

Active Learning and Optimized Information Gathering

Lecture 6 – Gaussian Process Optimization

CS 101.2
Andreas Krause

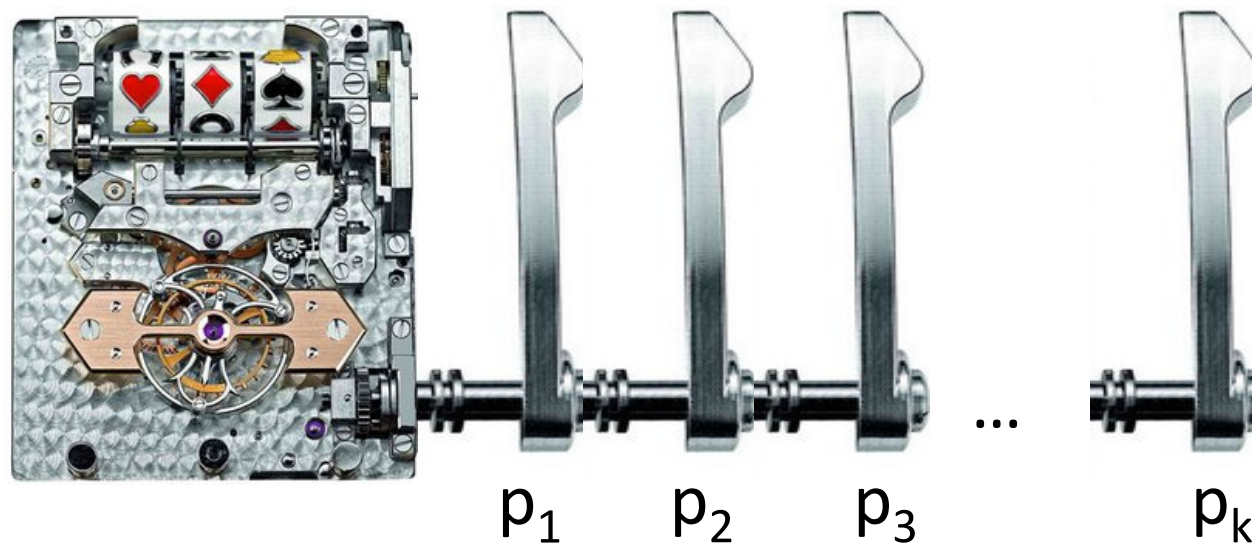
Announcements

- Homework 1: out tomorrow
 - Due Thu Jan 29
- Project
 - Proposal due Tue Jan 27
- Office hours
 - Come to office hours before your presentation!
 - Andreas: **Friday 1:30-3pm, 260 Jorgensen**
 - Ryan: Wednesday 4:00-6:00pm, 109 Moore

Course outline

1. Online decision making
2. Statistical active learning
3. Combinatorial approaches

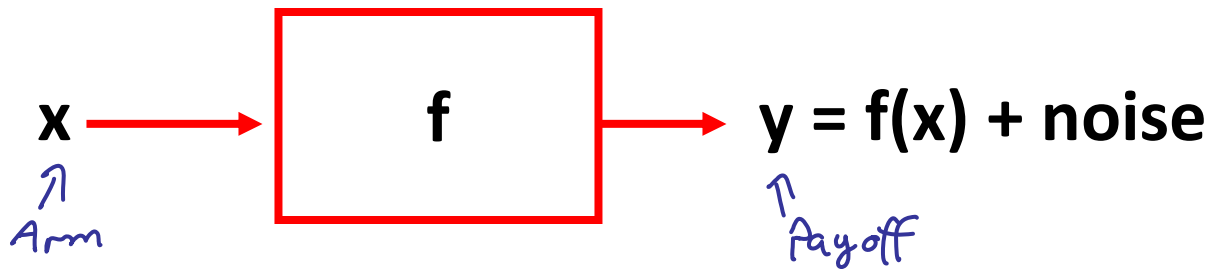
Recap Bandit problems



- K-arms
 - ϵ_n greedy, UCB1 have regret $O(\log(T) \mathbf{K})$
- What about infinite arms ($K=\infty$)
 - Have to make assumptions!

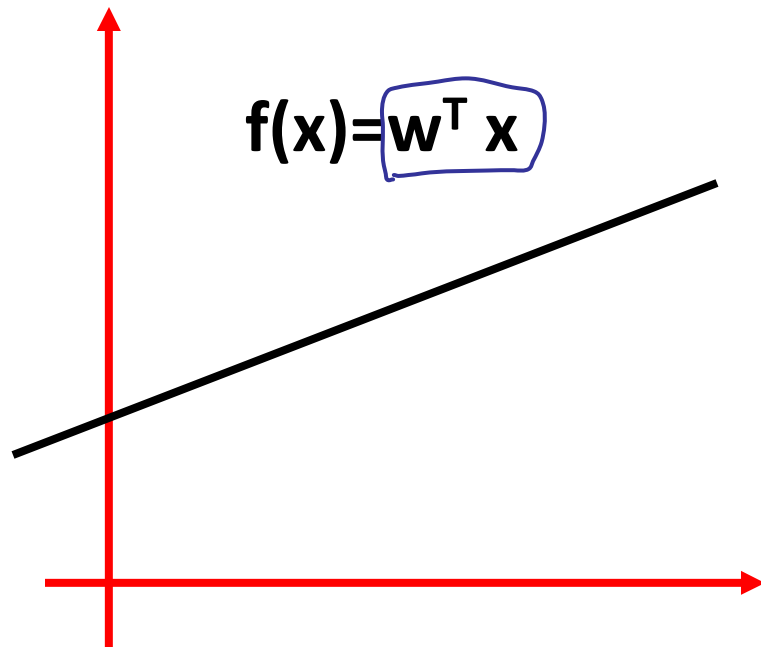
Bandits = Noisy function optimization

- We are given black box access to function f
 $f(x)$ = mean payoff for arm x

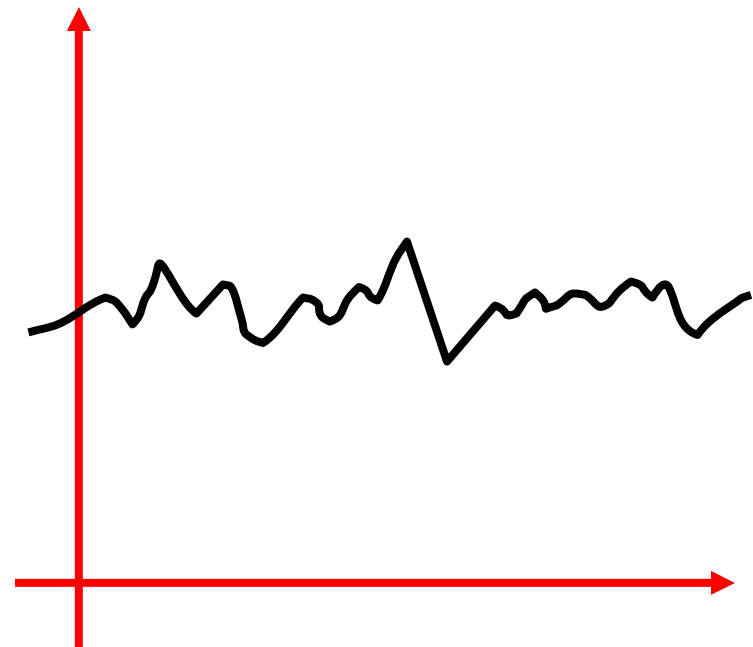


- Evaluating f is very expensive
- Want to (quickly) find $x^* = \operatorname{argmax}_x f(x)$

Bandits with ∞ -many arms



Linear



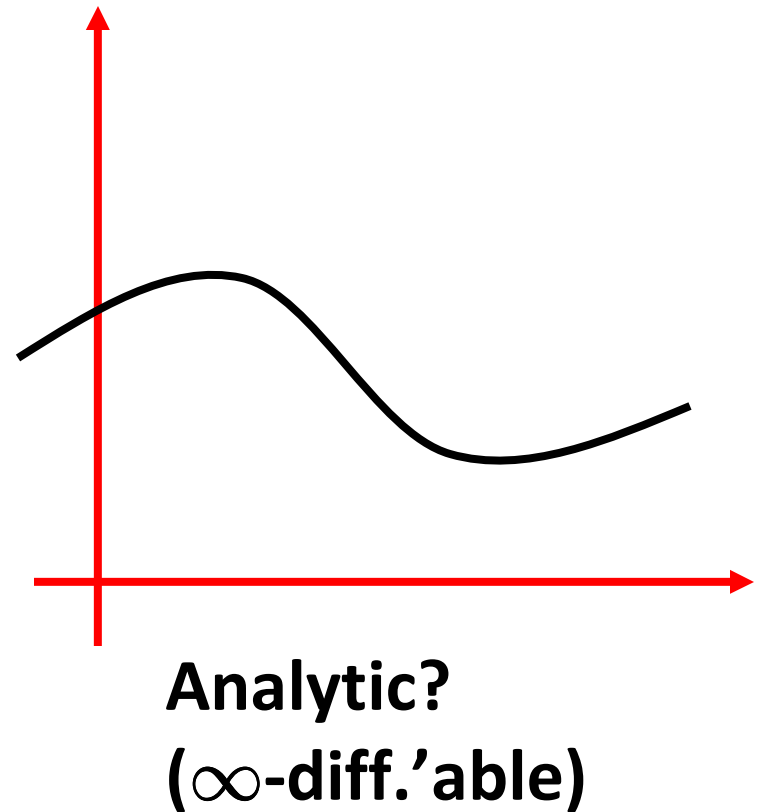
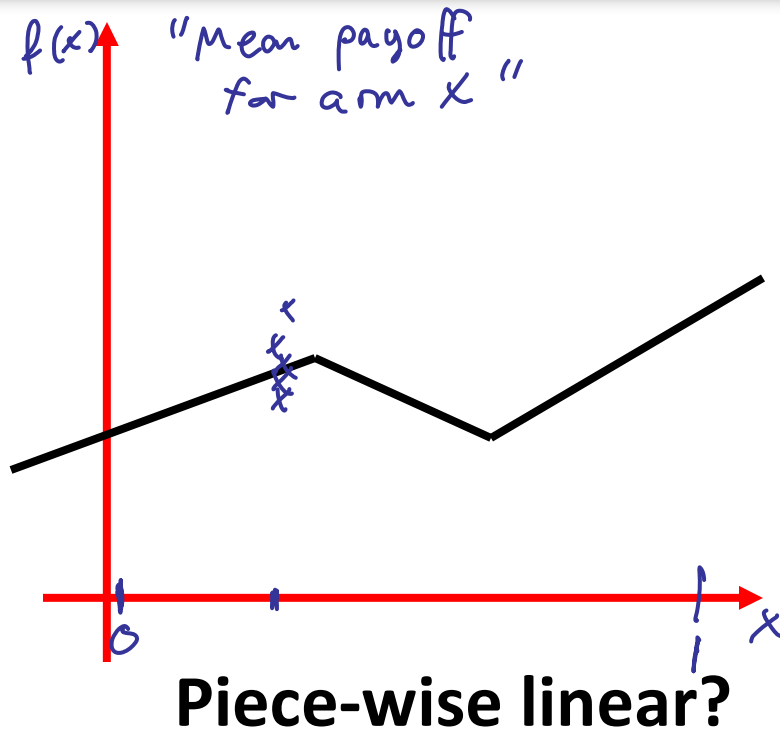
**Lipschitz-continuous
(bounded slope)**

- Can only hope to perform well if we make some assumptions

Regret depends on complexity

- Bandit linear optimization over \mathbb{R}^n
 - “strong” assumptions
 - Regret $O(T^{2/3}n)$
- Bandit problems for optimizing Lipschitz functions
 - “weak” assumptions
 - **Regret $O(C(n)T^{n/(n+1)})$**
 - Curse of dimensionality!
- Today: Flexible (Bayesian) approach for encoding assumptions about function complexity

What if we believe, the function looks like:



Want flexible way to encode assumptions about functions!

Bayesian inference

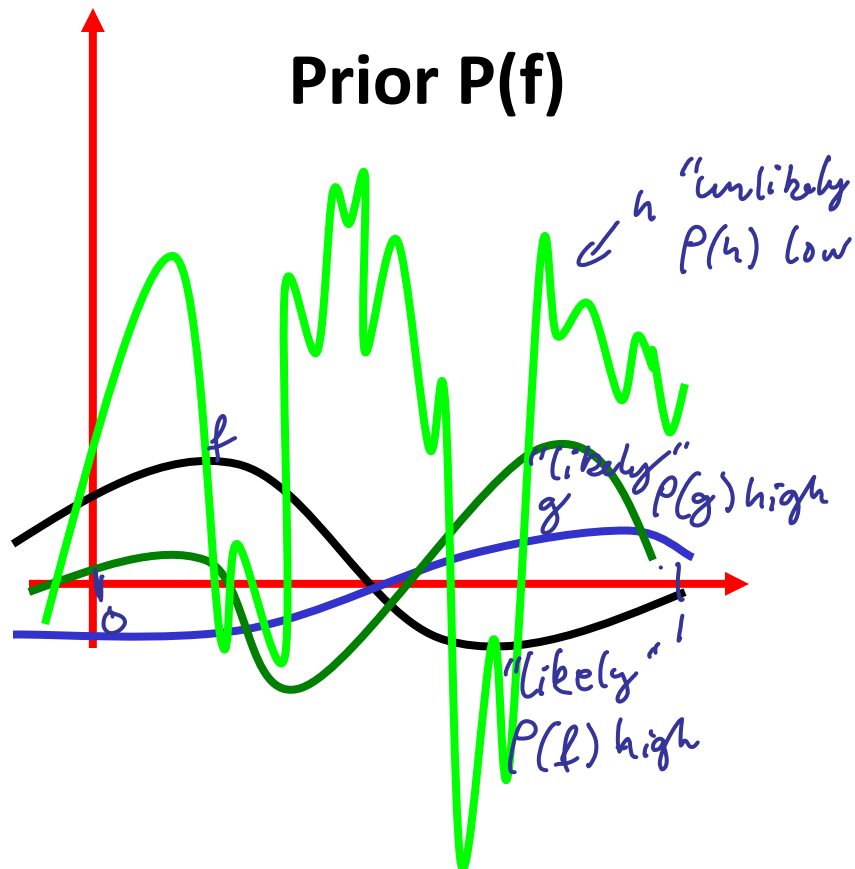
- Two Bernoulli variables A(larm), B(urglar)
- $P(B=1) = 0.1$; $P(A=1 | B=1) = 0.9$; $P(A=1 | B=0) = 0.1$
- What is $P(B | A)$?

$$P(B=1 | A=1) = \frac{P(B=1, A=1)}{P(A=1)} = \frac{P(A=1 | B=1) \cdot P(B=1)}{P(A=1 | B=1) P(B=1) + P(A=1 | B=0) P(B=0)}$$

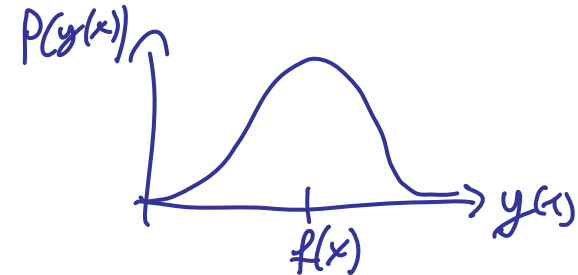
- $P(B)$ “prior”
- $P(A | B)$ “likelihood”
- $P(B | A)$ “posterior”

A Bayesian approach

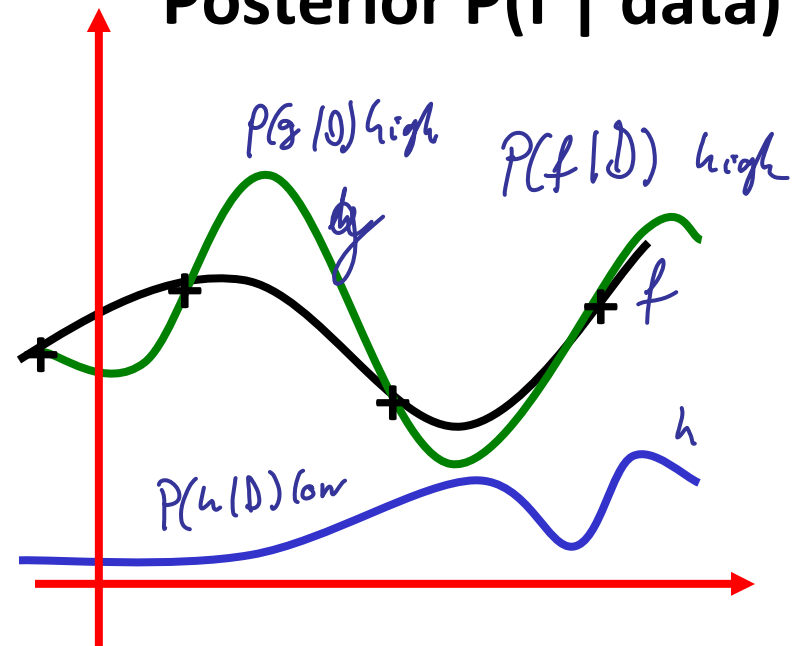
- Bayesian models for **functions**



Likelihood $P(\text{data} | f)$



Posterior $P(f | \text{data})$



- Uff... Why is this useful?

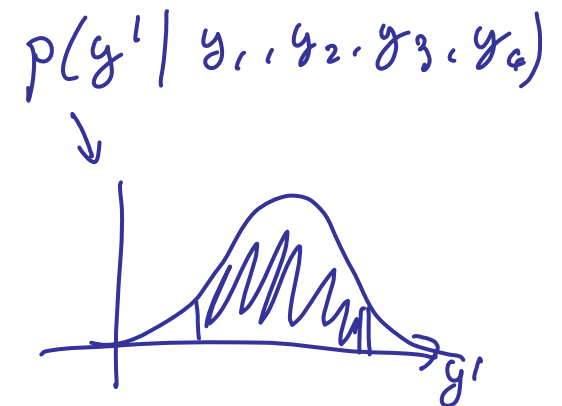
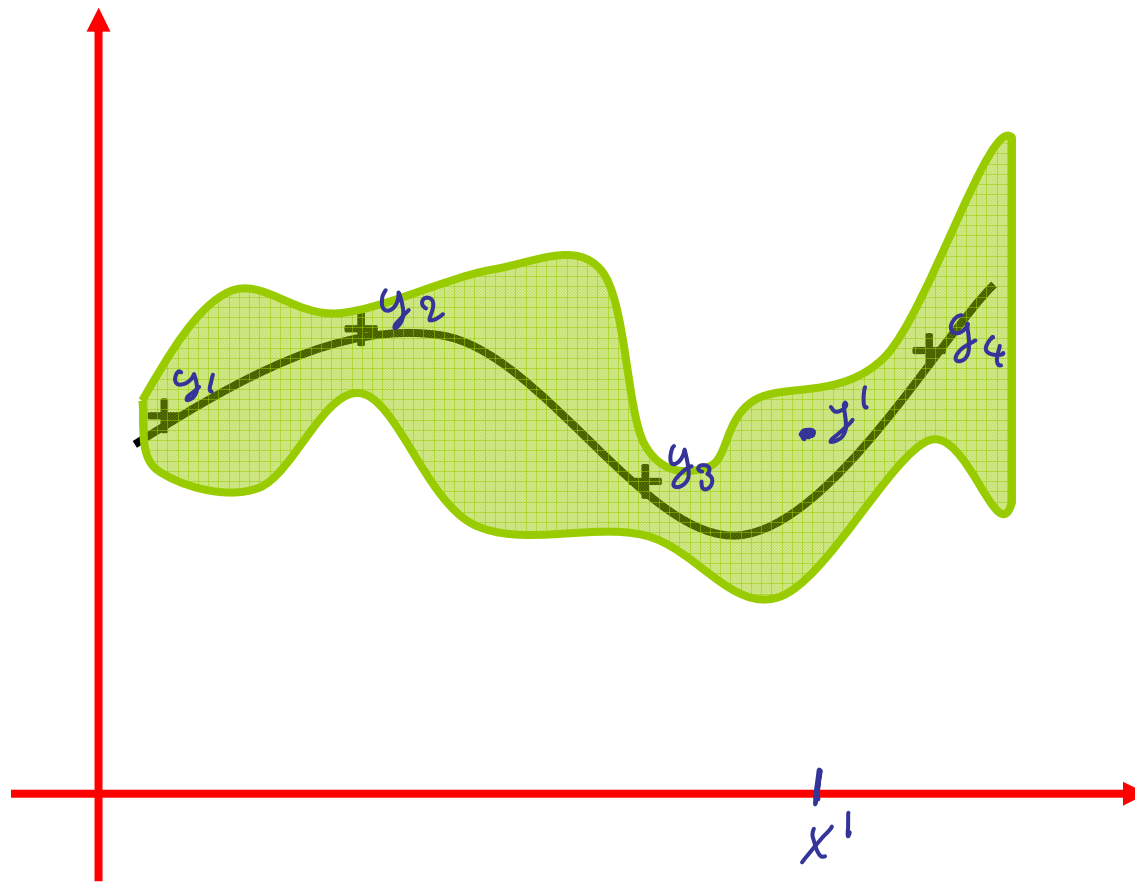
Probability of data

- $P(y_1, \dots, y_k) = \int P(f, y_1, \dots, y_k) d\rho(f) = \int \rho(f) \underset{\uparrow}{P(y_1, \dots, y_k | f)} d\rho(f)$

- Can compute

$$P(y' | y_1, \dots, y_k) = \frac{P(y', y_1, \dots, y_k)}{P(y_1, \dots, y_k)}$$

Regression with uncertainty about predictions!



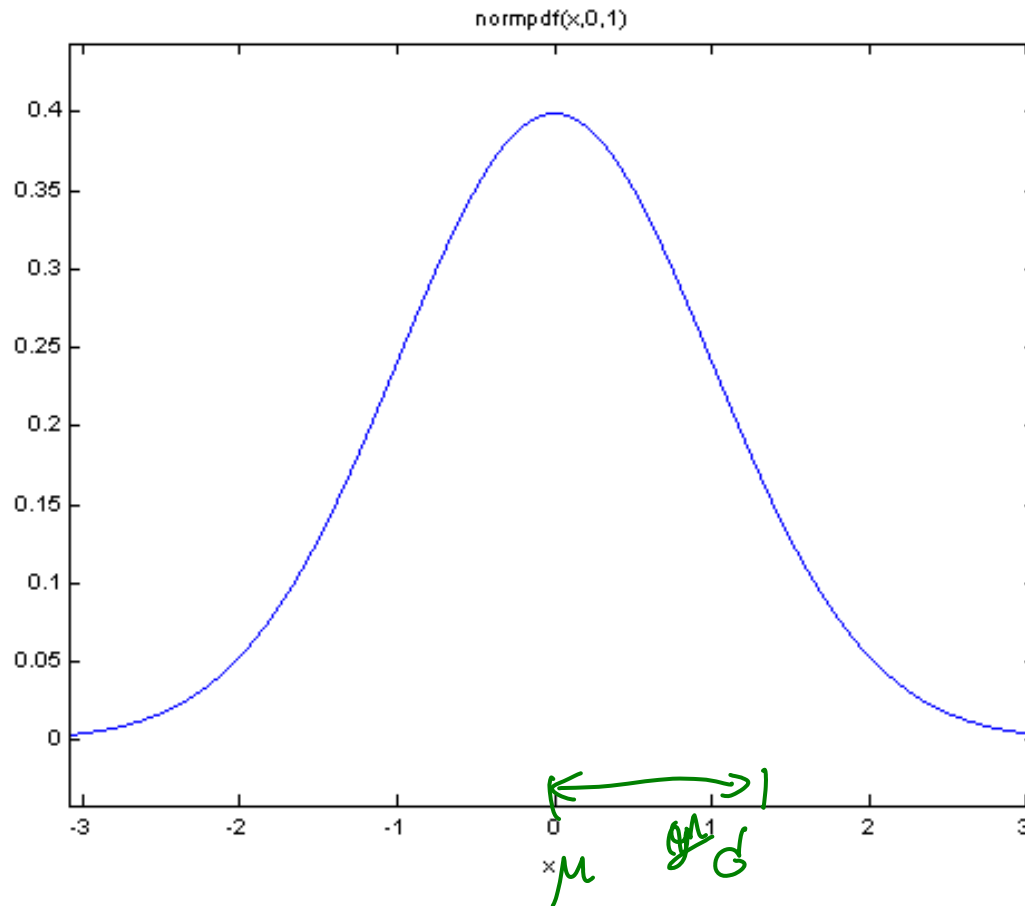
How can we do this?

- Want to compute $P(y' \mid y_1, \dots, y_k)$

$$P(y_1, \dots, y_k) = \int P(f, y_1, \dots, y_k) df$$

- Horribly complicated integral?? ☹️
 - Will see how we can compute this
 - (more or less) efficiently
 - In closed form!
- ... if $P(f)$ is a **Gaussian Process**

Gaussian distribution

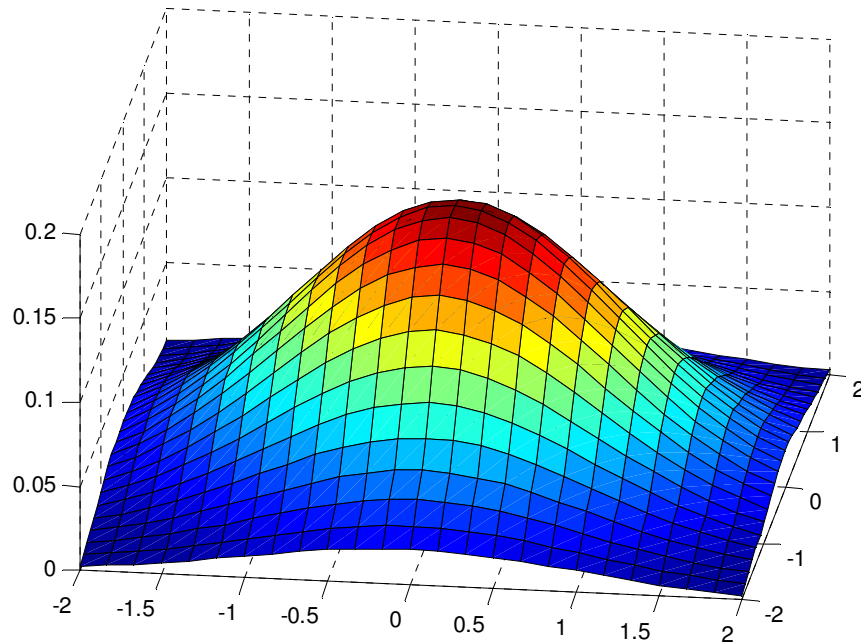


- σ = Standard deviation
 - μ = mean
- $$\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \mu)^2}{2\sigma^2}\right)$$

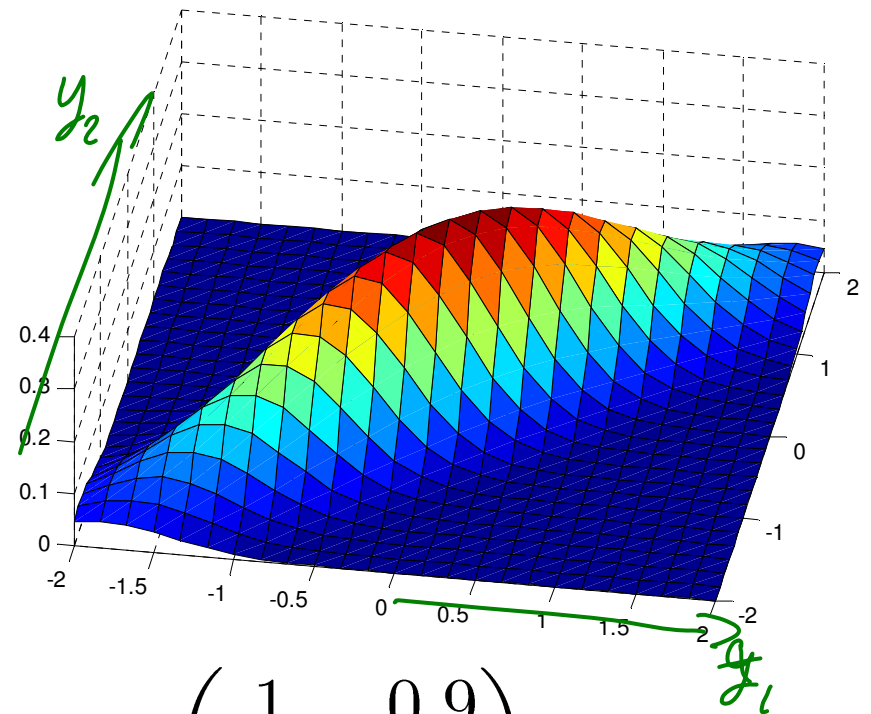
Bivariate Gaussian distribution

$$\frac{1}{2\pi\sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(y-\mu)^T \Sigma^{-1}(y-\mu)\right)$$

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$$



$$\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$



$$\Sigma = \begin{pmatrix} 1 & 0.9 \\ 0.9 & 1 \end{pmatrix}$$

Multivariate Gaussian distribution

$$\mathcal{N}(y; \Sigma, \mu) = \frac{1}{(2\pi)^{n/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(y - \mu)^T \Sigma^{-1} (y - \mu)\right)$$

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1n} \\ \vdots & & & \vdots \\ \sigma_{n1} & \sigma_{n2} & \dots & \sigma_n^2 \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix} \quad y \in \mathbb{R}^n$$

$\Sigma \in \mathbb{R}^{n \times n}$

- Joint distribution over n random variables $P(Y_1, \dots, Y_n)$

- $\sigma_{jk} = \underline{E[(Y_j - \mu_j)(Y_k - \mu_k)]}$

$\sigma_{jk} > 0$: high values of Y_j
 \Rightarrow high values of Y_k

- Y_j and Y_k independent $\Leftrightarrow \sigma_{jk} = 0$

only true for Gaussians!

Marginalization

- Suppose $(Y_1, \dots, Y_n) \sim \mathcal{N}(\mu, \Sigma)$

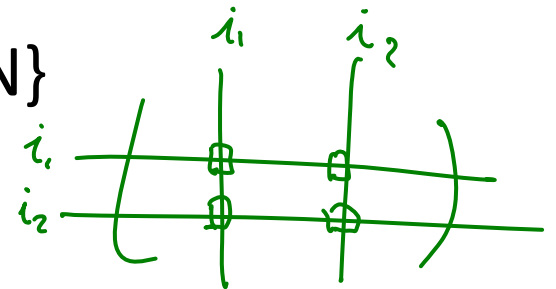
- What is $\mathcal{P}(Y_1)$??

$$Y_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$$

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1n} \\ \vdots & \vdots & \vdots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \dots & \sigma_n^2 \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix}$$

- More generally: Let $A = \{i_1, \dots, i_k\} \subseteq \{1, \dots, N\}$

- Write $Y_A = (Y_{i_1}, \dots, Y_{i_k})$



- $Y_A \sim \mathcal{N}(\mu_A, \Sigma_{AA})$

$$\Sigma_{AA} = \begin{pmatrix} \sigma_{i_1 i_1}^2 & \sigma_{i_1 i_2} & \dots & \sigma_{i_1 i_k} \\ \vdots & \vdots & \vdots & \vdots \\ \sigma_{i_k i_1} & \sigma_{i_k i_2} & \dots & \sigma_{i_k i_k}^2 \end{pmatrix} \quad \mu_A = \begin{pmatrix} \mu_{i_1} \\ \mu_{i_2} \\ \vdots \\ \mu_{i_k} \end{pmatrix}$$

Conditioning

- Suppose $(Y_1, \dots, Y_n) \sim N(\mu, \Sigma)$
- Decompose as (Y_A, Y_B)
- What is $P(Y_A | Y_B)$??

$$\mu = \begin{pmatrix} \mu_A \\ \mu_B \end{pmatrix} \quad \Sigma = \begin{pmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{pmatrix}$$

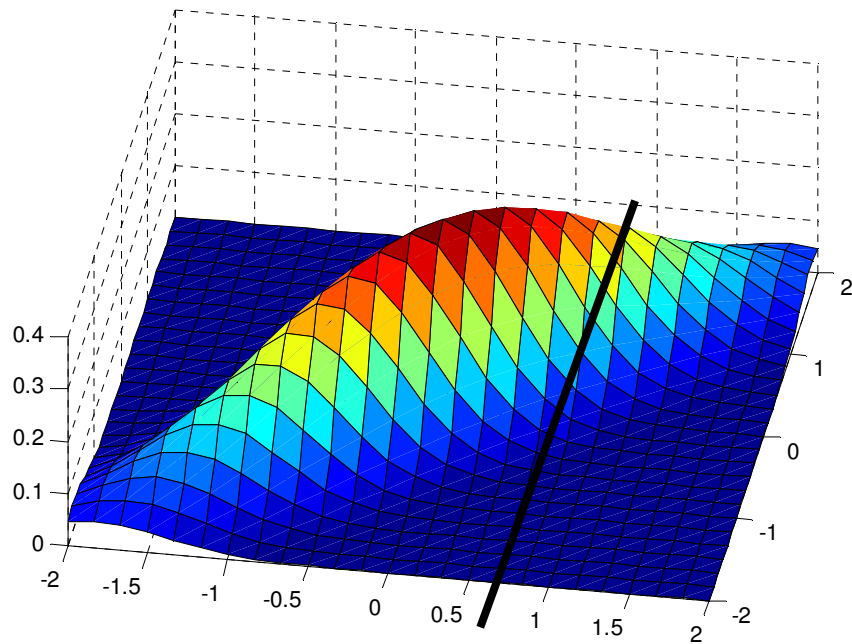
- $P(Y_A = y_A | Y_B = y_B) = N(y_A; \mu_{A|B}, \Sigma_{A|B})$ where

$$\mu_{A|B} = \mu_A + \Sigma_{AB} \Sigma_{BB}^{-1} (y_B - \mu_B)$$
$$\Sigma_{A|B} = \Sigma_{AA} - \Sigma_{AB} \Sigma_{BB}^{-1} \Sigma_{BA}$$

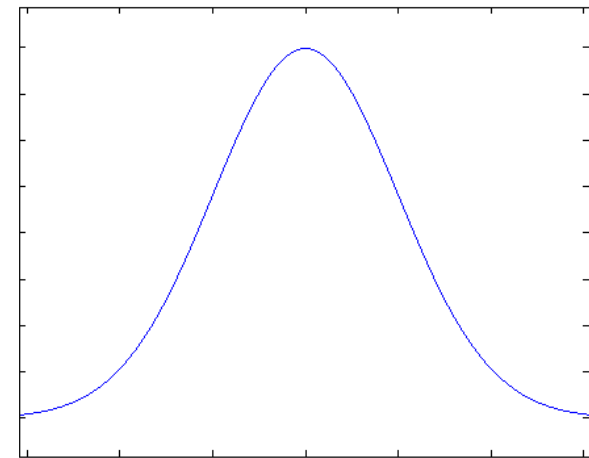
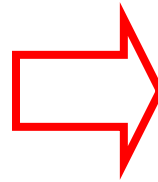
posterior mean → $\mu_{A|B}$
prior mean → μ_A

- Computable using linear algebra!

Conditioning



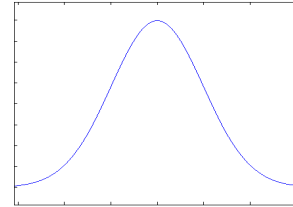
$Y_1 = 0.75$



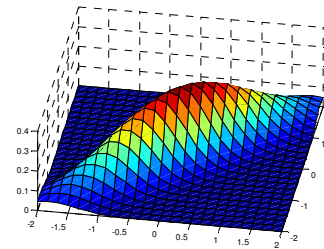
$P(Y_2 | Y_1 = 0.75)$

High dimensional Gaussians

- Gaussian



- Bivariate Gaussian



- Multivariate Gaussian

$$\mathcal{N}(y; \Sigma, \mu) = \frac{1}{(2\pi)^{n/2} \sqrt{|\Sigma|}} \exp \left(-\frac{1}{2} (y - \mu)^T \Sigma^{-1} (y - \mu) \right)$$

- Gaussian Process = “ ∞ -variate Gaussian”

Gaussian process

- A Gaussian Process (GP) is a
 - (infinite) set of random variables, indexed by some set V i.e., for each $x \in V$ there's a RV Y_x
 - Let $A \subseteq V$, $|A| = \{x_1, \dots, x_k\} < \infty$Then

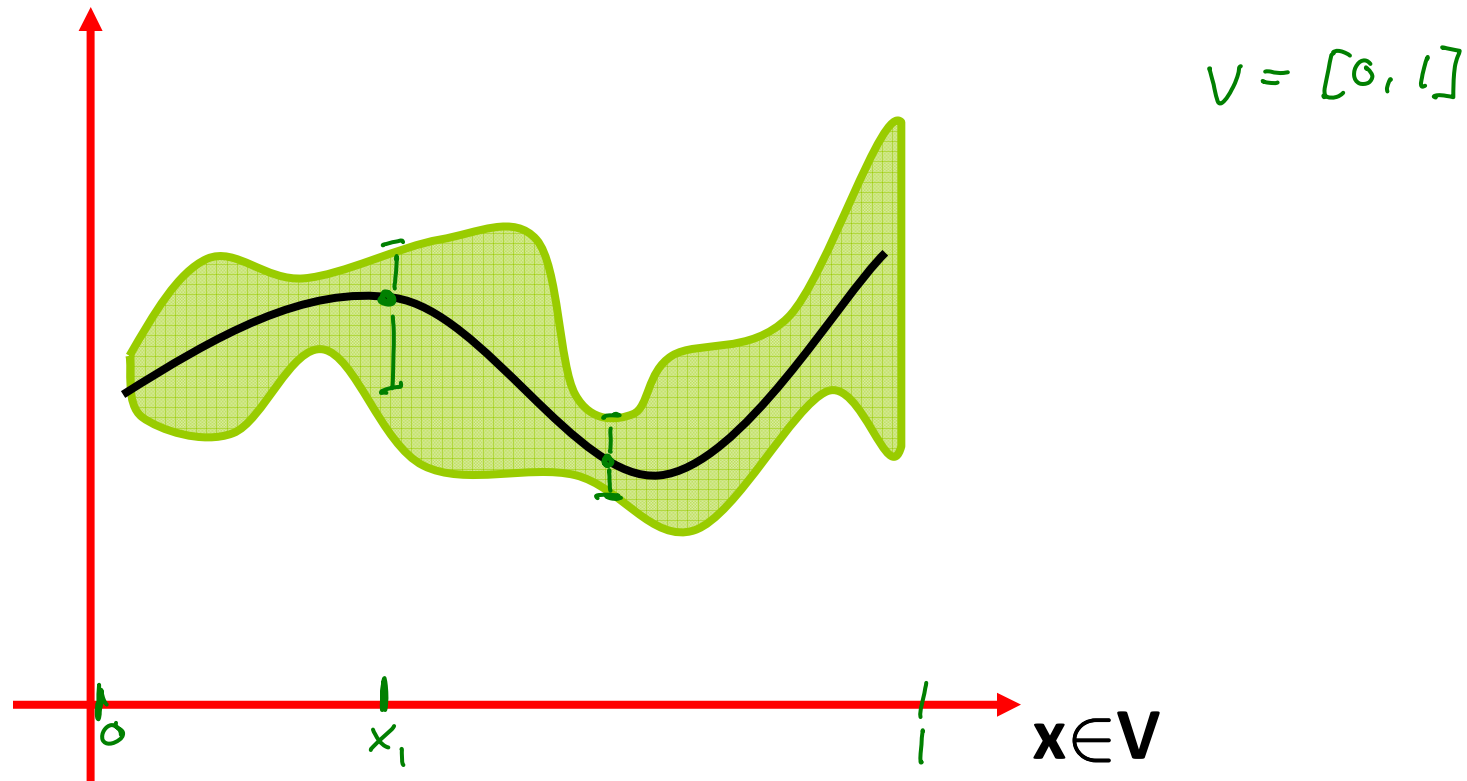
$$Y_A \sim N(\mu_A, \Sigma_{AA})$$

where

$$\Sigma_{AA} = \begin{pmatrix} \mathcal{K}(x_1, x_1) & \mathcal{K}(x_1, x_2) & \dots & \mathcal{K}(x_1, x_n) \\ \vdots & & & \vdots \\ \mathcal{K}(x_k, x_1) & \mathcal{K}(x_k, x_2) & \dots & \mathcal{K}(x_k, x_k) \end{pmatrix} \quad \mu_A = \begin{pmatrix} \mu(x_1) \\ \mu(x_2) \\ \vdots \\ \mu(x_k) \end{pmatrix}$$

- $\mathcal{K}: V \times V \rightarrow \mathbb{R}$ is called **kernel** (covariance) function
- $\mu: V \rightarrow \mathbb{R}$ is called **mean** function

Visualizing GPs



- Typically, only care about “marginals”, i.e.,
 $P(y) = N(y; \mu(x), K(x,x))$

Mean functions

- Can encode prior knowledge
- Typically, one simply assumes

$$\mu(\mathbf{x}) = 0$$

- Will do that here to simplify notation

Not a strong assumption

Kernel functions

- K must be **symmetric**

$$K(x, x') = K(x', x) \text{ for all } x, x' \in V$$

- K must be **positive definite**

For all A: Σ_{AA} is positive definite matrix

Matrix $M \in \mathbb{R}^{n \times n}$ pos. def.
 $\Leftrightarrow \forall x \in \mathbb{R}^n \quad x^T M x \geq 0$
 \Leftrightarrow All eigenvalues > 0

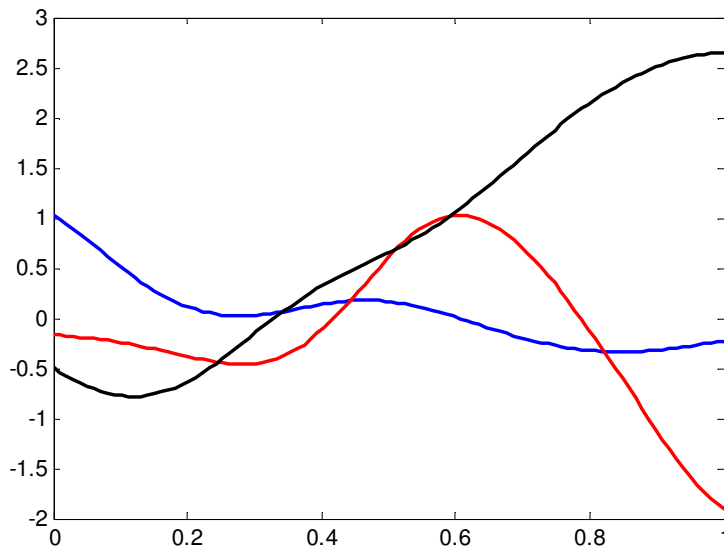
- Kernel function K: assumptions about correlation!

Kernel functions: Examples

- Squared exponential kernel

$$K(x, x') = \exp(-(x-x')^2/h^2)$$

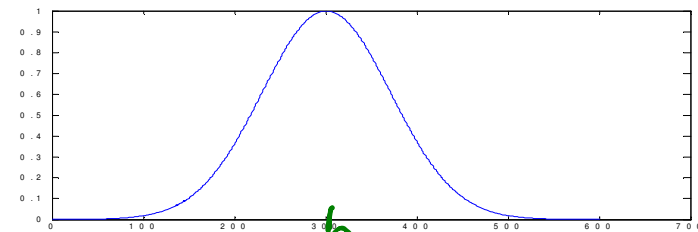
Samples from $P(f)$



Samples are ∞ -differentiable with Prob 1

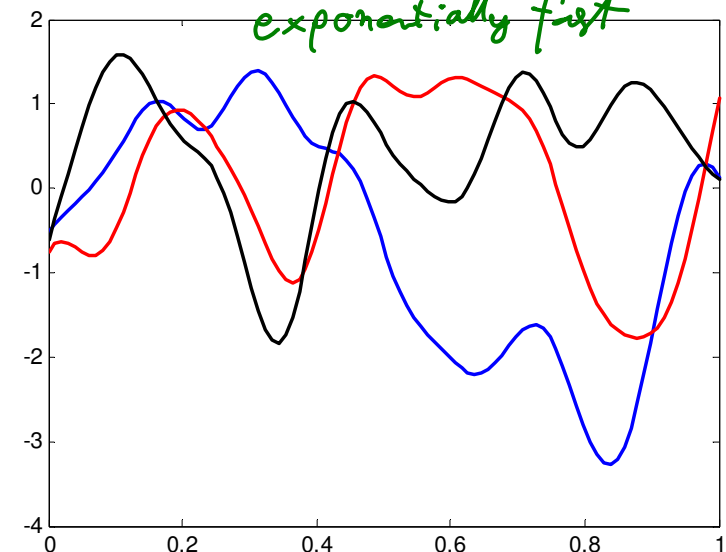
Bandwidth $h=0.3$

$V = \mathbb{R}$



Distance $|x-x'|$

"Correlation decays with distance"
↑
exponentially fast

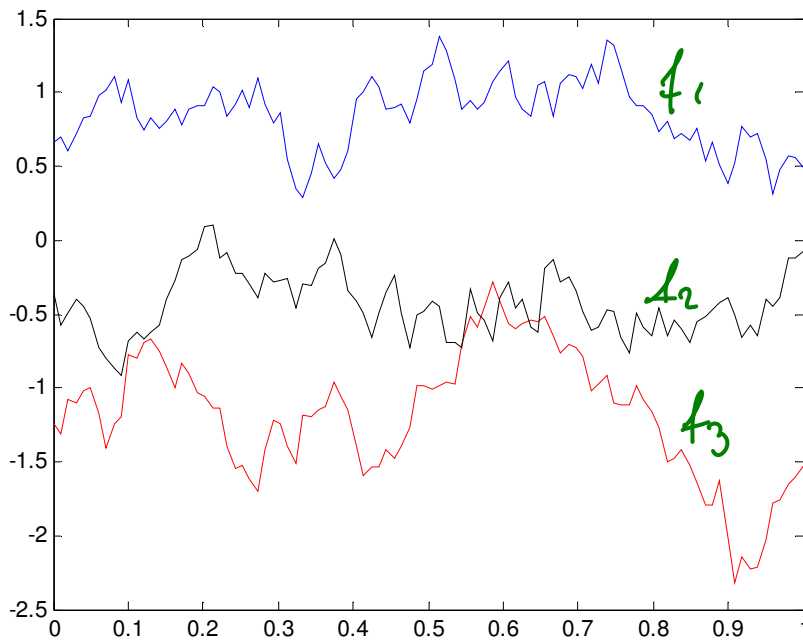
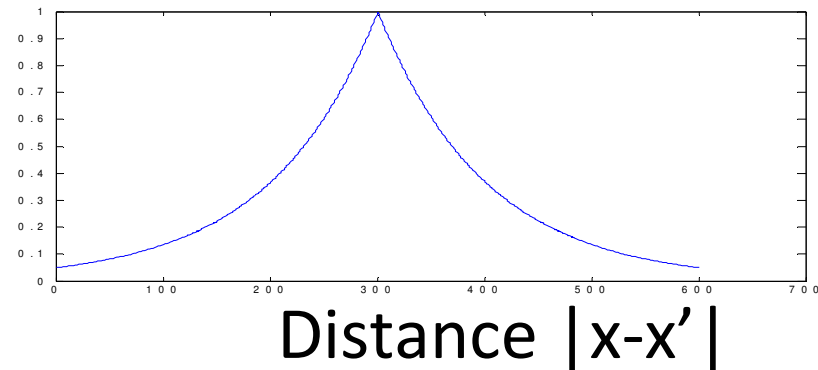


Bandwidth $h=0.1$

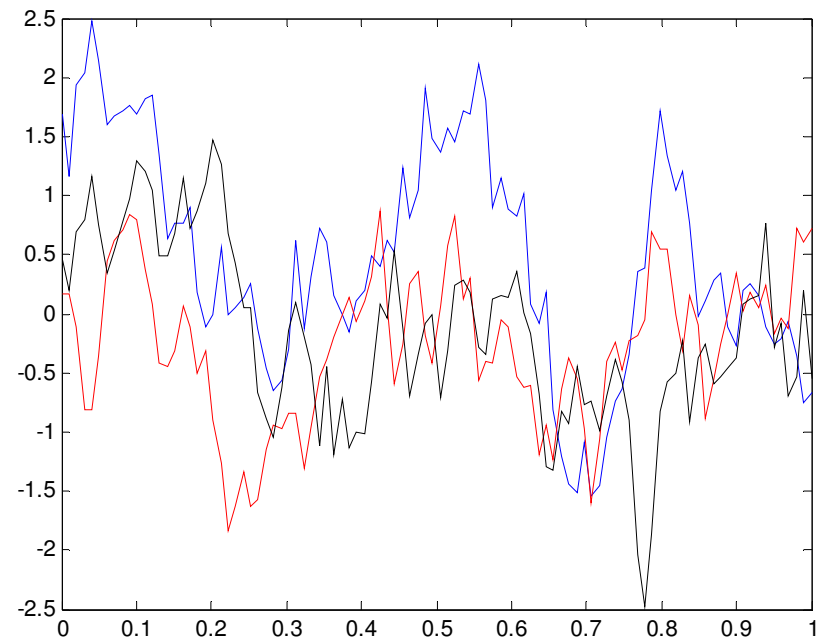
Kernel functions: Examples

- Exponential kernel
 $K(x, x') = \exp(-|x-x'|/h)$

$$f_1 \sim \text{GP}(\mu, \kappa)$$



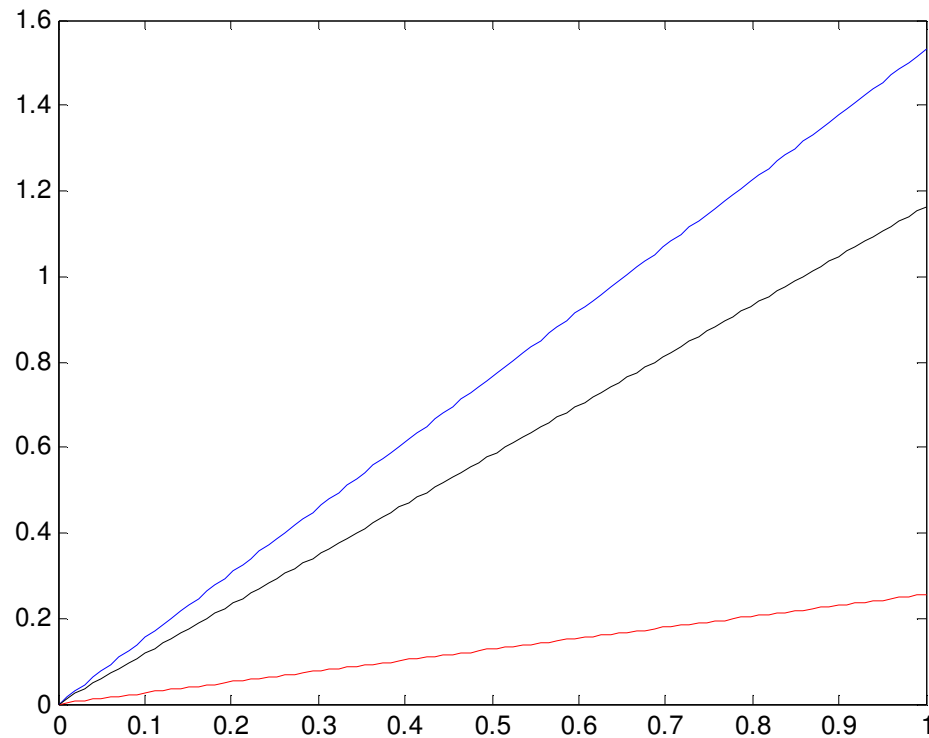
Bandwidth $h=1$



Bandwidth $h=.3$

Kernel functions: Examples

- Linear kernel:
 $K(x, x') = x^T x'$



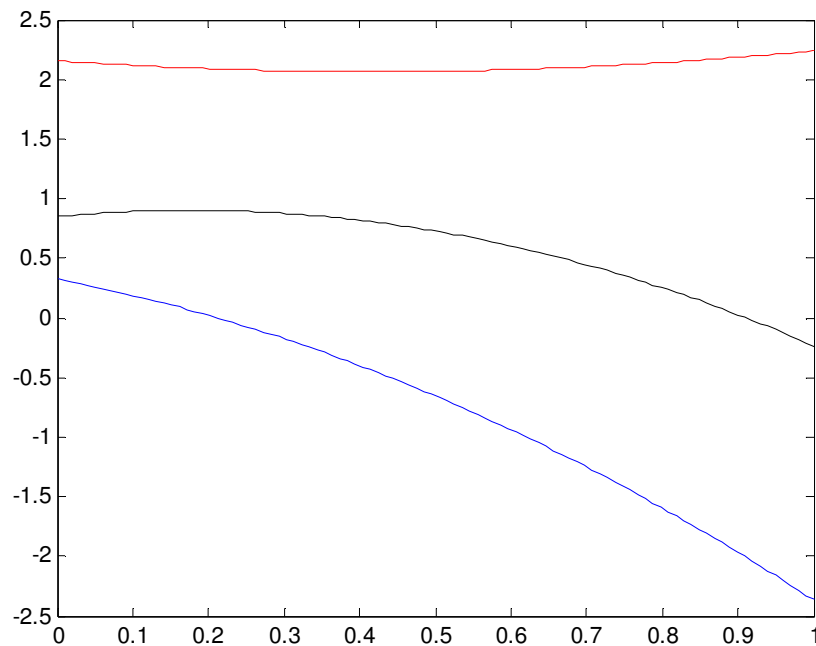
- Corresponds to linear regression!

Kernel functions: Examples

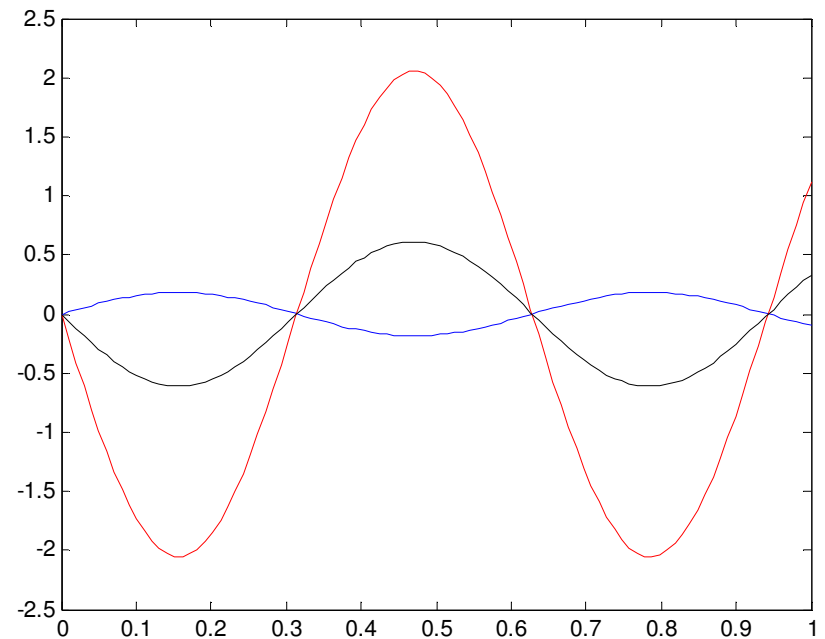
- Linear kernel with features:

$$K(x, x') = \Phi(x)^\top \Phi(x')$$

E.g., $\Phi(x) = [1, x, x^2]$



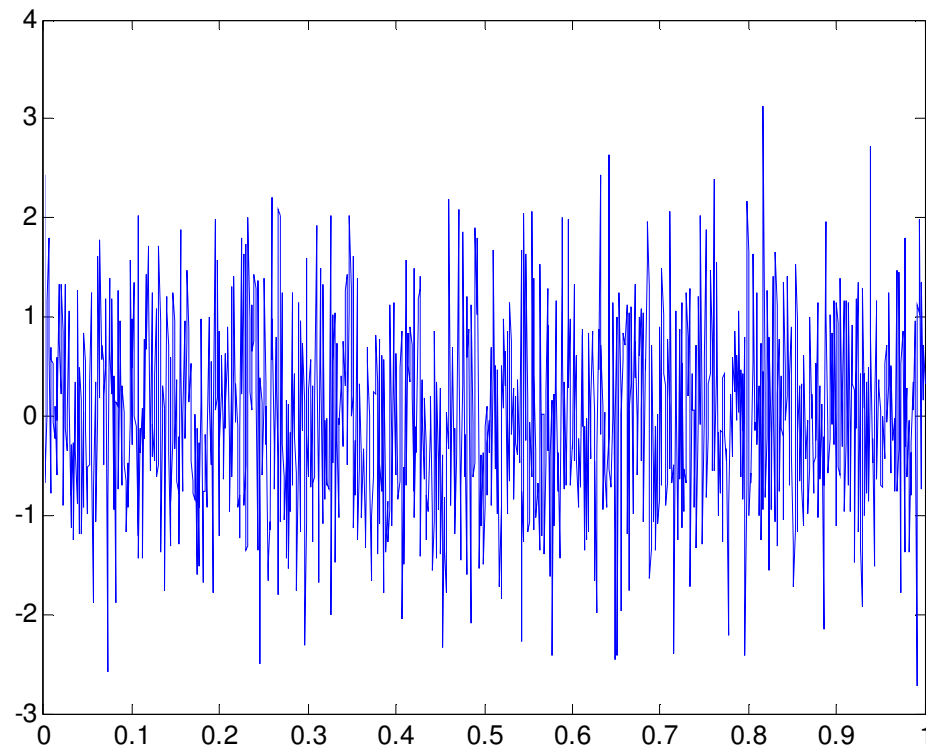
E.g., $\Phi(x) = \sin(x)$



Kernel functions: Examples

- White noise:

$$K(x,x) = 1; K(x,x') = 0 \text{ for } x' \neq x$$



Constructing kernels from kernels

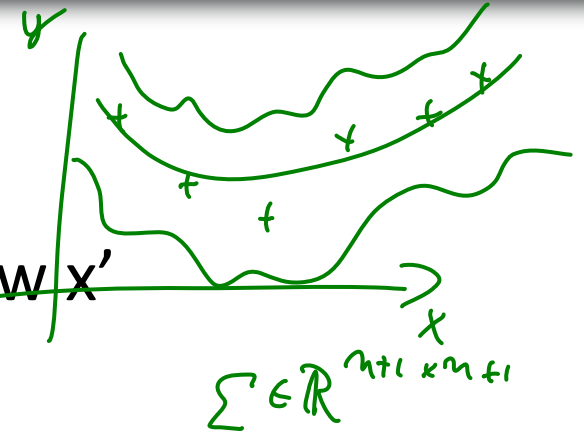
If $K_1(x,x')$ and $K_2(x,x')$ are kernel functions then

$\alpha K_1(x,x') + \beta K_2(x,x')$ is a kernel for $\alpha, \beta > 0$

$K_1(x,x') * K_2(x,x')$ is a kernel

GP Regression

- Suppose we know kernel function K
- Get data $(x_1, y_1), \dots, (x_n, y_n)$
- Want to predict $y' = f(x')$ for some new x'



$$P(y' | y_1, \dots, y_n)$$

1. Construct $P(y' | y_1, \dots, y_n) = \mathcal{N}(\cdot; \mu, \Sigma)$

2. Condition:

$$P(y' | y_1, \dots, y_n) = \mathcal{N}(y'; \mu_{x'|D}, \sigma_{x'|D}^2)$$

$$\mu_{x'|D} = \underbrace{\mu_{x'}}_0 + \sum_{x_i \in D} \Sigma_{DD}^{-1} (y_D - \mu_D)$$

$$\sigma_{x'|D}^2 = \sigma_{x'}^2 + \sum_{x_i \in D} \Sigma_{DD}^{-1} \Sigma_{D, x'}$$

$$\mu = \begin{pmatrix} \mu(x') \\ \mu(x_1) \\ \vdots \\ \mu(x_n) \end{pmatrix} = 0$$

$$\Sigma = \begin{pmatrix} k(x', x') & k(x', x_1) & \dots & k(x', x_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_m, x') & \dots & \dots & k(x_n, x') \end{pmatrix}$$

"Read off" Σ

Linear prediction

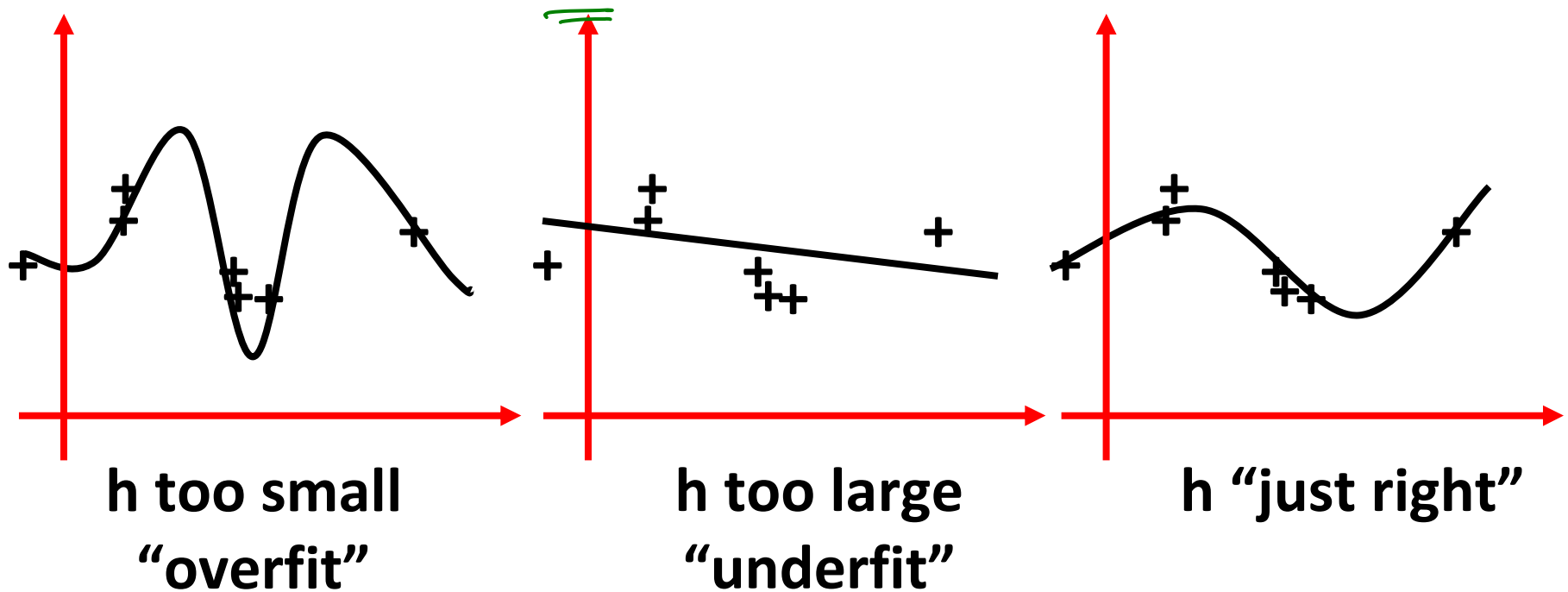
- Posterior mean $\mu_{x'|D} = \underbrace{\Sigma_{x',D} \Sigma_{D,D}^{-1}}_{\in \mathbb{R}^{k \times m}} y_D$

$$\Sigma = \begin{pmatrix} \Sigma_{AA} & \Sigma_{AD} \\ \Sigma_{DA} & \Sigma_{DD} \end{pmatrix}$$

$$A = \{x'\}, D = D$$
- Hence, $\mu_{x'|D} = \sum_{i=1}^n \underbrace{w_i}_{\text{Regression coefficients}} y_i = \underbrace{w^T}_{\text{depends on } k} y$
- Prediction $\mu_{x'|D}$ depends **linearly** on inputs y_i !
- For fixed data set $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$, can *precompute* weights w_i
- Like linear regression, but number of parameters w_i grows with training data
 - “Nonparametric regression”
 - Can fit any data set!! 😊

Learning parameters

- Example: $K(x, x') = \exp(-(x-x')^2/h^2)$
Need to specify h !



- In general, kernel function has parameters θ
- Want to learn θ from data

Learning parameters

- Pick parameters that make data most likely!

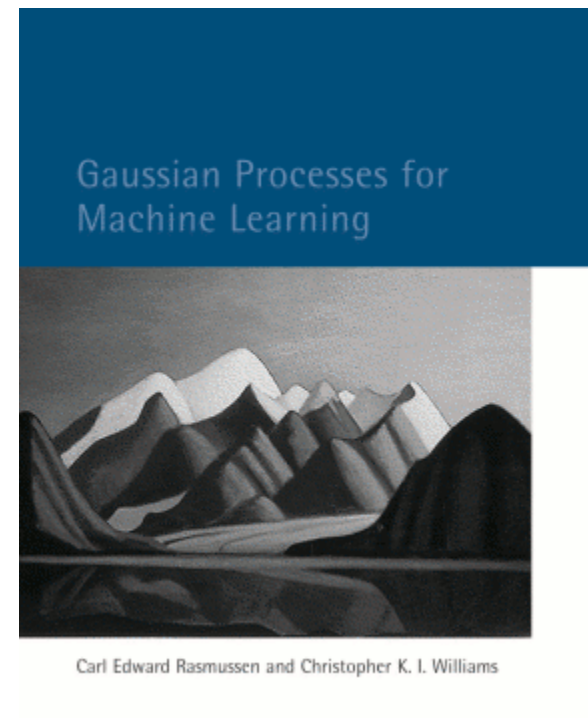
$$P(y; x, \theta) = \frac{1}{(2\pi)^{n/2} \sqrt{|K_\theta|}} \exp\left(-\frac{1}{2} y^T K_\theta^{-1} y\right)$$

$$\begin{aligned} \theta^* \operatorname{argmax}_{\theta} P(y; x, \theta) &= \operatorname{argmax}_{\theta} \log P(y; x, \theta) \\ &= \operatorname{const} - \underbrace{\frac{1}{2} \log |K_\theta|}_{\text{complexity}} - \underbrace{\frac{1}{2} y^T K_\theta^{-1} y}_{\text{goodness of}} \end{aligned}$$

- $\log P(y | \theta)$ differentiable if $K(x, x')$ is!
 - Can do gradient descent, conjugate gradient, etc.
- Tends to work well (not over- or underfit) in practice!

Matlab demo

- [Rasmussen & Williams, Gaussian Processes for Machine Learning]
- <http://www.gaussianprocess.org/gpml/>



Gaussian process

- A Gaussian Process (GP) is a
 - (infinite) set of random variables, indexed by some set V i.e., for each $x \in V$ there's a RV Y_x
 - Let $A \subseteq V$, $|A| = \{x_1, \dots, x_k\} < \infty$

Then

$$Y_A \sim N(\mu_A, \Sigma_{AA})$$

where

$$\Sigma_{AA} = \begin{pmatrix} \mathcal{K}(x_1, x_1) & \mathcal{K}(x_1, x_2) & \dots & \mathcal{K}(x_1, x_n) \\ \vdots & & & \vdots \\ \mathcal{K}(x_k, x_1) & \mathcal{K}(x_k, x_2) & \dots & \mathcal{K}(x_k, x_k) \end{pmatrix} \quad \mu_A = \begin{pmatrix} \mu(x_1) \\ \mu(x_2) \\ \vdots \\ \mu(x_k) \end{pmatrix}$$

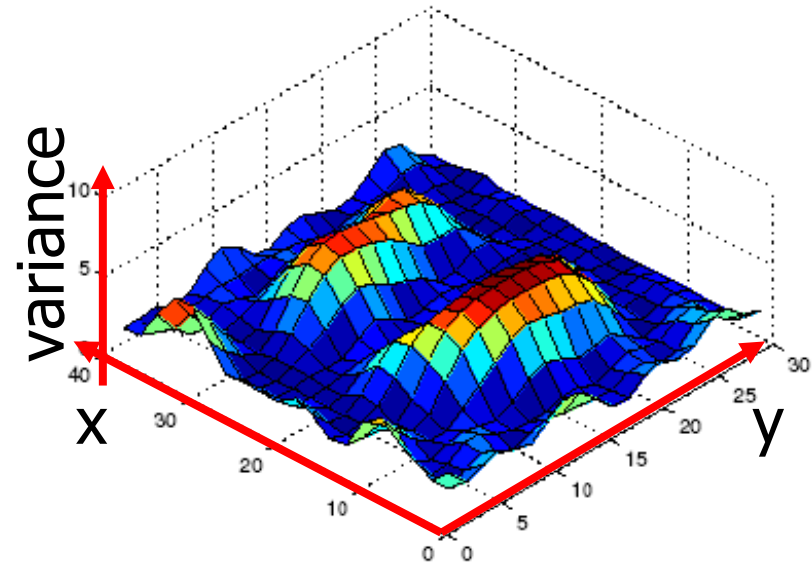
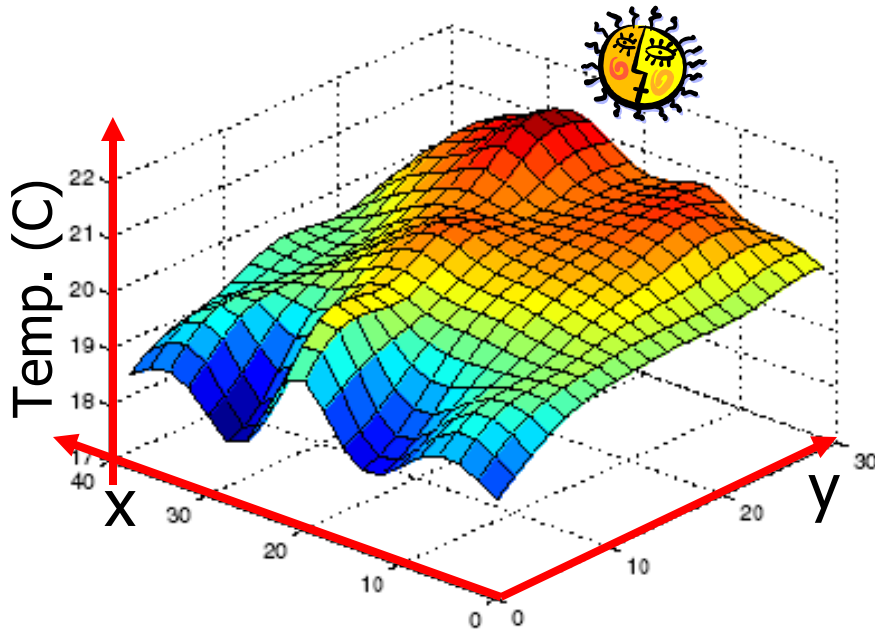
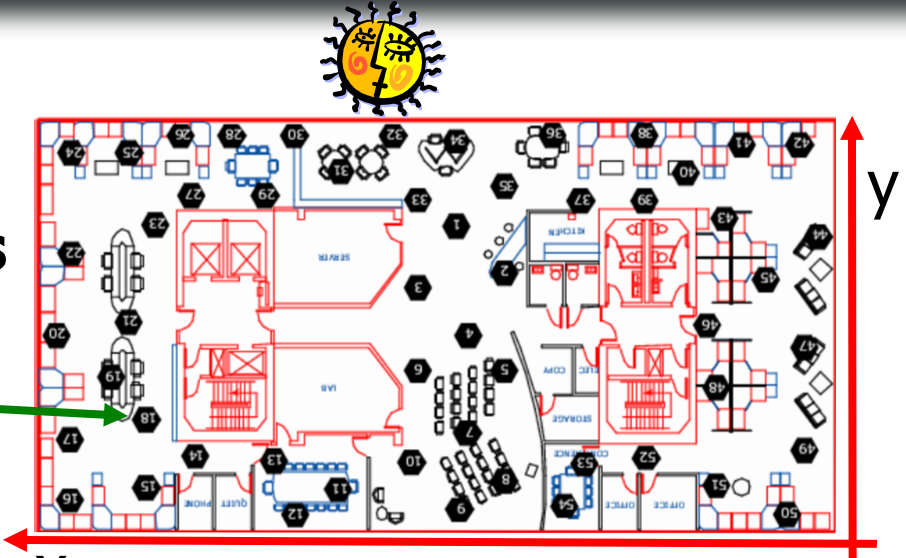
- $\mathcal{K}: V \times V \rightarrow \mathbb{R}$ is called **kernel** (covariance) function
- $\mu: V \rightarrow \mathbb{R}$ is called **mean** function

GPs over other sets

- GP is collection of random variables, indexed by set V
- So far: Have seen GPs over $V = \mathbb{R}$
- Can define GPs over
 - Text (strings)
 - Graphs
 - Sets
 - ...
- Only need to choose appropriate kernel function

Example: Using GPs to model spatial phenomena

**Real deployment
of temperature sensors**
measurements from 52 sensors
(black dots)



Predicted temperature
throughout the space

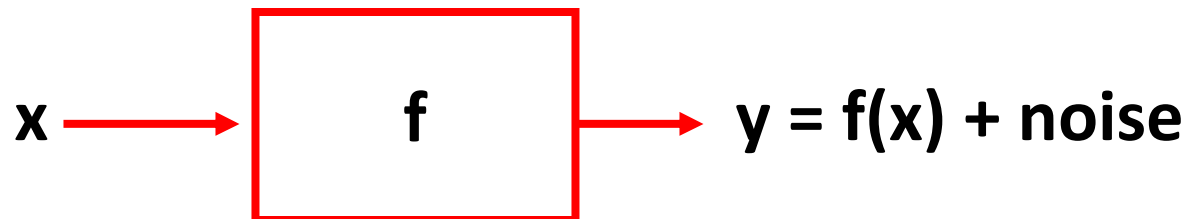
Other extensions (won't cover here)

- GPs for classification
 - Nonparametric generalization of logistic regression
 - Like SVMs (but give confidence on predicted labels!)
- GPs for modeling non-Gaussian phenomena
 - Model count data over space, ...
- Active set methods for fast inference
- ...

Still active research area in machine learning

Bandits = Noisy function optimization

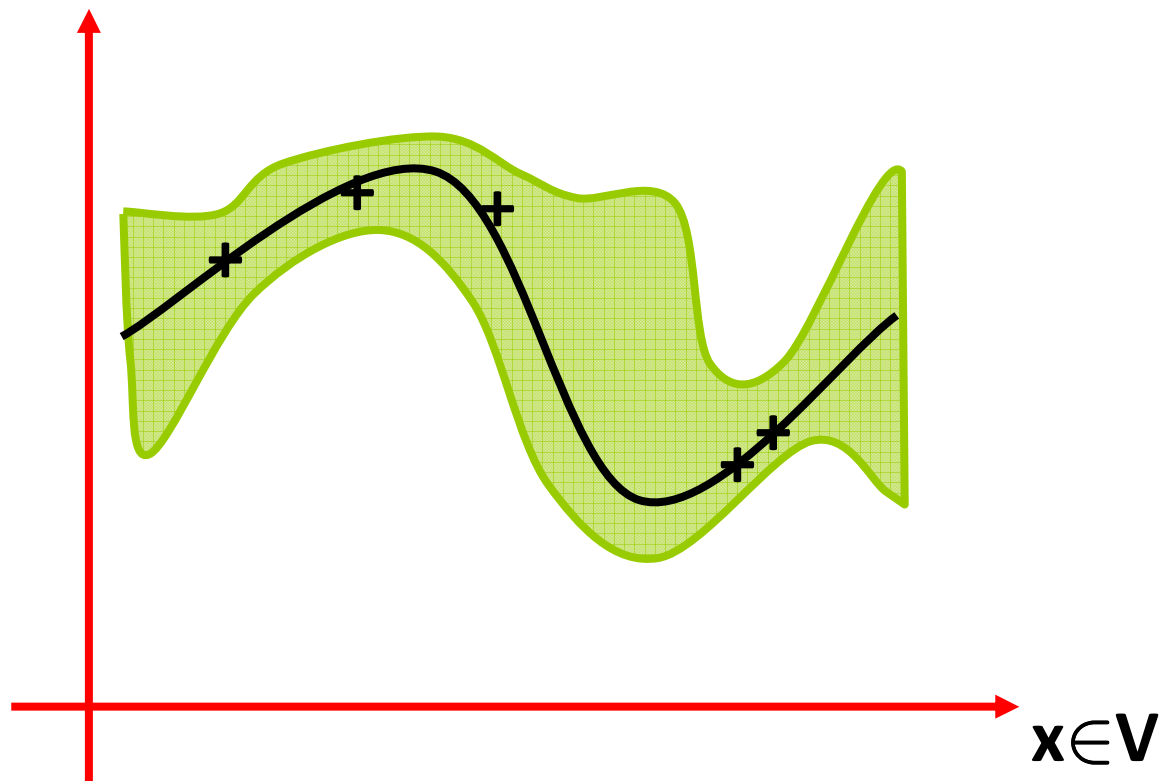
- We are given black box access to function f



- Evaluating f is very expensive
- Want to (quickly) find $x^* = \operatorname{argmax}_x f(x)$
- Idea: Assume f is a sample from a Gaussian Process!
→ Gaussian Process optimization
(a.k.a.: Response surface optimization)

Upper confidence bound approach

- $UCB(x | D) = \underbrace{\mu(x | D)} + \underbrace{2^* \sigma(x | D)}$
- Pick point $x^* = \operatorname{argmax}_x UCB(x | D)$



Matlab demo

Properties

- Implicitly trades off exploration and exploitation
- Exploits prior knowledge about function
- Can converge to optimal solution very quickly! 😊
- Seems to work well in many applications

- Can perform poorly if our prior assumptions are wrong 😞

What you need to know

- GPs =
 - Nonparametric generalization of linear regression
 - Flexible ways to encode prior assumptions about mean payoffs
- Definition of GPs
- Properties of multivariate Gaussians (marginalization, conditioning)
- Gaussian Process optimization
 - Combination of regression and optimization
 - Use confidence bands for selecting samples