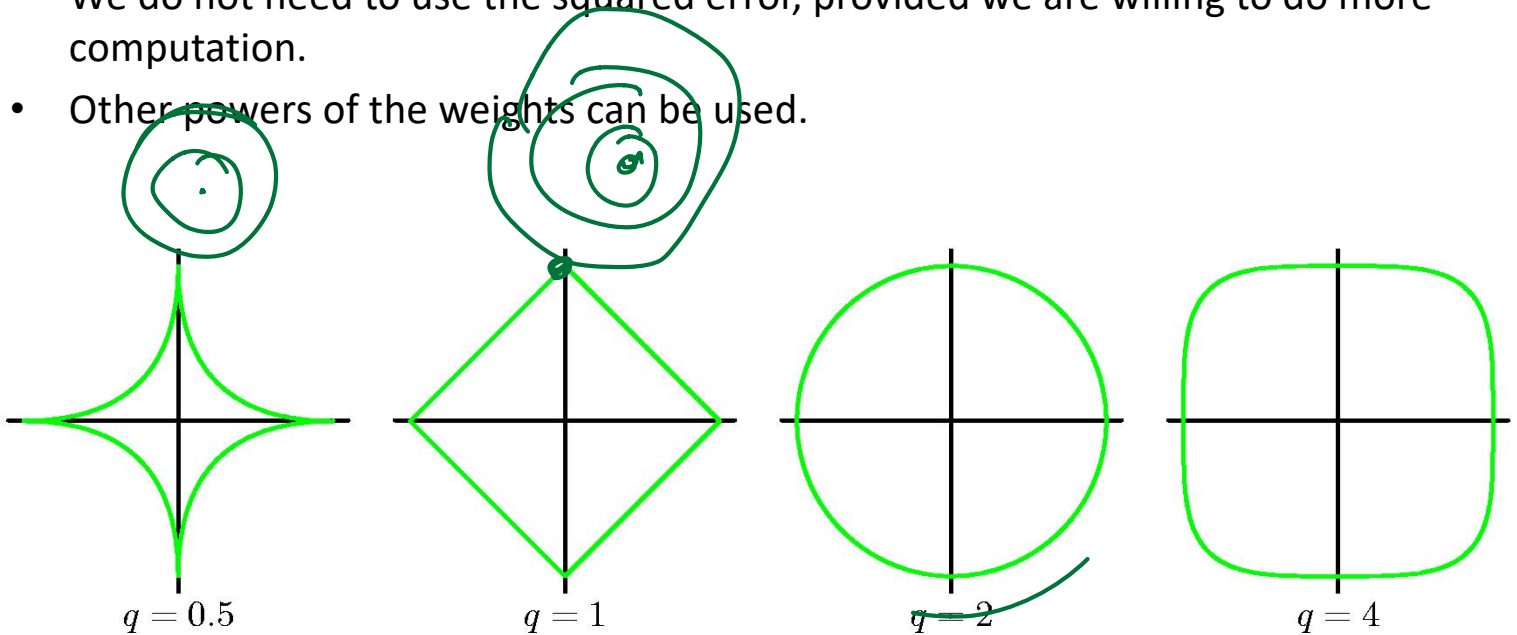
- The overall cost function is the sum of two parabolic bowls.
- The sum is also a parabolic bowl.
- The combined minimum lies on the line between the minimum of the squared error and the origin.
- The L2 regularizer just **shrinks** the weights.





# Other regularizers

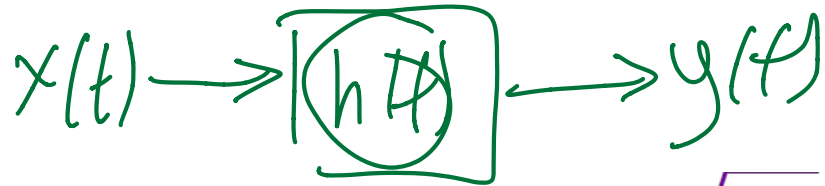
- We do not need to use the squared error, provided we are willing to do more computation.
- Other powers of the weights can be used.



# Transfer function reconstruction for outdoor sound field control

Diego Caviedes Nozal

27/03/2019

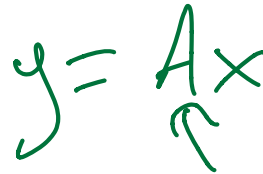


$$y(t) = x(t) * h(t)$$

A collage of mathematical symbols and expressions. It includes a Taylor series expansion:  $f(x+\Delta x) = \sum_{i=0}^{\infty} \frac{(\Delta x)^i}{i!} f^{(i)}(x)$ . Other symbols include a definite integral  $\int_a^b \epsilon \Theta$ , a constant  $\sqrt{17}$ , a summation  $\sum$ , a factorial  $!$ , and a complex number  $\{2.7182818284\}$ .

## Related to the course

Linear models for regression (Lecture 4)

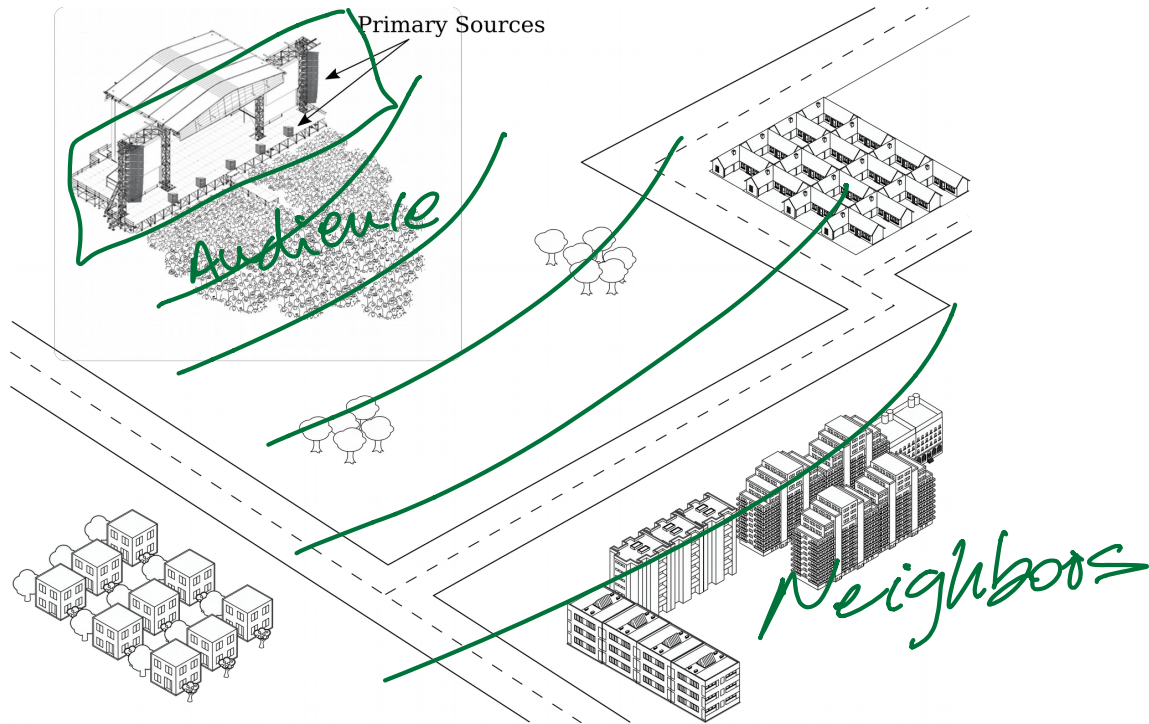
$$y = Ax$$


Non-linear models for regression

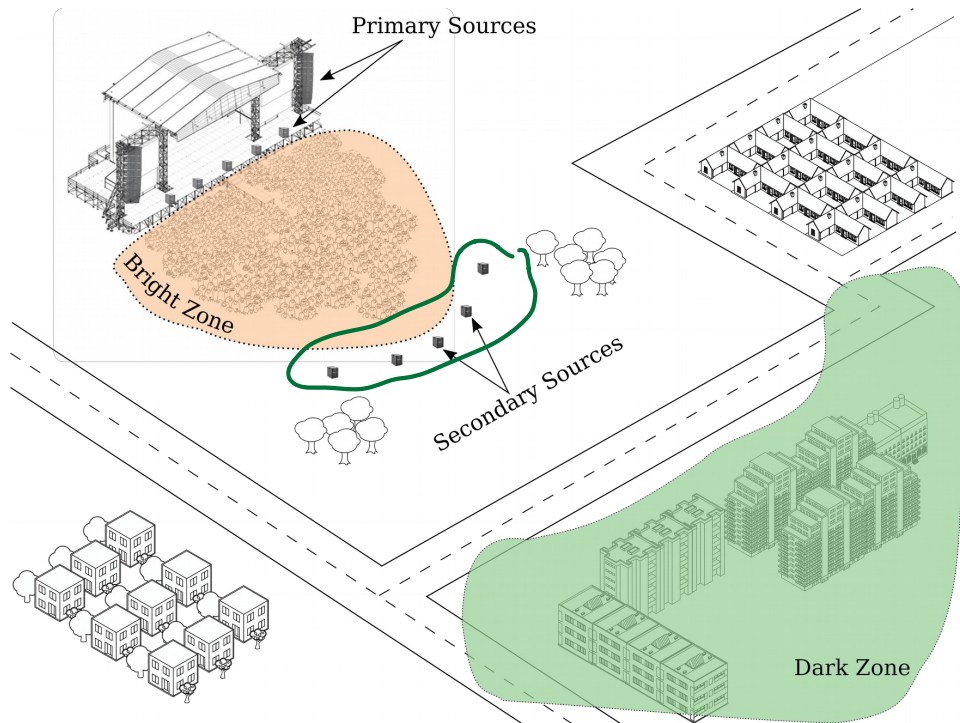
Gaussian Processes (Lecture 3)

Bayes rule (Lecture 1)

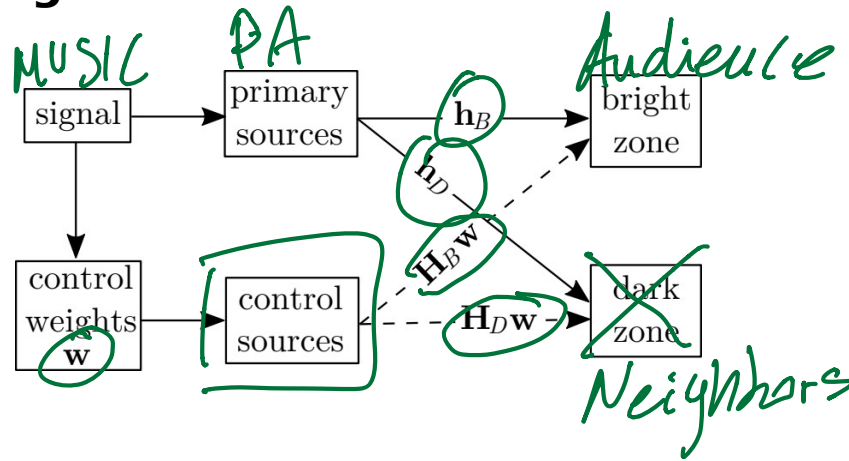
# The problem



# Our goal



# Sound zoning

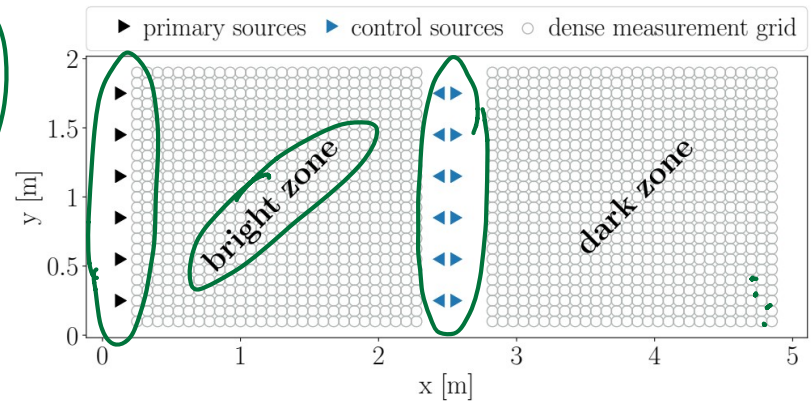
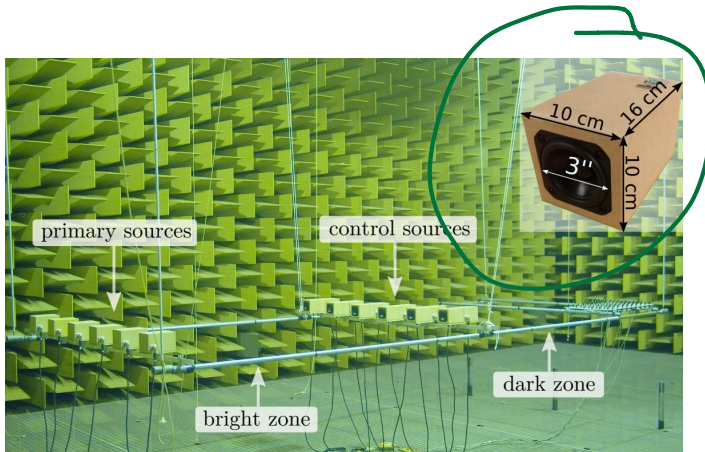


## Objectives

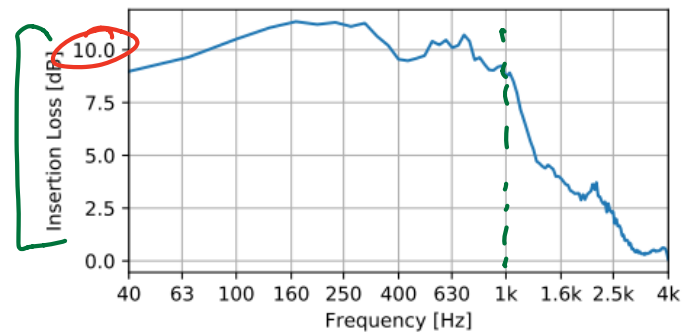
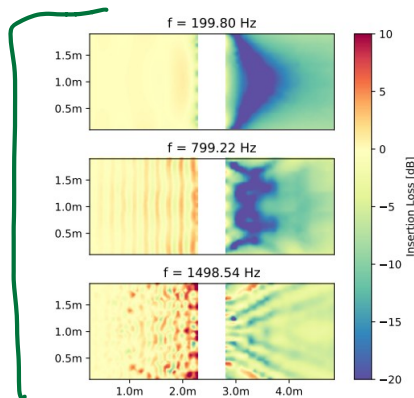
- 1) Cancellation of sound from the primary sources in a *dark zone* using a set of secondary control sources.
- 2) Minimization of the sound radiated by the control sources into the *bright zone*.

$$\text{minimize}_{\mathbf{w}} \quad \underbrace{\kappa}_{2)} \underbrace{\|\mathbf{H}_B \mathbf{w}\|_2^2}_{1)} + (1 - \kappa) \underbrace{\|\mathbf{H}_D \mathbf{w} + \mathbf{h}_D\|_2^2}_{1)} \quad \underbrace{\kappa \in [0, 1]}_{1)}$$

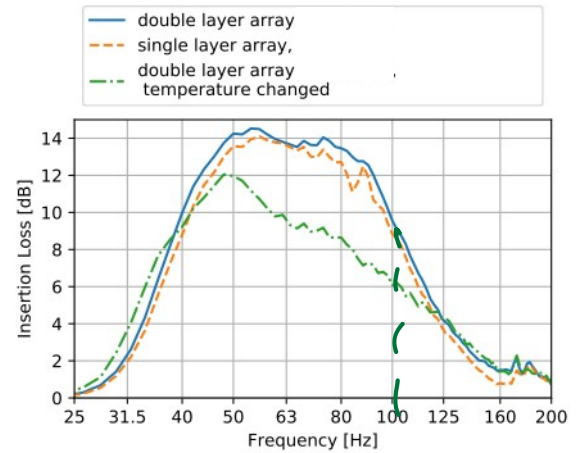
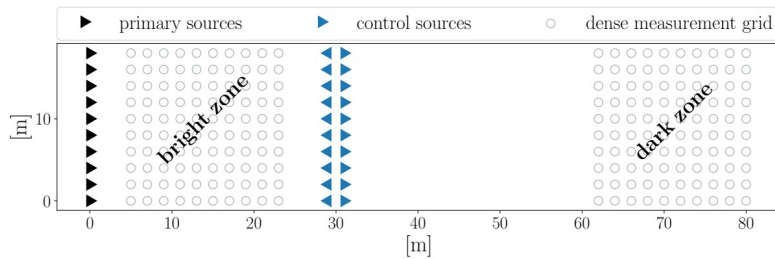
# Real experiments: Anechoic conditions



$$IL = 10 \log_{10} \left( \frac{\|h_{D1}\|^2}{\|H_{D1} + h_{D1}\|^2} \right)$$



# Real experiments: Outdoor conditions





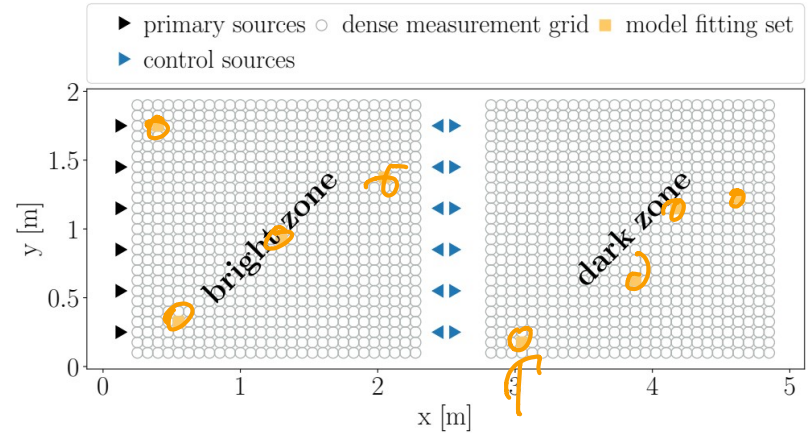
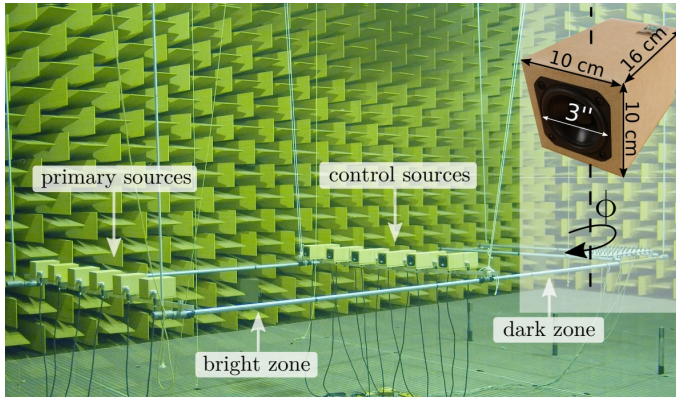
# Measuring transfer functions: Too many issues

- Measuring hundreds of transfer functions to sample the control zones is not possible in real open air concerts.
- The acoustic transfer functions must be representative of the conditions that the sound field control is applied.

Different approach: **Sound propagation models to estimate the transfer functions.**

- Use sparse measurements to fit the model.

# Modeling: Anechoic conditions

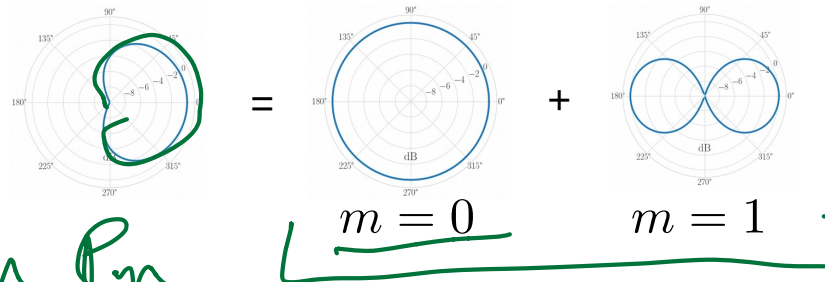


## Source model: spherical harmonics

$$\hat{h}(k, \mathbf{r}) = \sum_{m=0}^{M-1} a_m \eta_m^{(2)}(kr) P_m(\cos(\phi))$$

$$y = W^T \phi(x)$$

$$h = a^T S \quad S = h_m P_m$$



# Modeling: Anechoic conditions

Considering no ~~sensor mismatch~~ neither between mics nor loudspeakers...

...the recorded transfer functions at a single frequency between  $N_L$  sources and  $N_M$  positions

$$\mathbf{h} = \hat{\mathbf{h}} + \mathbf{n}$$

where

$$\hat{\mathbf{h}} = \mathbf{S}\mathbf{a}$$

with  $\hat{\mathbf{h}} \in \mathbb{C}^{N_L N_M}$ ,  $\mathbf{a} \in \mathbb{C}^M$  and  $\mathbf{S} \in \mathbb{C}^{N_L N_M \times M}$  with elements

$$s_{mi} = h_m(kr_i)P_m(\cos(\phi_i))$$

# Modeling: Anechoic conditions

How do we find  $\mathbf{a}$ ?  
Bayesian Inference

$$\pi(\mathbf{a} | \mathbf{h}) \propto \pi(\mathbf{h} | \mathbf{a})\pi(\mathbf{a})$$

Where priors

$$\mathbf{n} \sim \mathcal{CN}(0, \tau^{-1}\mathbf{I})$$

$$\mathbf{a} \sim \mathcal{CN}(0, \delta^{-1}\mathbf{I})$$

$$\tau \sim \mathcal{G}(\alpha, \beta), \delta \sim \mathcal{G}(\alpha, \beta)$$

Likelihood

$$\pi(\mathbf{h} | \mathbf{a}, \tau) \sim \mathcal{CN}(\hat{\mathbf{h}}, \tau^{-1}\mathbf{I}) \propto \exp(-\tau^2 \|\mathbf{S}\mathbf{a} - \mathbf{h}\|^2)$$

And posterior

$$\pi(\mathbf{a}, \tau, \delta | \mathbf{h}) \propto \pi(\mathbf{h} | \mathbf{a}, \tau)\pi(\mathbf{a} | \delta)\pi(\tau)\pi(\delta)$$

$$(\mathbf{a}, \tau, \delta)_{\text{MAP}} = \underset{\mathbf{a}, \tau, \delta}{\text{argmax}} \pi(\mathbf{a}, \tau, \delta | \mathbf{h})$$

$$\hat{h}(k, \mathbf{r}_*) = \mathbf{s}_* \mathbf{a}_{\text{MAP}}^T$$

# Bayesian Languages (e.g. STAN)

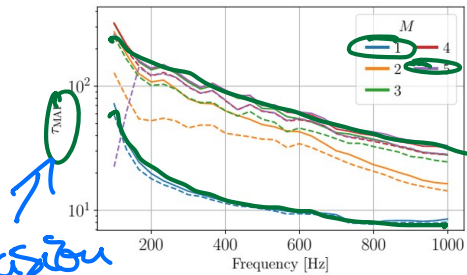
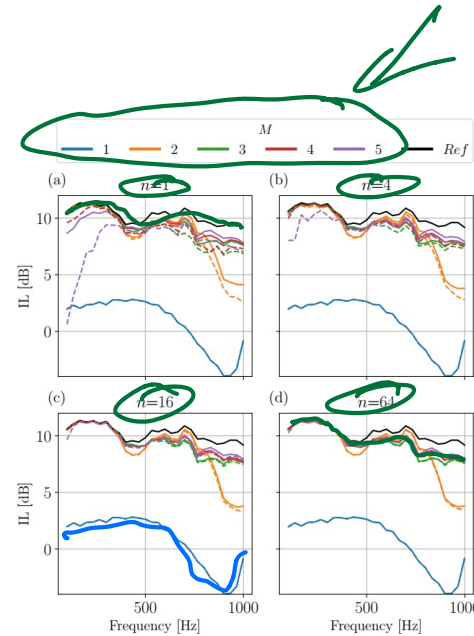
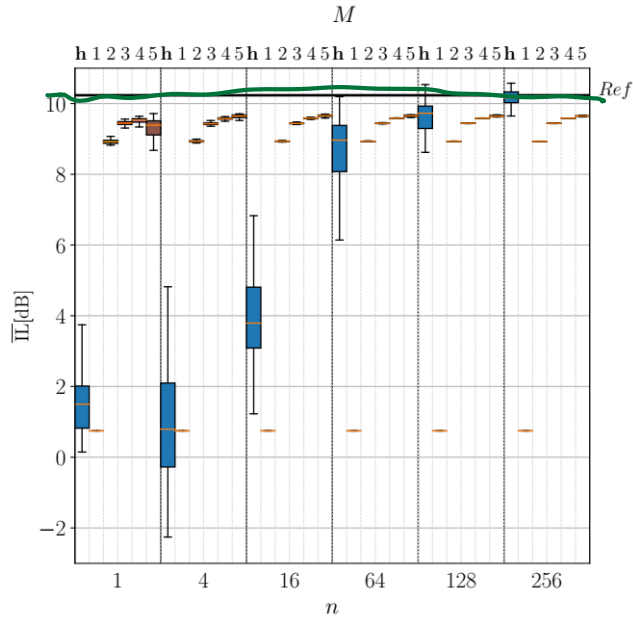
```

1 // ----- Spherical Harmonics -----
2 data {
3   int<lower=0> N_meas;           // Total Number of independent measurements.
4   int<lower=0> M;               // Number of spherical harmonics modes.
5   vector[N_meas] d;           // Distance loudspeaker-measurement.
6   vector[N_meas] pt_real;     // Measured Pressure at receiver. Real part
7   vector[N_meas] pt_imag;     // Measured Pressure at receiver. Imaginary part
8   vector[N_meas] legendre[M]; // Legendre polynomials
9   vector[N_meas] bessel[M];   // Spherical bessel functions
10  vector[N_meas] neumann[M];   // Spherical Neumann functions
11  real a;                       // Prior Hyperparameters
12  real b;                       // Prior Hyperparameters
13 }
14 parameters {
15   vector[M] Ar;                // real part of the source strength
16   vector[M] Ai;                // imaginary part of the source strength
17   real<lower=0> tau;
18   real<lower=0> delta;
19 }
20 transformed parameters{
21   vector[N_meas] mu_real;      // mean of the real part of the pressure
22   vector[N_meas] mu_imag;     // mean of the imaginary part of the pressure
23   real<lower=0> inv_delta;
24   real<lower=0> inv_tau;
25   inv_tau = 1/tau;
26   inv_delta = 1/delta;
27   // Loop over modes
28   mu_real = 0 * d;
29   mu_imag = 0 * d;
30   for (m in 1:M){
31     mu_real += d .* (Ar[m] * bessel[m] + Ai[m] * neumann[m]) .* legendre[m];
32     mu_imag += d .* (Ai[m] * bessel[m] - Ar[m] * neumann[m]) .* legendre[m];
33   }
34 }
35 model {
36   tau ~ gamma(a, b);
37   delta ~ gamma(a, b);
38   Ar ~ normal(0, inv_delta);
39   Ai ~ normal(0, inv_delta);
40   pt_real ~ normal(mu_real, inv_tau);
41   pt_imag ~ normal(mu_imag, inv_tau);
42 }

```

# Modeling: Anechoic conditions

Insertion loss

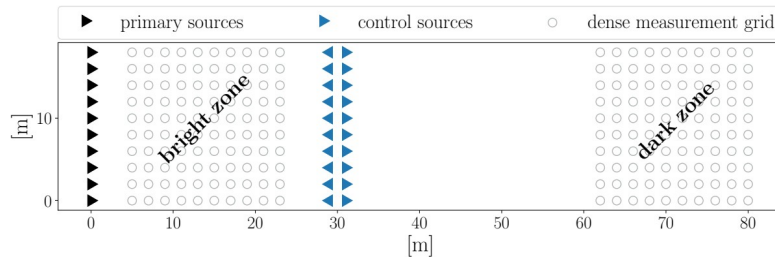


precision  
T

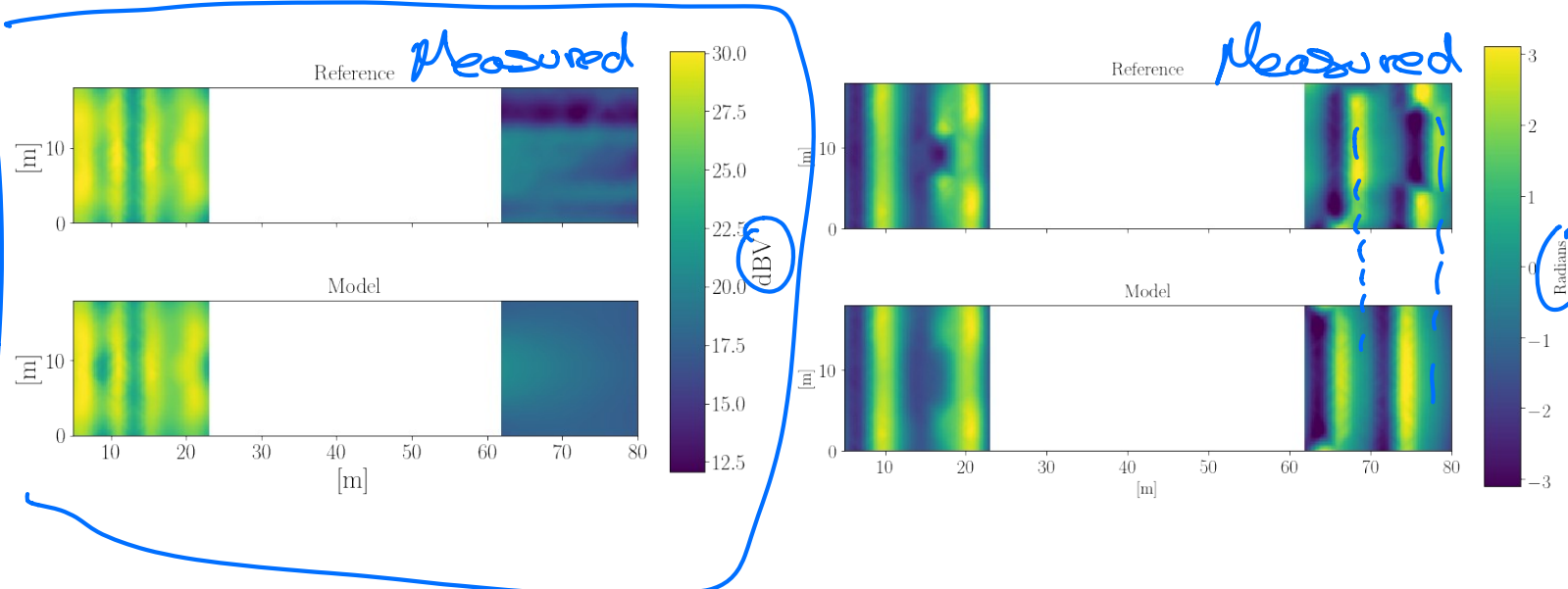
# Modeling: Outdoor conditions



- Complex scenario (geometry).
- Complex medium (refraction, turbulences...).



# Modeling: Outdoor conditions





# Semiparametric GP

Reformulate the model (a little bit)

$$\mathbf{h} = \hat{\mathbf{h}} + \mathbf{d} + \mathbf{n}$$

Physics  $\rightarrow$

$$\mathbf{h} \sim \mathcal{N}(\hat{\mathbf{h}}, \mathbf{K}_n)$$

Where

$$\mathbf{d} \sim \mathcal{GP}(\mathbf{0}, \mathbf{K})$$

Data  $\rightarrow$

$$\mathbf{K}_n = \mathbf{K} + (\mathbf{I}^{-1})^2 \mathbf{I}$$

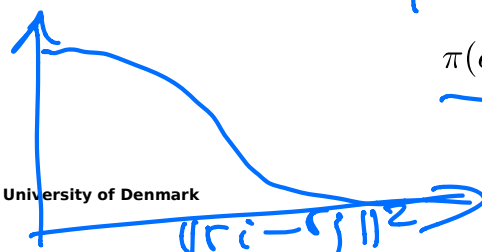
It introduces flexibility to the covariance matrix.

The main issue is to define a kernel that makes sense for the problem (probably spatially periodic?).

We started with something easy: radial basis function

$$\mathbf{K} = \kappa(\mathbf{R}, \mathbf{R}) = \begin{bmatrix} \kappa(\mathbf{r}_1, \mathbf{r}_1) & \dots & \kappa(\mathbf{r}_1, \mathbf{r}_n) \\ \vdots & \kappa(\mathbf{r}_i, \mathbf{r}_i) & \vdots \\ \kappa(\mathbf{r}_N, \mathbf{r}_1) & \dots & \kappa(\mathbf{r}_n, \mathbf{r}_n) \end{bmatrix}$$

$$\kappa(\mathbf{r}_i, \mathbf{r}_j) = \alpha^2 \exp\left(-\frac{1}{2\rho^2} \|\mathbf{r}_i - \mathbf{r}_j\|^2\right)$$



$$\pi(\alpha) \sim \mathcal{N}(0, \sigma_\alpha^2), \quad \pi(\rho) \sim \mathcal{N}(0, \sigma_\rho^2)$$

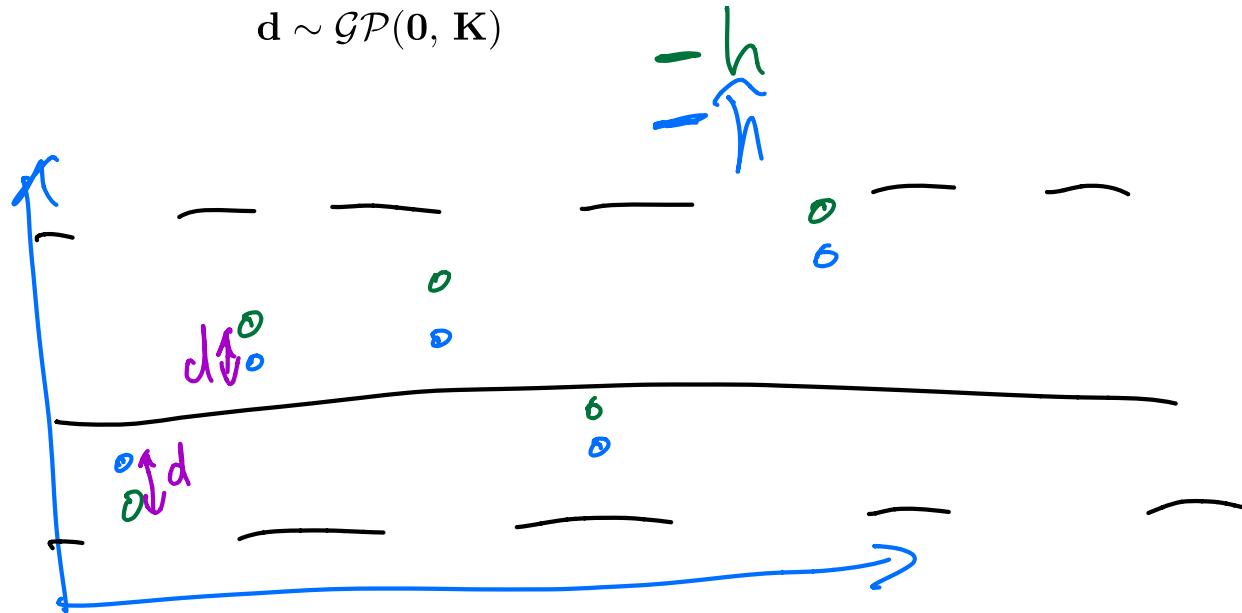
# Semiparametric GP

Reformulate the model (a little bit)

$$\mathbf{h} = \hat{\mathbf{h}} + \mathbf{d} + \mathbf{n}$$

Where

$$\mathbf{d} \sim \mathcal{GP}(\mathbf{0}, \mathbf{K})$$



# Transfer function prediction

How do we predict elsewhere? A bit more complicated than before

$$\pi(\mathbf{a}, \tau, \delta | \mathbf{h}) \propto \pi(\mathbf{h} | \mathbf{a}, \tau) \pi(\mathbf{a} | \delta) \pi(\tau) \pi(\delta)$$

$$(\mathbf{a}, \tau, \delta)_{\text{MAP}} = \underset{\mathbf{a}, \tau, \delta}{\text{argmax}} \pi(\mathbf{a}, \tau, \delta | \mathbf{h})$$

$$\hat{h}_{\text{MAP}}(k, \mathbf{r}_*) = \mathbf{s}_* \mathbf{a}_{\text{MAP}}^T$$

$$\bar{\mathbf{h}}_* = \hat{\mathbf{h}}_*^{\text{MAP}} + (\mathbf{K}_*^{\text{MAP}})^T (\mathbf{K}_n^{\text{MAP}})^{-1} (\mathbf{h} - \hat{\mathbf{h}}_{\text{MAP}})$$

$$\text{COV}(\mathbf{h}_*) = \mathbf{K}_{**}^{\text{MAP}} - (\mathbf{K}_*^{\text{MAP}})^T (\mathbf{K}_n^{\text{MAP}})^{-1} \mathbf{K}_*^{\text{MAP}}$$

$$\mathbf{K}_n^{\text{MAP}} = \mathbf{K}^{\text{MAP}} + \frac{1}{\tau^2} \mathbf{I}$$

# Transfer function prediction

