

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

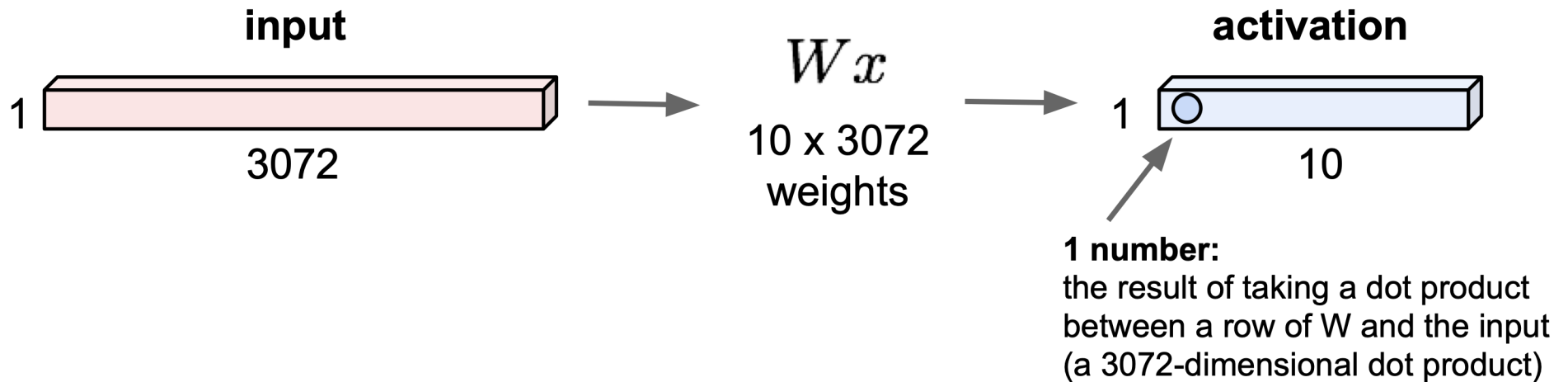
Stanford CNN, Linear models for regression/classification (Bishop 3 and 4),

CNN

Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1

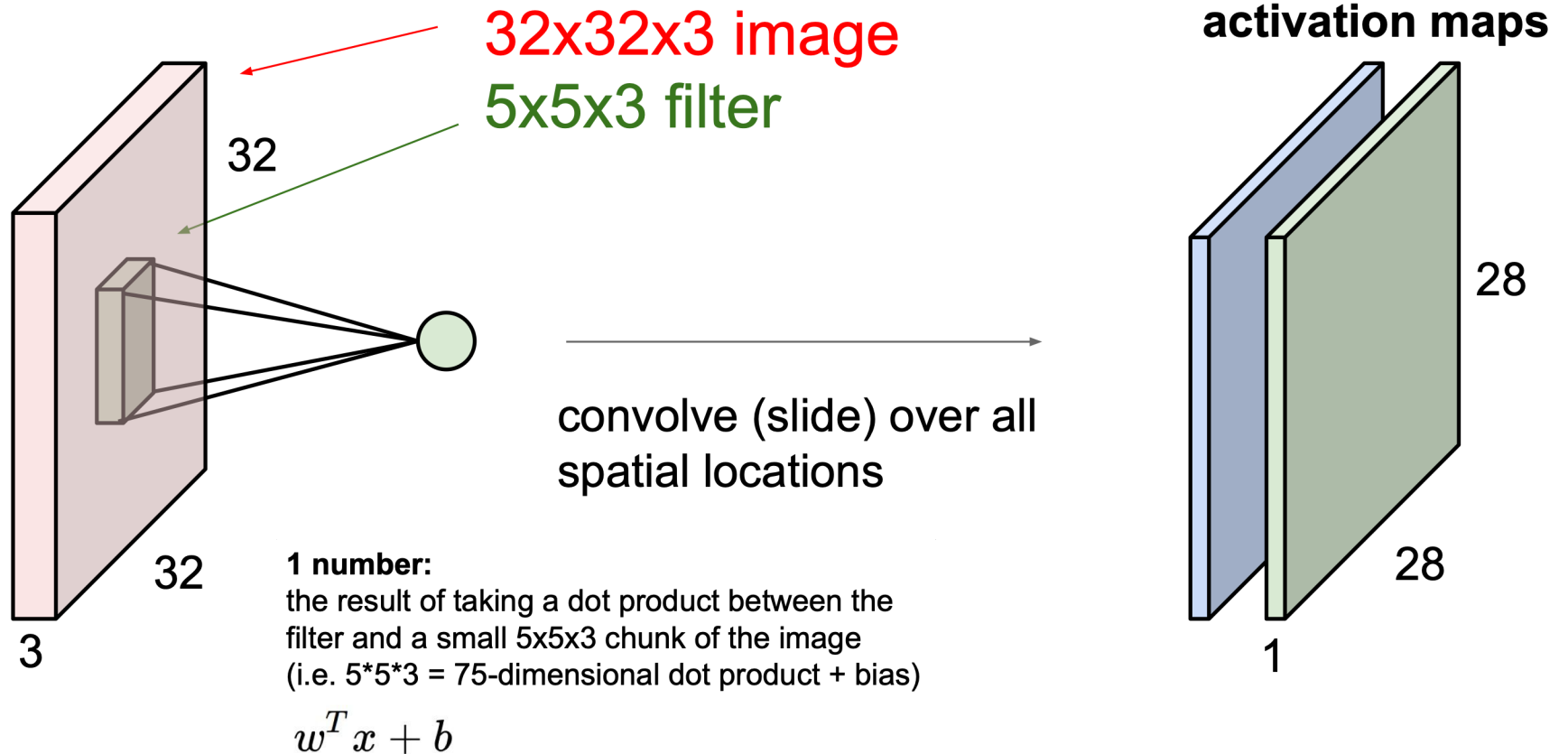
Each neuron
looks at the full
input volume



CNN

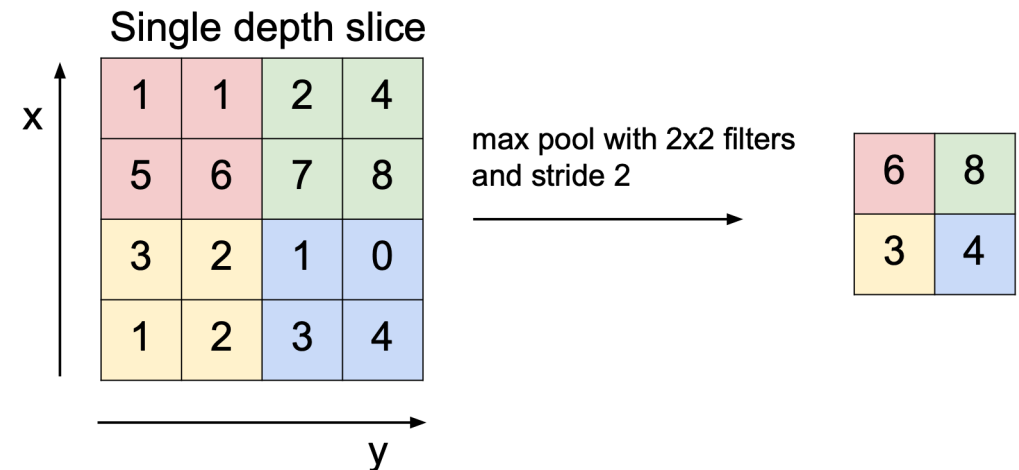
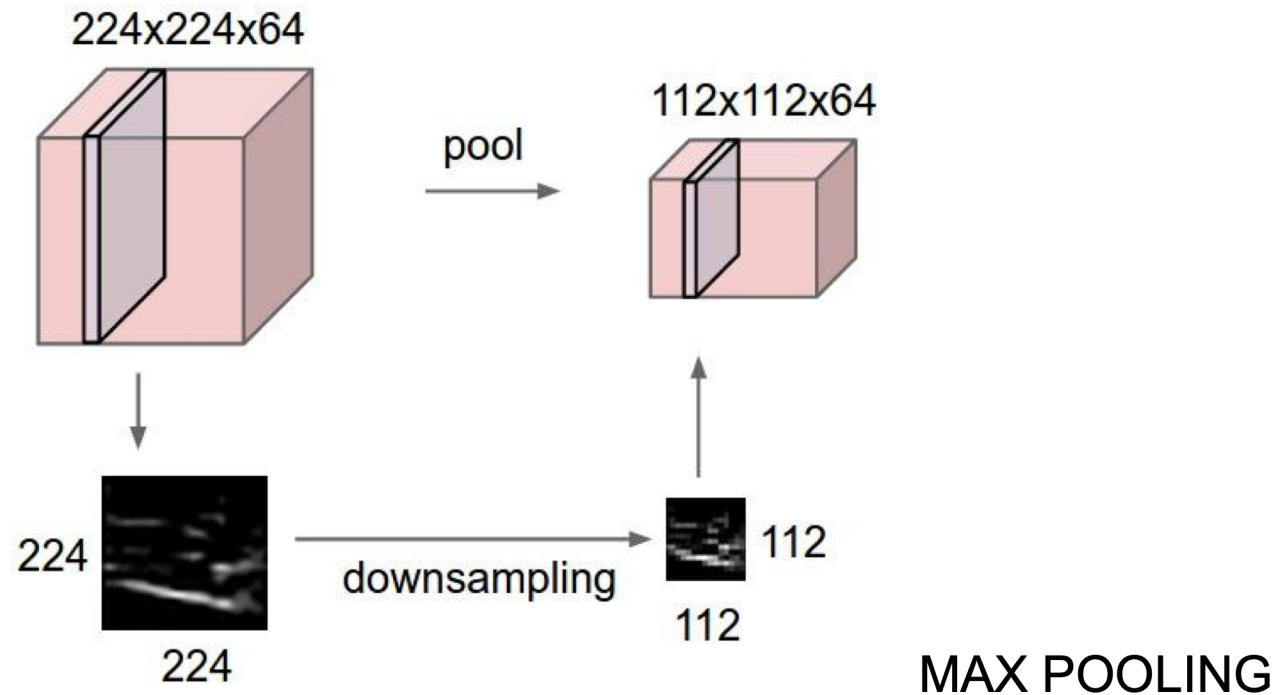
Convolution Layer

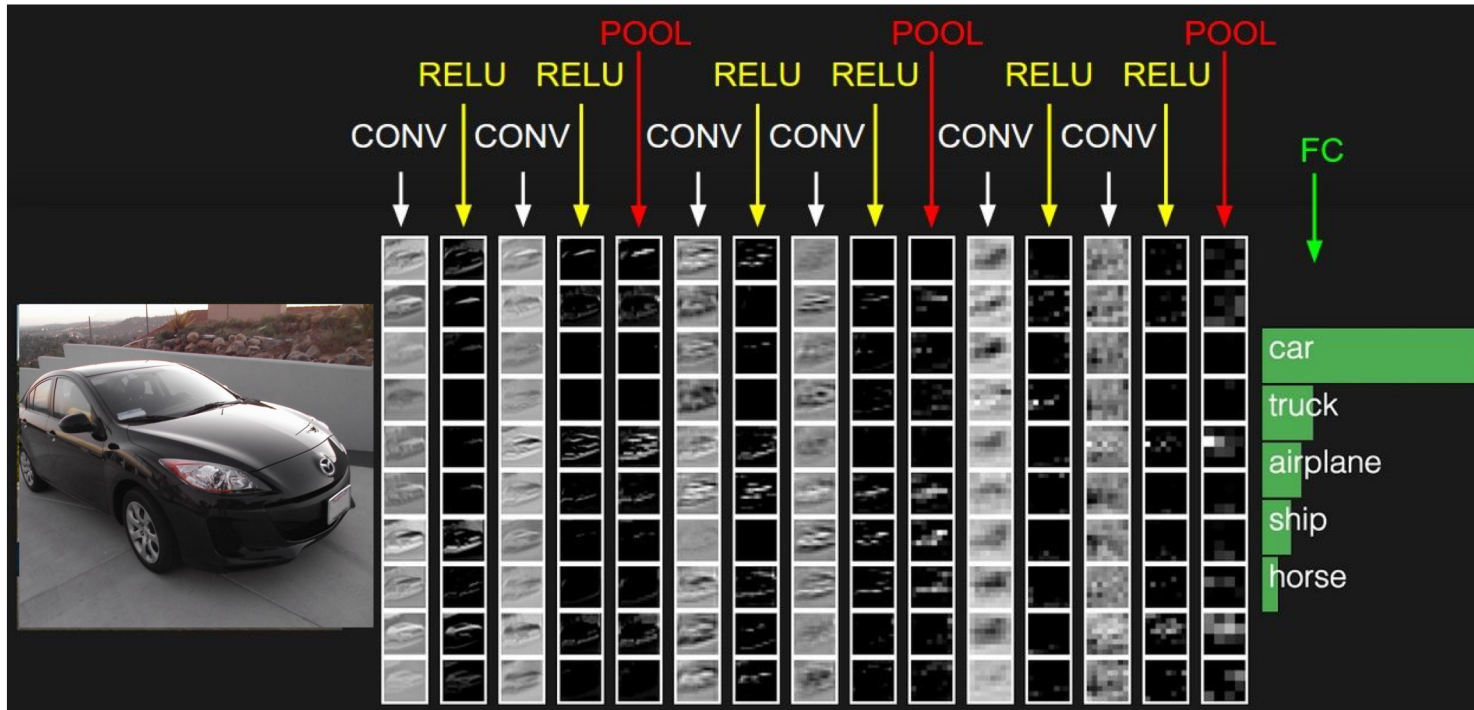
consider a second, **green** filter



Pooling layer

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





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 **$[(\text{CONV-RELU})^N\text{-POOL?}]^M\text{-(FC-RELU)}^K, \text{SOFTMAX}$** 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

Transfer function reconstruction for outdoor sound field control

Diego Caviedes Nozal

27/03/2019

$$f(x+\Delta x) = \sum_{i=0}^{\infty} \frac{(\Delta x)^i}{i!} f^{(i)}(x)$$

$\int_a^b \epsilon \Theta + \Omega \int \delta e^{i\pi} = \{2.7182818284\}$
 χ^2
 Σ
 ∞
 $\sqrt{17}$
 Δ
 ϵ
 Θ
 Ω
 δ
 $e^{i\pi}$
 $\{2.7182818284\}$
 χ^2
 Σ
 ∞
 $\sqrt{17}$

Related to the course

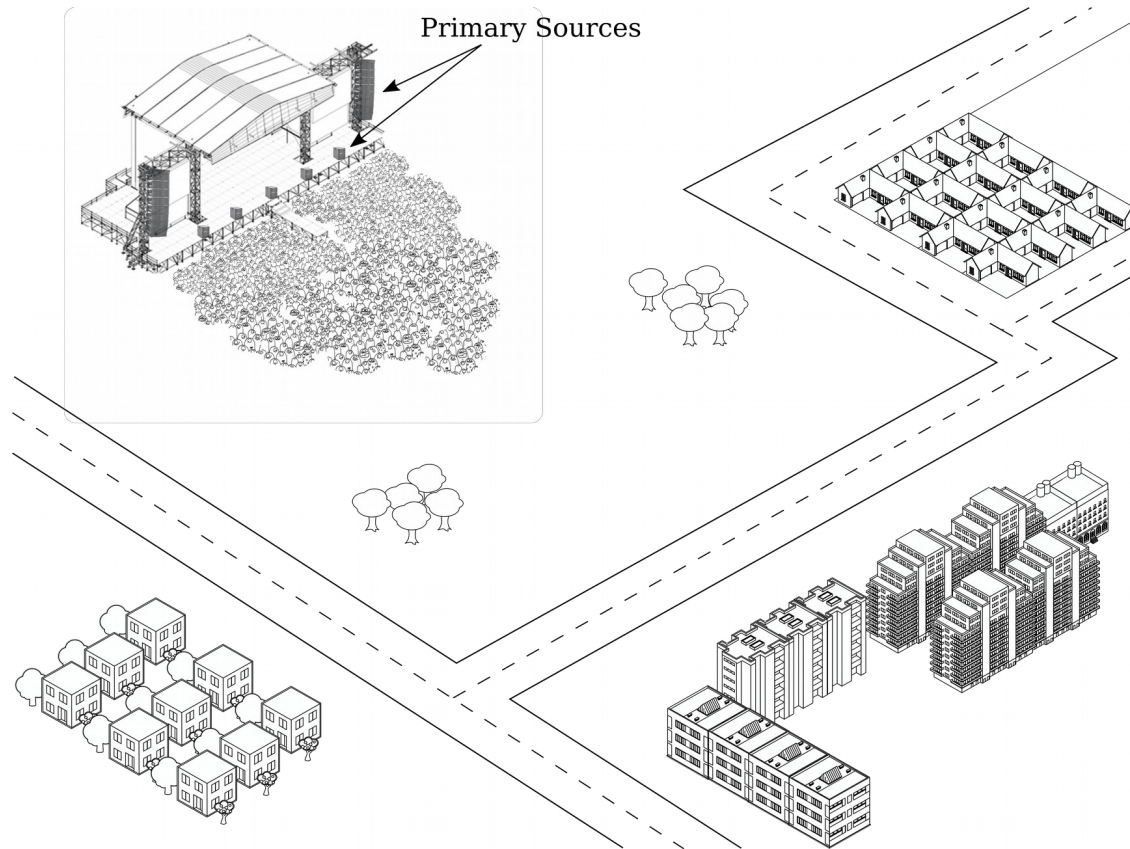
Linear models for regression (Lecture 4)

Non-linear models for regression

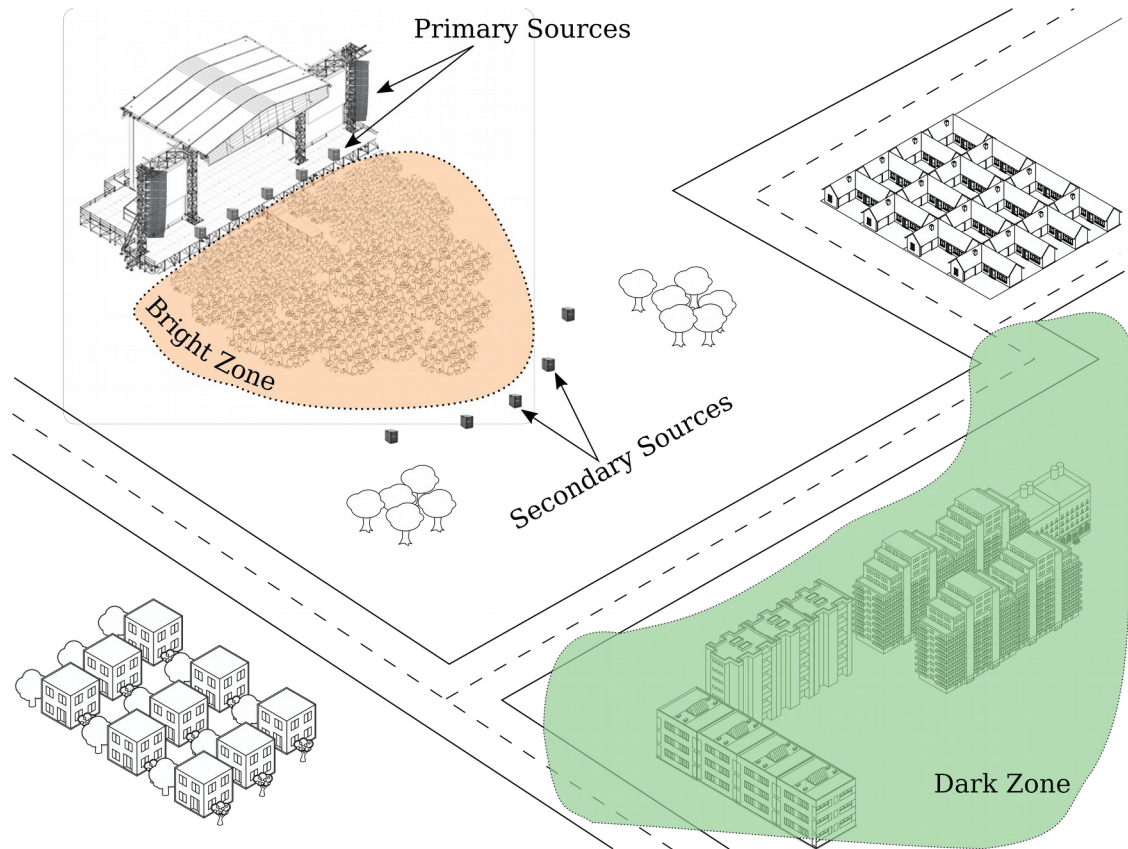
Gaussian Processes (Lecture 3)

Bayes rule (Lecture 1)

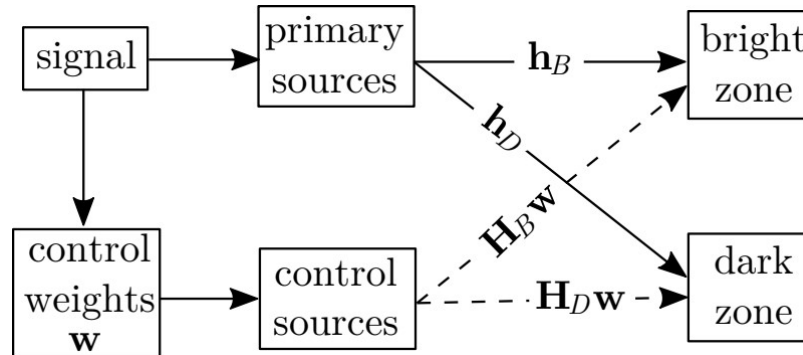
The problem



Our goal



Sound zoning

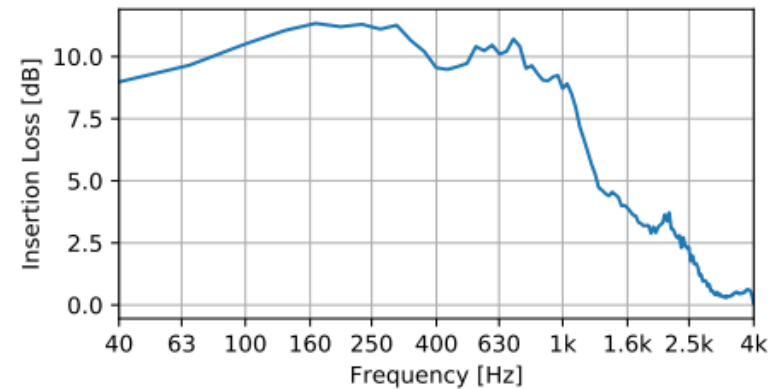
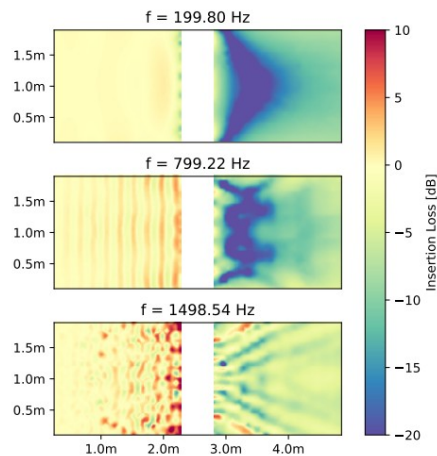
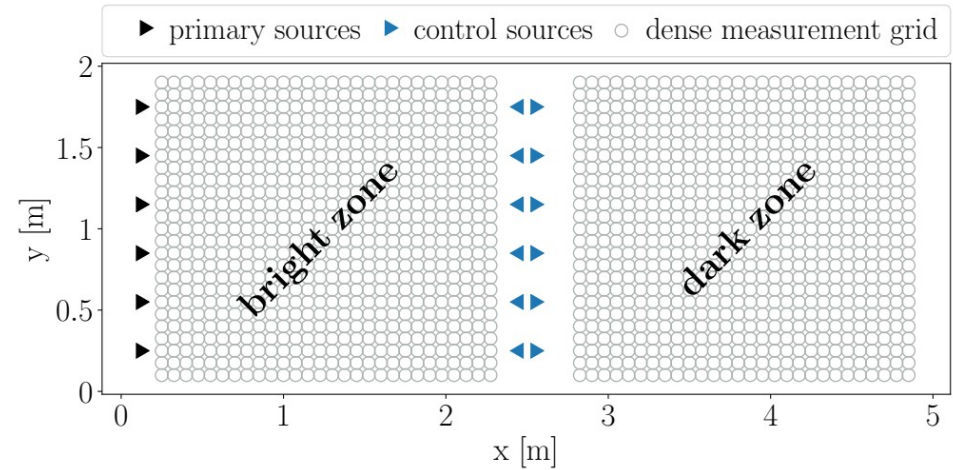
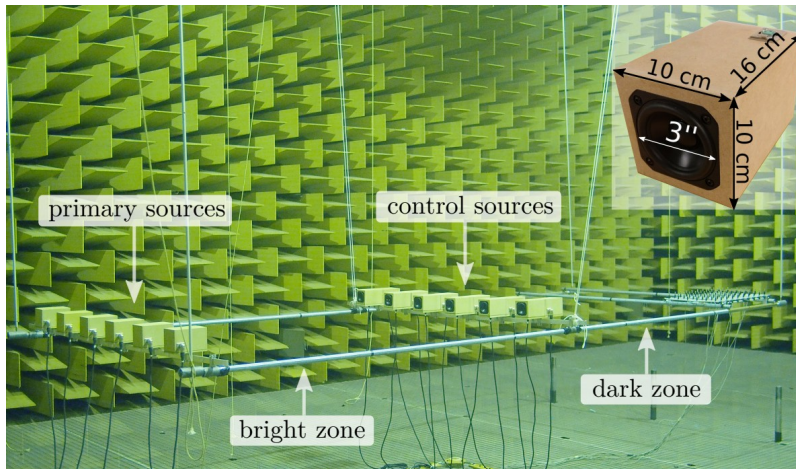


Objectives

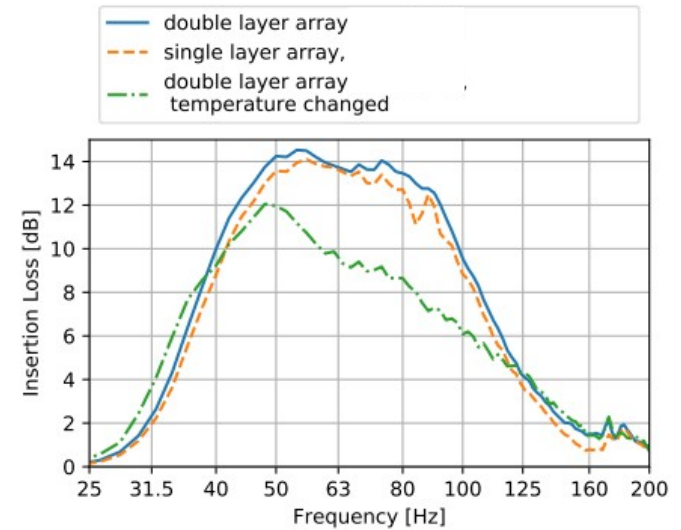
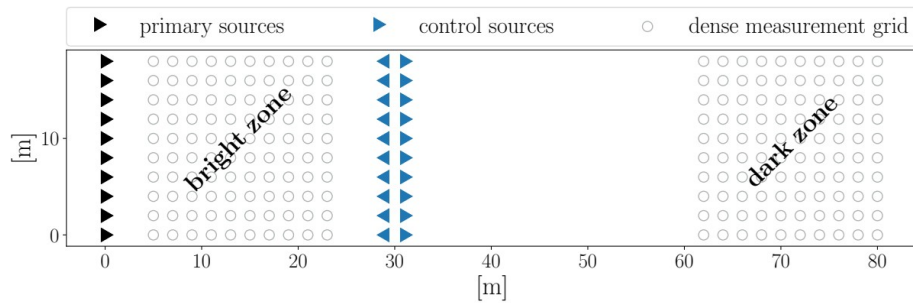
- 1) Cancellation of sound from the primary sources in a *dark zone* using a set of secondary control sources.
- 2) Minimization of the sound radiated by the control sources into the *bright zone*.

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

Real experiments: Anechoic conditions



Real experiments: Outdoor conditions



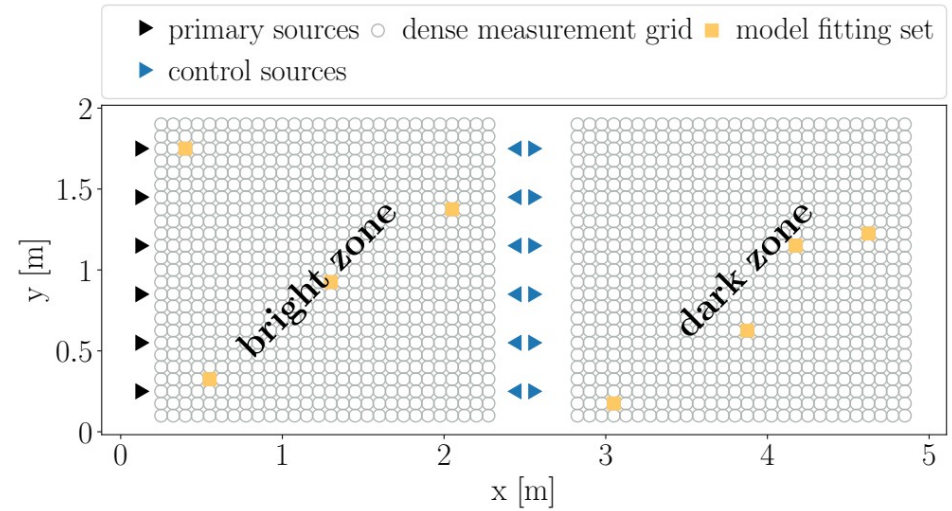
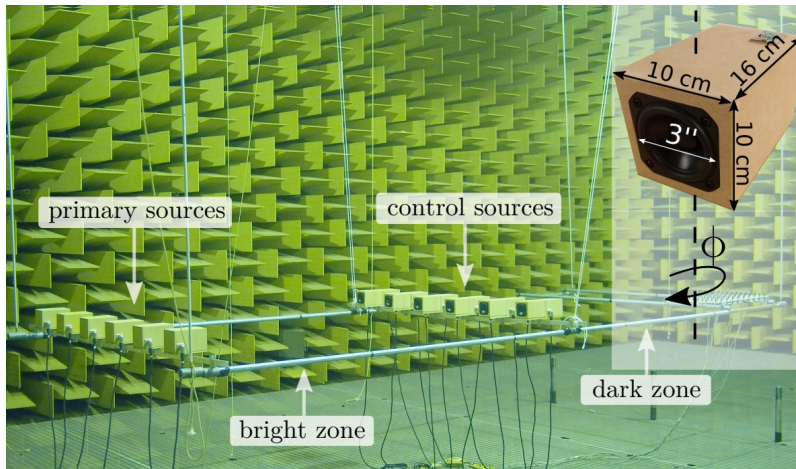
Measuring transfer functions: Too many issues

- 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.**

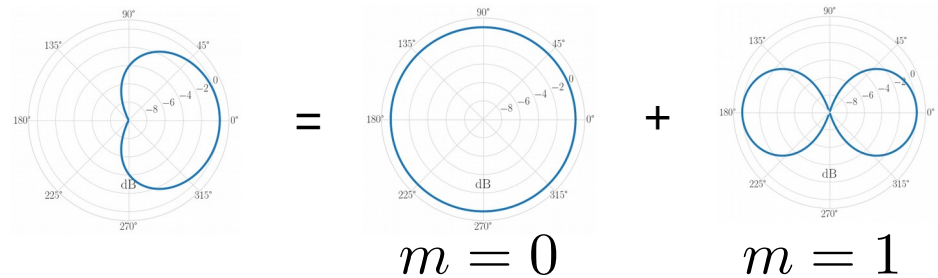
- Use sparse measurements to fit the model.

Modeling: Anechoic conditions



Source model: spherical harmonics

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



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

$$\begin{aligned} \mathbf{n} &\sim \mathcal{CN}(0, \tau^{-1}\mathbf{I}) & \tau &\sim \mathcal{G}(\alpha, \beta), \quad \delta \sim \mathcal{G}(\alpha, \beta) \\ \mathbf{a} &\sim \mathcal{CN}(0, \delta^{-1}\mathbf{I}) \end{aligned}$$

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

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

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

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

Bayesian Languages (e.g. STAN)

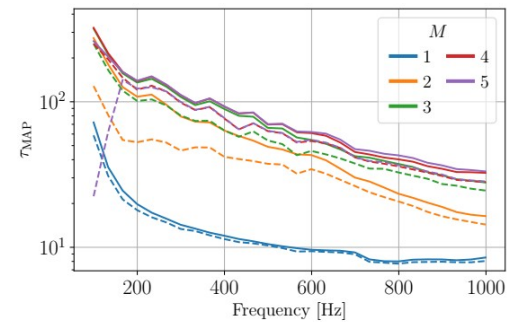
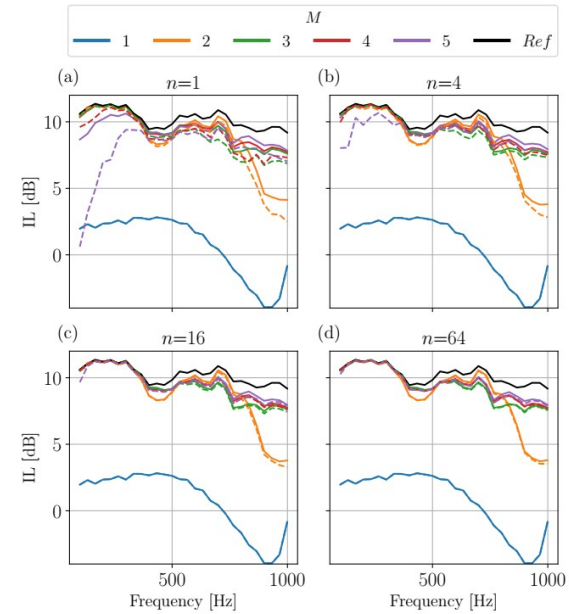
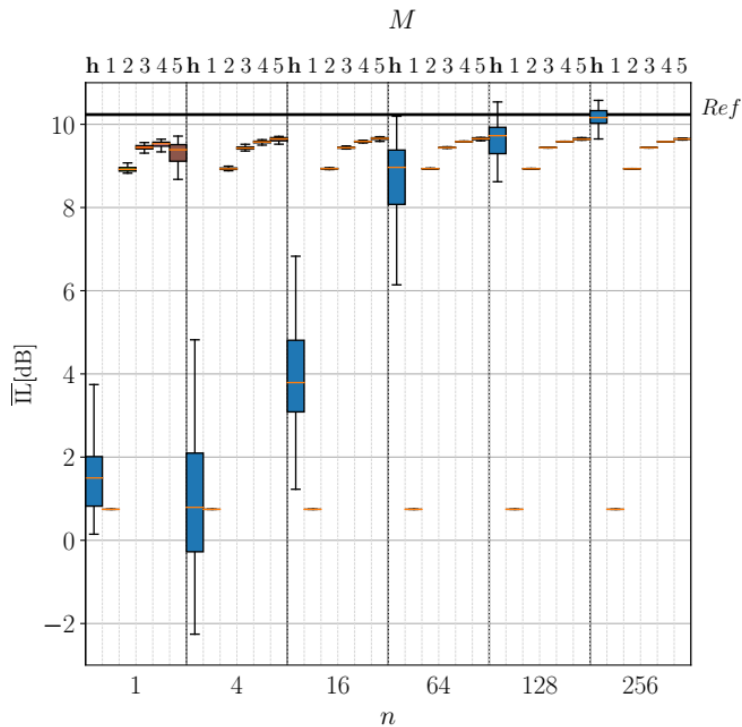
```

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] bessell[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] * bessell[m] + Ai[m] * neumann[m]) .* legendre[m];
32     mu_imag += d .* (Ai[m] * bessell[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

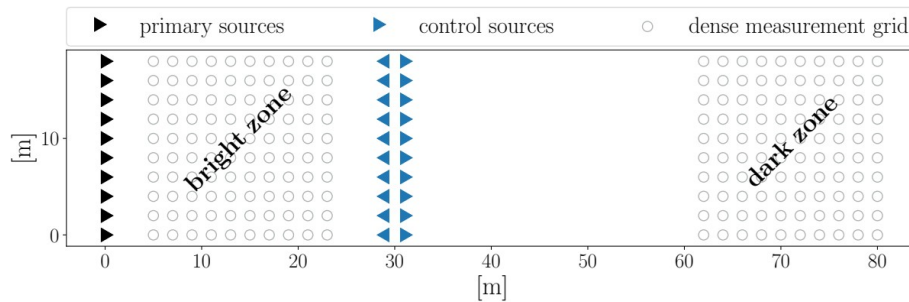
Insertion loss



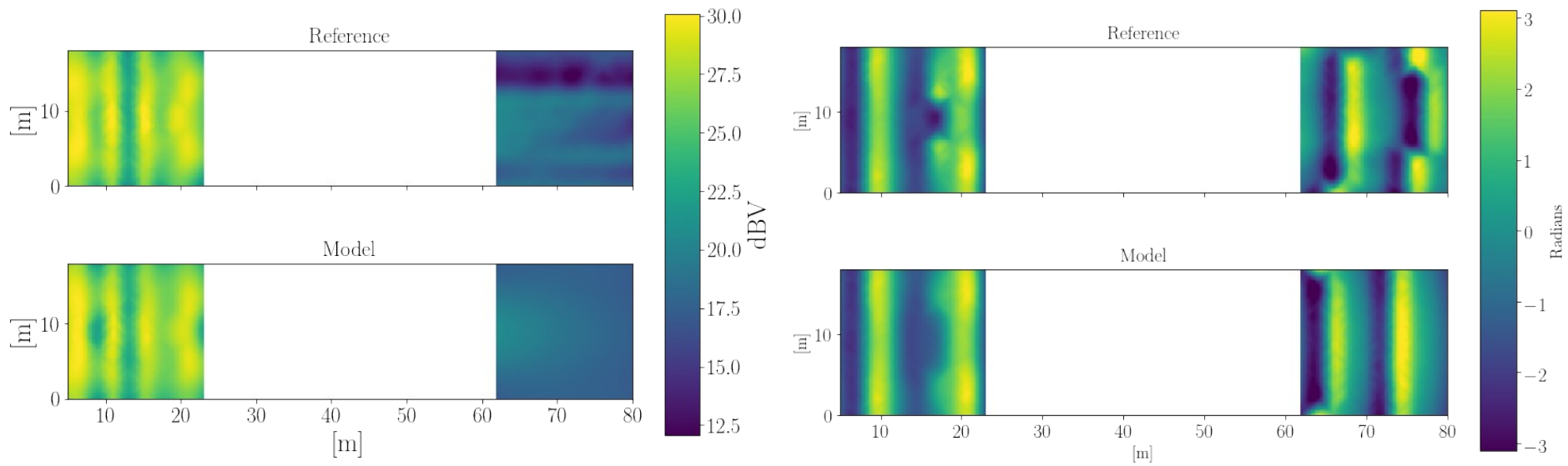
Modeling: Outdoor conditions



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



Modeling: Outdoor conditions



Semiparametric GP

Reformulate the model (a little bit)

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

Where

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

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

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

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

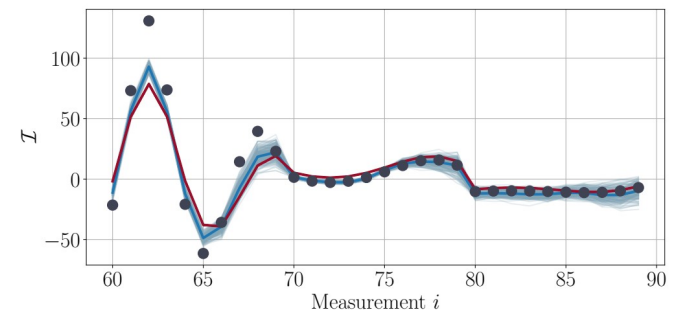
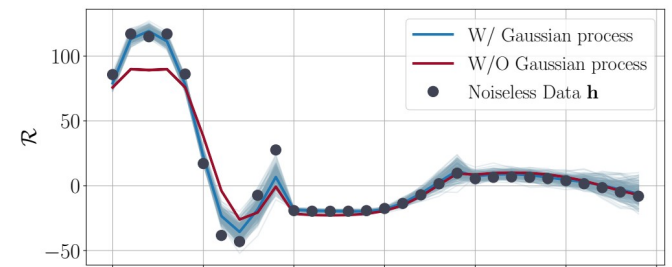
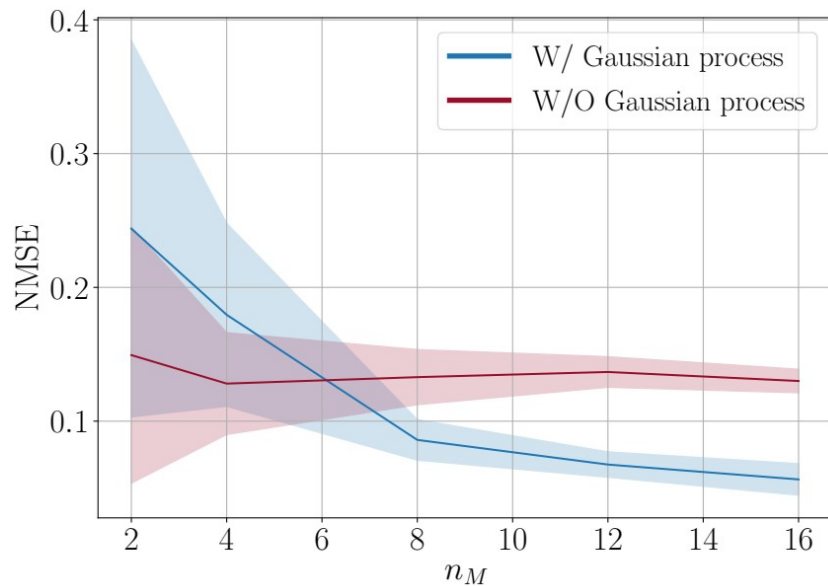
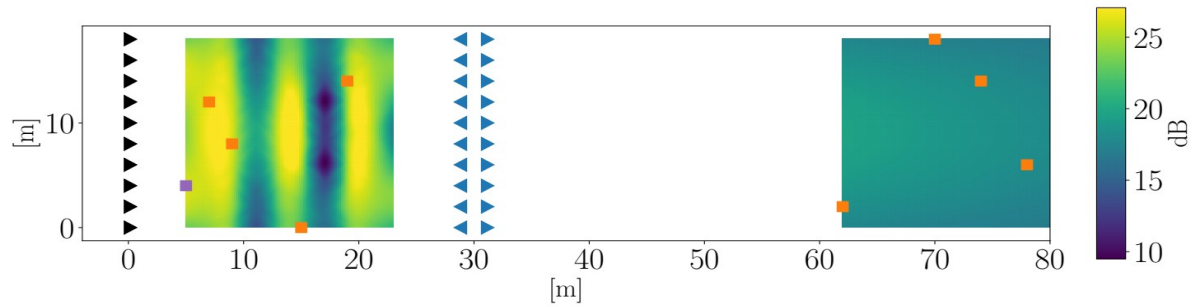
$$\hat{h}_{\text{MAP}}(k, \mathbf{r}_*) = \mathbf{s}_* \mathbf{a}_{\text{MAP}}^{\text{T}}$$

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

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

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

Transfer function prediction

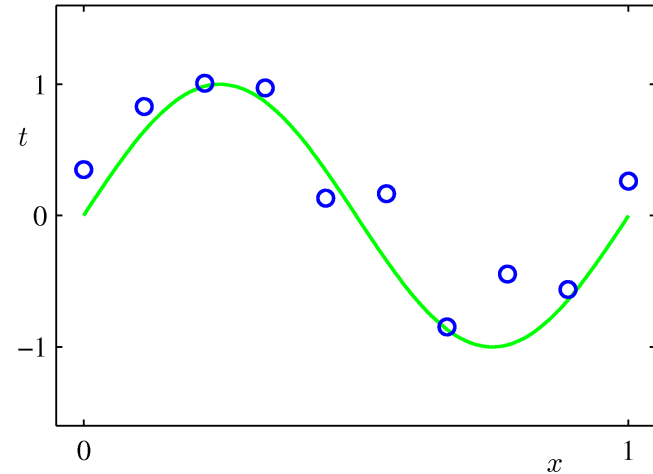


Linear regression: Linear Basis Function Models (1)

Generally

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

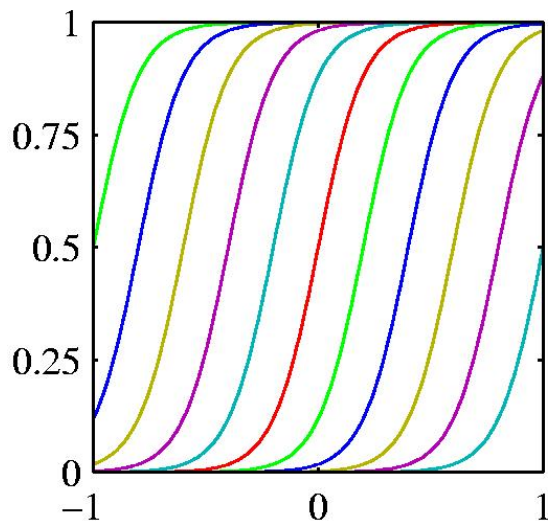
- where $\phi_j(x)$ are known as *basis functions*.
- Typically, $\phi_0(x) = 1$, so that w_0 acts as a bias.
- Simplest case is linear basis functions: $\phi_d(x) = x_d$.



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

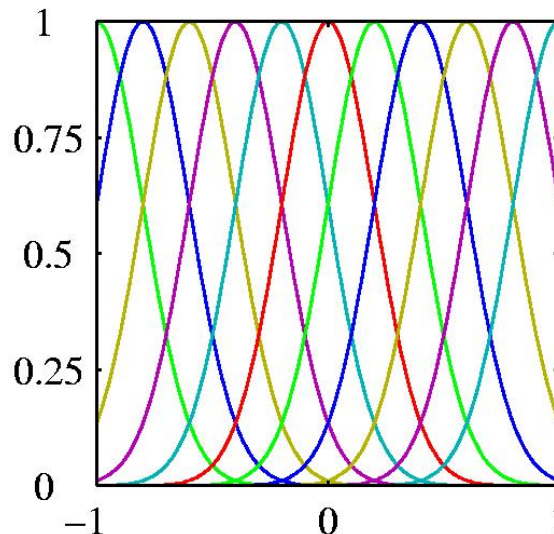
<http://playground.tensorflow.org/>

Some types of basis function in 1-D



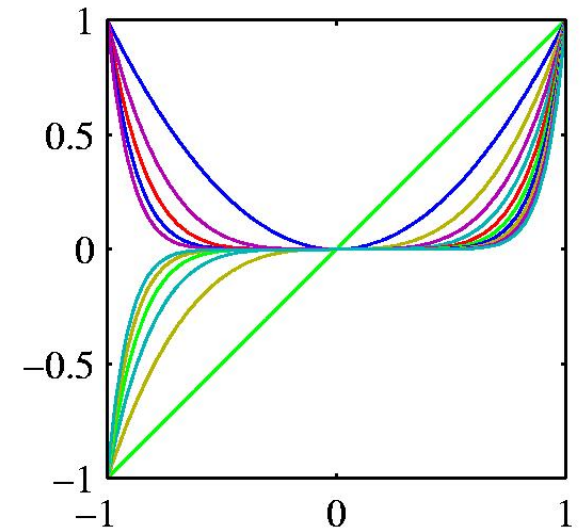
Sigmoids

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



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

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

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

- 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

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

Where the sum-of-squares error is

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2$$

Maximum Likelihood and Least Squares (3)

Computing the gradient and setting it to zero yields

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

Solving for \mathbf{w} ,

where

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

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

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

Maximum Likelihood and Least Squares (4)

Maximizing with respect to the bias, w_0 , alone,

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

We can also maximize with respect to β , giving

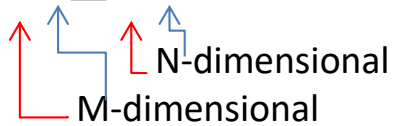
$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{n=1}^N \{t_n - \mathbf{w}_{\text{ML}}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2$$

Geometry of Least Squares

Consider

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

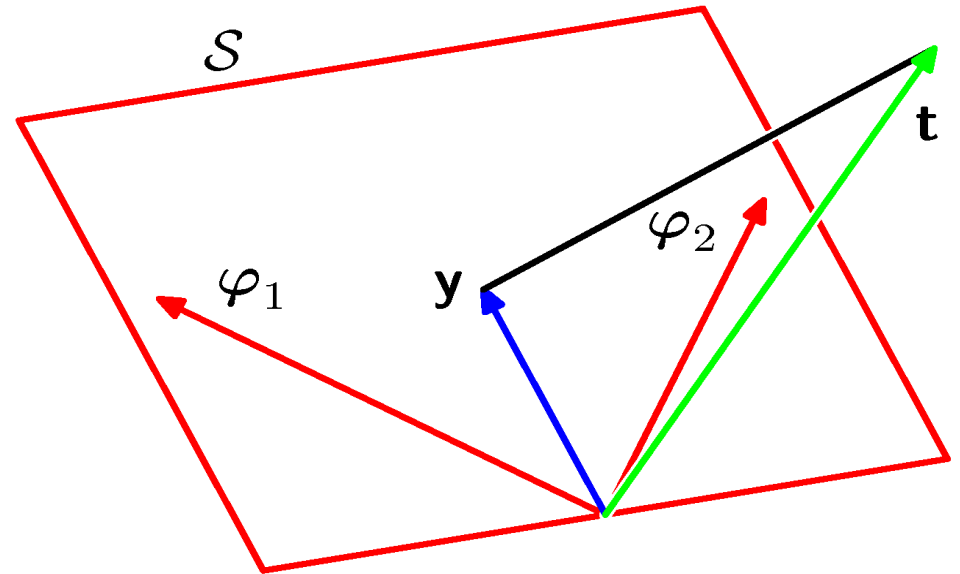
$$\mathbf{y} \in \mathcal{S} \subseteq \mathcal{T} \quad \mathbf{t} \in \mathcal{T}$$



\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} .



Least mean squares: An alternative approach for big datasets

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

↑ weights after seeing training case tau+1

↑ learning rate

↑ 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.

$$\frac{\partial E}{\partial \mathbf{w}} = \sum_{n=1}^N \sigma_n (t_n - \mathbf{w}^T \boldsymbol{\theta}_n)$$

Regularized least squares

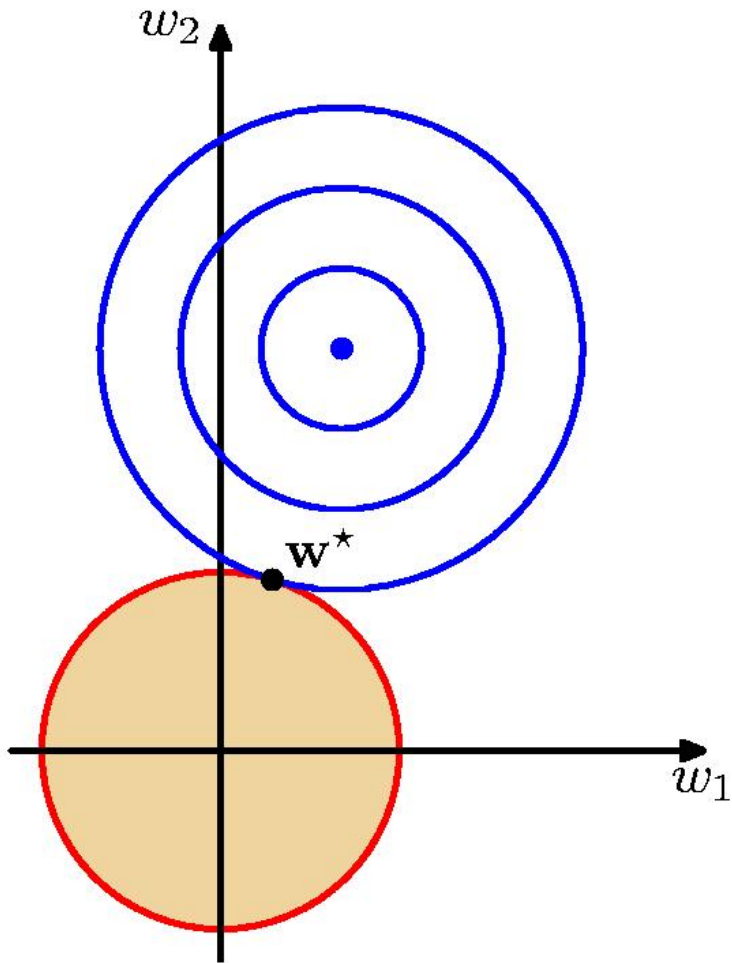
$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

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

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

↑
identity matrix

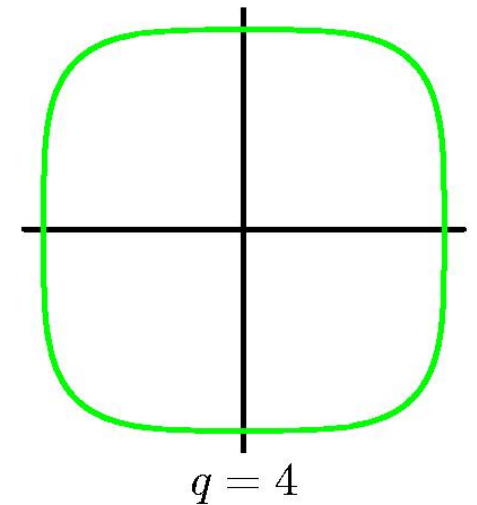
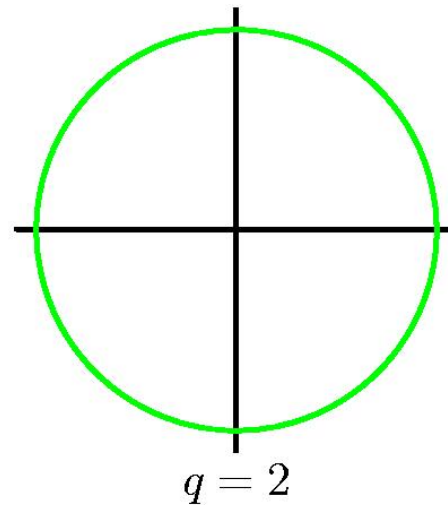
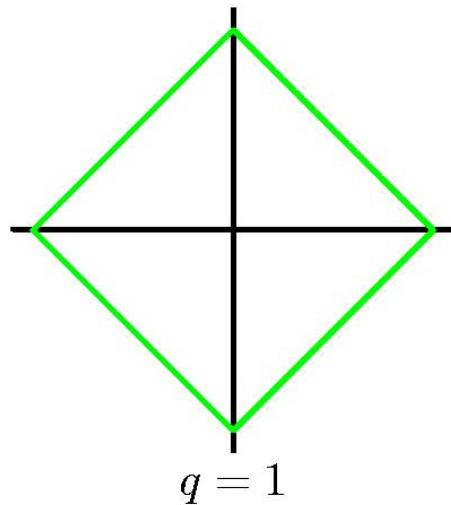
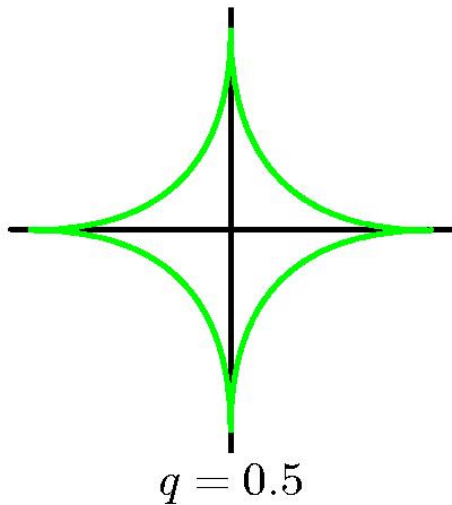
A picture of the effect of the regularizer



- 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.



Minimizing the absolute error

$$\min_{\text{over } \mathbf{w}} \sum_n |t_n - \mathbf{w}^T \mathbf{x}_n|$$

- This minimization involves solving a linear programming problem.
- It corresponds to maximum likelihood estimation if the output noise is modeled by a Laplacian instead of a Gaussian.

$$p(t_n | y_n) = a e^{-a |t_n - y_n|}$$

$$-\log p(t_n | y_n) = -a |t_n - y_n| + \text{const}$$

The bias-variance decomposition

model estimate for
testcase n trained
on dataset D

average target
value for test
case n

“Bias” term is the squared error of the average,
over training datasets D, of the estimates.

Bias: average between prediction and desired.

$$\left\langle \left\{ y(\mathbf{x}_n; D) - \bar{t}_n \right\}^2 \right\rangle_D = \left\langle \left\{ \left\langle y(\mathbf{x}_n; D) \right\rangle_D - \bar{t}_n \right\}^2 \right\rangle_D + \left\langle \left\{ y(\mathbf{x}_n; D) - \left\langle y(\mathbf{x}_n; D) \right\rangle_D \right\}^2 \right\rangle_D$$

$\langle . \rangle$ means
expectation over D

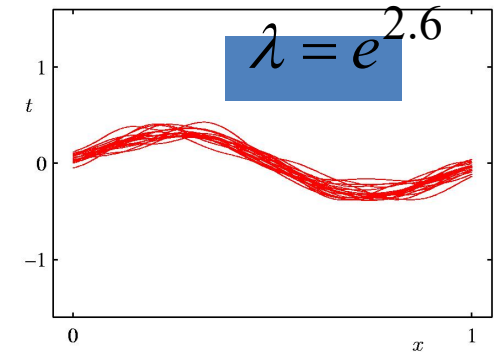
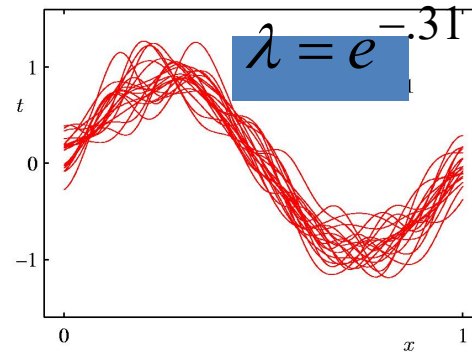
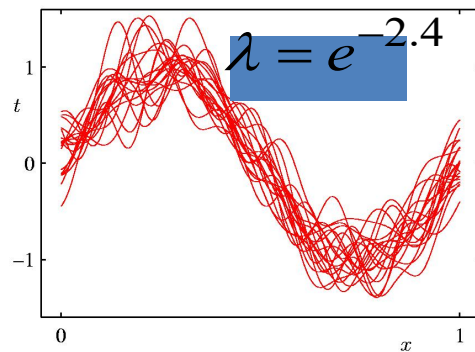
“Variance” term: variance over training datasets D,
of the model estimate.

Regularization parameter affects the bias and variance terms

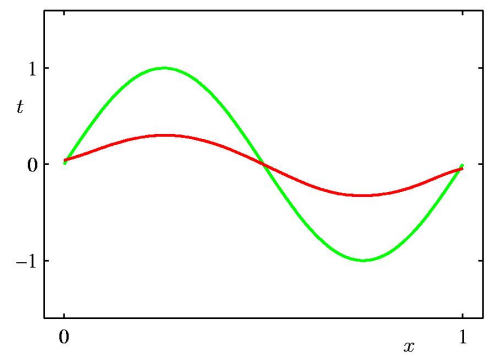
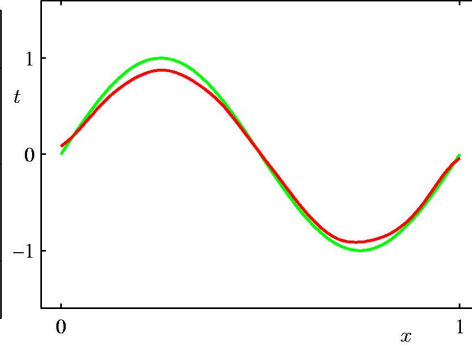
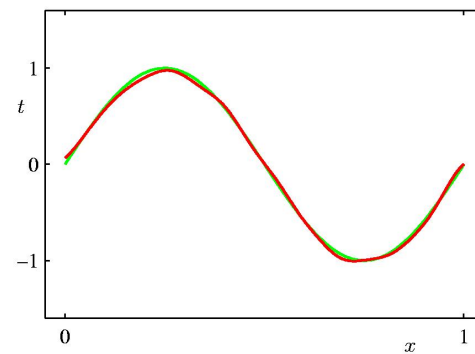
high variance

low variance

20 realizations



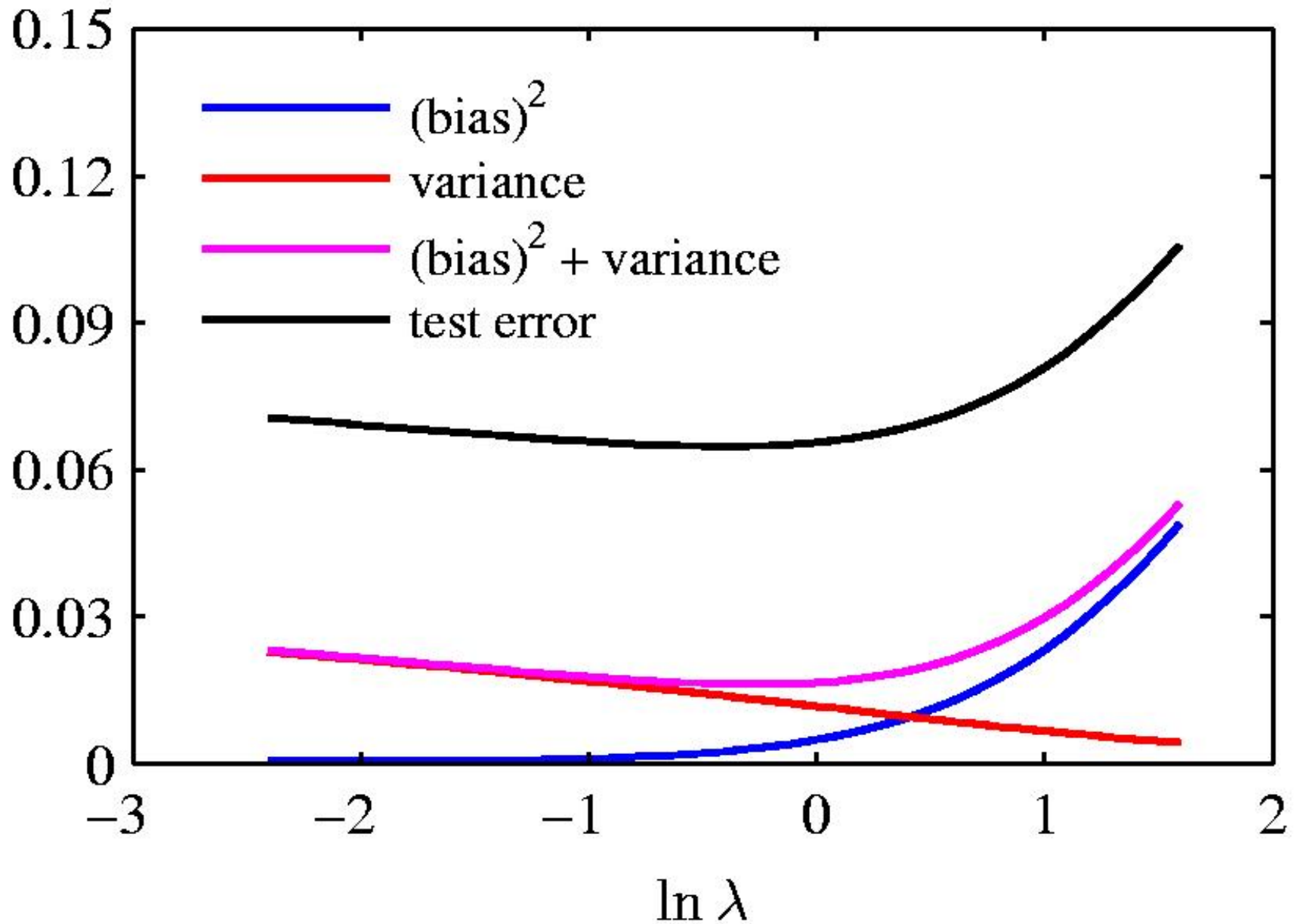
True model
average



low bias

high bias

An example of the bias-variance trade-off



Beating the bias-variance trade-off

- Reduce the variance term by averaging lots of models trained on different datasets.
 - Seems silly. For lots of different datasets it is better to combine them into one big training set.
 - More training data has much less variance.
- **Weird idea:** We can create different datasets by bootstrap sampling of our single training dataset.
 - This is called “**bagging**” and it works surprisingly well.
- If we have enough computation its better doing it **Bayesian:**
 - Combine the predictions of many models using the posterior probability of each parameter vector as the combination weight.