

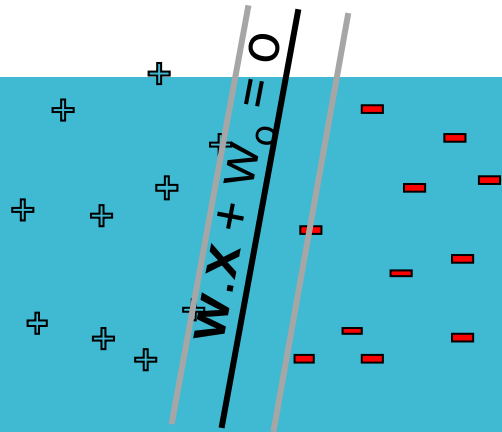
# 10-701: Introduction to Machine Learning Lecture 25 – Kernels & Gaussian Processes

Pradeep Ravikumar & Geoff Gordon

Spring 2026

# Front Matter

- Announcements:
  - HW6
    - due on Wednesday, April 22
    - late due date on Friday, April 24
  - Project
    - Video Showcase due on Wednesday, April 22
    - Final report due on Thursday, April 23
    - Project video showcase this Friday, April 24! Attendance required, refreshments will be served!
    - As a reminder, no late days can be used for project deliverables.
- Recommended Readings:
  - Murphy, Chapters 14.1-14.5



# Recap: SVMs, Primal-Dual Optimization

$$\begin{aligned} &\text{minimize } \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ &\text{subject to } y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1 \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D} \end{aligned} \quad \left. \vphantom{\begin{aligned} &\text{minimize } \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ &\text{subject to } y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1 \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D} \end{aligned}} \right\} \text{Primal}$$

$\Leftrightarrow$

$$\begin{aligned} &\text{maximize } -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} \mathbf{x}^{(i)T} \mathbf{x}^{(j)} + \sum_{i=1}^N \alpha^{(i)} \\ &\text{subject to } \sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0 \\ &\quad \alpha^{(i)} \geq 0 \quad \forall i \in \{1, \dots, N\} \end{aligned} \quad \left. \vphantom{\begin{aligned} &\text{maximize } -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} \mathbf{x}^{(i)T} \mathbf{x}^{(j)} + \sum_{i=1}^N \alpha^{(i)} \\ &\text{subject to } \sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0 \\ &\quad \alpha^{(i)} \geq 0 \quad \forall i \in \{1, \dots, N\} \end{aligned}} \right\} \text{Dual}$$

# Recap: SVMs, Primal-Dual Optimization

- Primal
  - Directly returns the weights,  $[\hat{w}_0, \hat{w}]$
  - Support vectors are all  $(\mathbf{x}^{(s)}, y^{(s)}) \in \mathcal{D}$  s.t.

$$y^{(s)} (\hat{w}^T \mathbf{x}^{(s)} + \hat{w}_0) = 1$$

- Dual
  - Returns the vector,  $\hat{\alpha}$

$$\hat{w} = \sum_{i=1}^N \hat{\alpha}^{(i)} y^{(i)} \mathbf{x}^{(i)}$$

$$\hat{w}_0 = y^{(s)} - \hat{w}^T \mathbf{x}^{(s)} \text{ for any } s \text{ s.t. } \hat{\alpha}^{(s)} > 0$$

- Support vectors are all  $(\mathbf{x}^{(s)}, y^{(s)}) \in \mathcal{D}$  s.t.  $\hat{\alpha}^{(s)} > 0$

# Primal-Dual Optimization

- Primal

- $\hat{y} = \text{sign}(\hat{\mathbf{w}}^T \vec{x} + \hat{w}_0)$

- Dual

- $\hat{y} = \text{sign}(\hat{\mathbf{w}}^T \vec{x} + \hat{w}_0)$

$$= \text{sign} \left( \left( \sum_{i=1}^N \hat{\alpha}^{(i)} y^{(i)} \mathbf{x}^{(i)} \right)^T \mathbf{x} + \hat{w}_0 \right)$$

$$= \text{sign} \left( \sum_{i: \hat{\alpha}^{(i)} > 0} \hat{\alpha}^{(i)} y^{(i)} \mathbf{x}^{(i)T} \mathbf{x} + \hat{w}_0 \right)$$

# Primal-Dual Soft-Margin SVMs

$$\begin{aligned} &\text{minimize} && \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi^{(i)} \\ &\text{subject to} && y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1 - \xi^{(i)} \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D} \\ &&& \xi^{(i)} \geq 0 \quad \forall i \in \{1, \dots, N\} \end{aligned} \quad \left. \vphantom{\begin{aligned} &\text{minimize} \\ &\text{subject to} \end{aligned}} \right\} \text{Primal}$$

$$\begin{aligned} &\text{minimize} && \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} \mathbf{x}^{(i)T} \mathbf{x}^{(j)} - \sum_{i=1}^N \alpha^{(i)} \\ &\text{subject to} && \sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0 \\ &&& 0 \leq \alpha^{(i)} \leq C \quad \forall i \in \{1, \dots, N\} \end{aligned} \quad \left. \vphantom{\begin{aligned} &\text{minimize} \\ &\text{subject to} \end{aligned}} \right\} \text{Dual}$$

# Primal-Dual Soft-Margin SVMs

$$\begin{aligned} &\text{minimize} && \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi^{(i)} \\ &\text{subject to} && y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1 - \xi^{(i)} \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D} \\ &&& \xi^{(i)} \geq 0 \quad \forall i \in \{1, \dots, N\} \end{aligned} \quad \left. \vphantom{\begin{aligned} &\text{minimize} \\ &\text{subject to} \end{aligned}} \right\} \text{Primal}$$

$$\begin{aligned} &\text{minimize} && \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} \mathbf{x}^{(i)T} \mathbf{x}^{(j)} - \sum_{i=1}^N \alpha^{(i)} \\ &\text{subject to} && \sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0 \\ &&& 0 \leq \alpha^{(i)} \leq C \quad \forall i \in \{1, \dots, N\} \end{aligned} \quad \left. \vphantom{\begin{aligned} &\text{minimize} \\ &\text{subject to} \end{aligned}} \right\} \text{Dual}$$

## Recall: Nonlinear Transforms

- Decide on some transformation  $\Phi: \mathcal{X} \rightarrow \mathcal{Z}$
- Given  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$ , learn a hypothesis,  $\tilde{h}(\mathbf{z})$ ,  
using  $\tilde{\mathcal{D}} = \{(\mathbf{z}^{(i)} = \Phi(\mathbf{x}^{(i)}), y^{(i)})\}_{i=1}^N$
- Return the corresponding predictor in the original space:  
 $h(\mathbf{x}) = \tilde{h}(\Phi(\mathbf{x}))$

# Nonlinear SVMs

- Decide on some transformation  $\Phi: \mathcal{X} \rightarrow \mathcal{Z}$
- Find a maximal-margin separating hyperplane in the transformed space,  $[\tilde{\mathbf{w}}, \tilde{w}_0]$ , by solving the QP:

$$\text{minimize } \frac{1}{2} \tilde{\mathbf{w}}^T \tilde{\mathbf{w}}$$

$$\text{subject to } y^{(i)} (\tilde{\mathbf{w}}^T \Phi(\mathbf{x}^{(i)}) + \tilde{w}_0) \geq 1 \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}$$

- Return the corresponding predictor in the original space:

$$h(\mathbf{x}) = \text{sign}(\tilde{\mathbf{w}}^T \Phi(\mathbf{x}) + \tilde{w}_0)$$

# Nonlinear Dual SVMs

- Decide on some transformation  $\Phi: \mathcal{X} \rightarrow \mathcal{Z}$
- Find a maximal-margin separating hyperplane in the transformed space by solving the QP:

$$\text{minimize } \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} \Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x}^{(j)}) - \sum_{i=1}^N \alpha^{(i)}$$

$$\text{subject to } \sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0$$

$$0 \leq \alpha^{(i)} \leq C \quad \forall i \in \{1, \dots, N\}$$

- Return the corresponding predictor in the original space:

$$h(\mathbf{x}) = \text{sign} \left( \sum_{i: \hat{\alpha}^{(i)} > 0} \hat{\alpha}^{(i)} y^{(i)} \Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x}) + \widetilde{w}_0 \right)$$

# Efficiency

- Depending on the transformation  $\Phi$  and the dimensionality of the original input space  $\mathcal{X}$ , computing  $\Phi(\mathbf{x})$  can be prohibitively computationally expensive
  - Computing  $\Phi_2(\mathbf{x}) = [x_1, x_2, \dots, x_D, x_1^2, x_1x_2, \dots, x_D^2]$  requires  $D + \binom{D}{2} + D = \frac{D^2+3D}{2} = O(D^2)$  time
  - Computing  $\Phi_{10}(\mathbf{x})$  requires  $O(D^{10})$  time
- Tradeoff:
  - High-dimensional transformations can result in good hypotheses (as long as they don't overfit)
  - High-dimensional transformations are expensive

# Nonlinear Dual SVMs

- Insight:  $\Phi$  only appears in inner products!
- Find a maximal-margin separating hyperplane in the transformed space by solving the QP:

$$\text{minimize } \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} \Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x}^{(j)}) - \sum_{i=1}^N \alpha^{(i)}$$

$$\text{subject to } \sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0$$

$$0 \leq \alpha^{(i)} \leq C \quad \forall i \in \{1, \dots, N\}$$

- Return the corresponding predictor in the original space:

$$h(\mathbf{x}) = \text{sign} \left( \sum_{i: \hat{\alpha}^{(i)} > 0} \hat{\alpha}^{(i)} y^{(i)} \Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x}) + \widetilde{w}_0 \right)$$

# Nonlinear Dual SVMs

- Insight:  $\Phi$  only appears in inner products!
- Find a maximal-margin separating hyperplane in the transformed space by solving the QP:

$$\text{minimize } \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} \Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x}^{(j)}) - \sum_{i=1}^N \alpha^{(i)}$$

$$\text{subject to } \sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0$$

$$0 \leq \alpha^{(i)} \leq C \quad \forall i \in \{1, \dots, N\}$$

- Return the corresponding predictor in the original space:

$$h(\mathbf{x}) = \text{sign} \left( \sum_{i: \hat{\alpha}^{(i)} > 0} \hat{\alpha}^{(i)} y^{(i)} \Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x}) + \widetilde{w}_0 \right)$$

# The Kernel Trick

- Approach: instead of computing  $\Phi(\mathbf{x})$ , find some function  $K_\Phi$  s.t.  $K_\Phi(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x})^T \Phi(\mathbf{x}') \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}$ 
  - $K_\Phi(\mathbf{x}, \mathbf{x}')$  should be cheaper to compute than  $\Phi(\mathbf{x})$
- Example:  $\Phi'_2(\mathbf{x}) = [x_1, \dots, x_D, x_1^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_{D-1}x_D, x_D^2]$ 
$$\begin{aligned}\Phi'_2(\mathbf{x})^T \Phi'_2(\mathbf{x}') &= \sum_{i=1}^D x_i x'_i + \sum_{i=1}^D x_i^2 x_i'^2 + \sum_{i=1}^D \sum_{j>i} 2x_i x'_i x_j x'_j \\ &= \sum_{i=1}^D x_i x'_i + \left( \sum_{i=1}^D x_i x'_i \right)^2 = \mathbf{x}^T \mathbf{x}' + (\mathbf{x}^T \mathbf{x}')^2\end{aligned}$$
$$K_{\Phi'_2}(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}' + (\mathbf{x}^T \mathbf{x}')^2$$
- Computing  $\Phi'_2(\mathbf{x})^T \Phi'_2(\mathbf{x}')$  requires  $O(D^2)$  time whereas computing  $K_{\Phi'_2}(\mathbf{x}, \mathbf{x}')$  only takes  $O(D)$ !

# Common Kernels

- $K_{\Phi'_2}(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}' + (\mathbf{x}^T \mathbf{x}')^2$

- Implied feature transformation:

$$\Phi'_2(\mathbf{x}) = [x_1, \dots, x_D, x_1^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_{D-1}x_D, x_D^2]$$

- Implied dimensionality:  $\frac{D^2+3D}{2}$

- $K_{\Phi_2^{(\gamma)}}(\mathbf{x}, \mathbf{x}') = (1 + \gamma \mathbf{x}^T \mathbf{x}')^2 - 1$

- Implied feature transformation:

$$\Phi_2^{(\gamma)}(\mathbf{x}) = [\sqrt{2\gamma}x_1, \dots, \sqrt{2\gamma}x_D, \gamma x_1^2, \gamma x_1x_2, \dots, \gamma x_D^2]$$

- $\gamma$  affects the geometry of the transform

- Implied dimensionality:  $\frac{D^2+3D}{2}$

# Common Kernels

- Polynomial Kernel:  $K_{\Phi_Q^{(\gamma)}}(\mathbf{x}, \mathbf{x}') = (1 + \gamma \mathbf{x}^T \mathbf{x}')^Q - 1$ 
  - Implied dimensionality:  $O(D^Q)$
  - $\gamma$  affects the geometry of the transform
- Gaussian-RBF Kernel:  $K_{\Phi_r}(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x}-\mathbf{x}'\|^2}{2r}}$ 
  - Implied feature transformation:  $\Phi_r(\mathbf{x}) = \left[ \left[ e^{-\frac{x_1^2}{2r}} \sqrt{\frac{(x_1^d)^2}{d!r^d}}, \dots, e^{-\frac{x_D^2}{2r}} \sqrt{\frac{(x_1^d)^2}{d!r^d}} \right] : d \in \mathbb{N} \right]$
  - Implied dimensionality:  $\infty!$

# Nonlinear Dual SVMs

- Decide on a (valid) kernel function  $K_{\Phi}$
- Find a maximal-margin separating hyperplane in the transformed space by solving the QP:

$$\text{minimize } \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} K_{\Phi}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) - \sum_{i=1}^N \alpha^{(i)}$$

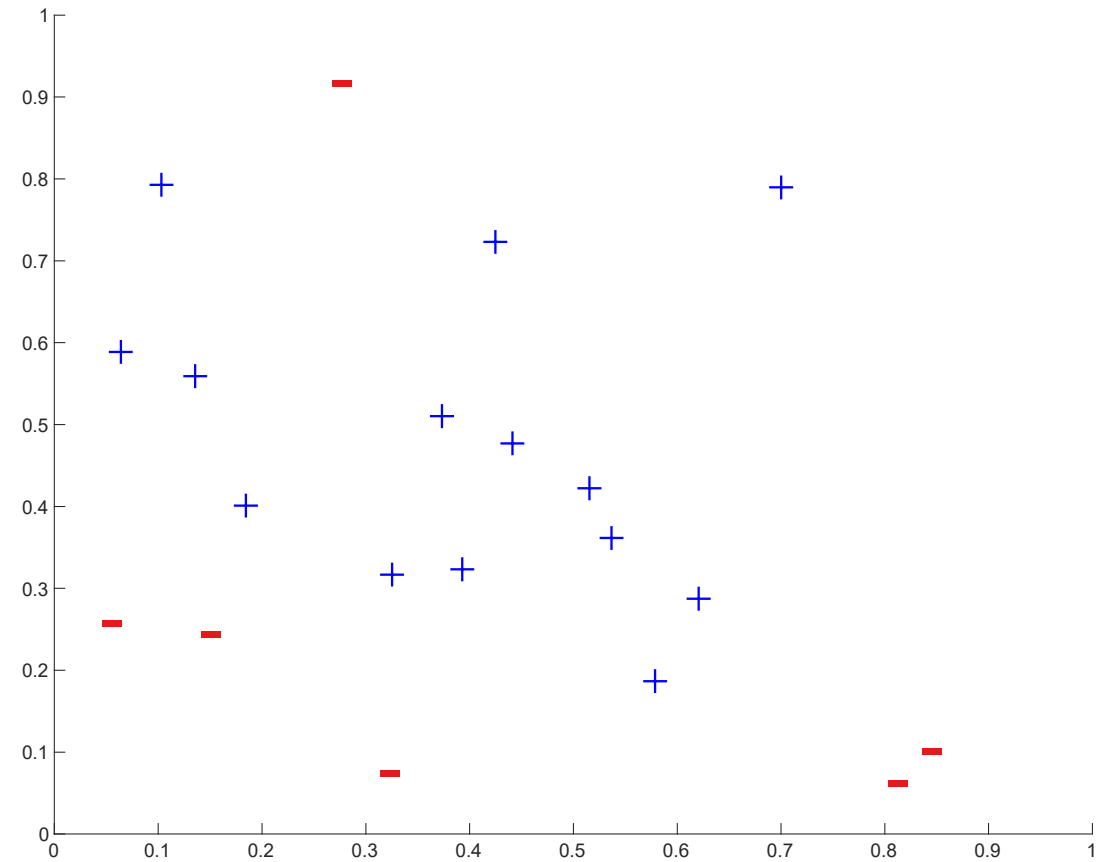
$$\text{subject to } \sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0$$

$$0 \leq \alpha^{(i)} \quad \forall i \in \{1, \dots, N\}$$

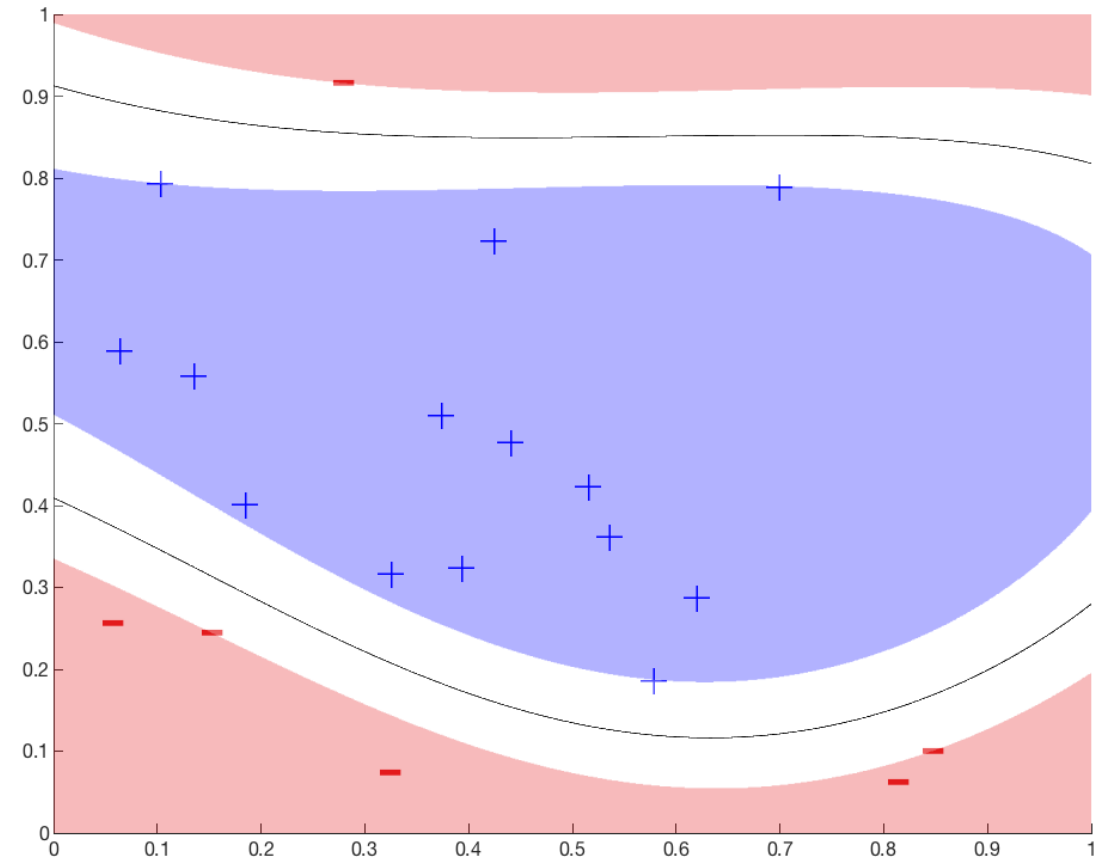
- Return the corresponding predictor in the original space:

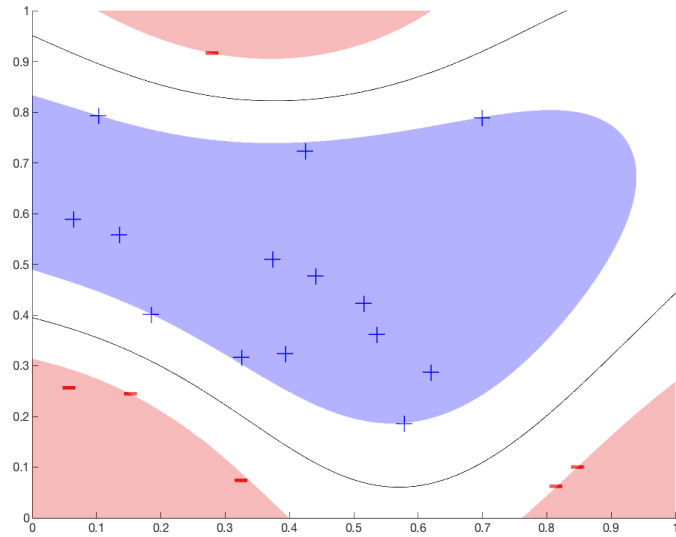
$$h(\mathbf{x}) = \text{sign} \left( \sum_{i: \hat{\alpha}^{(i)} > 0} \hat{\alpha}^{(i)} y^{(i)} K_{\Phi}(\mathbf{x}^{(i)}, \mathbf{x}) + \widetilde{w}_0 \right)$$

# Gaussian- RBF Kernel

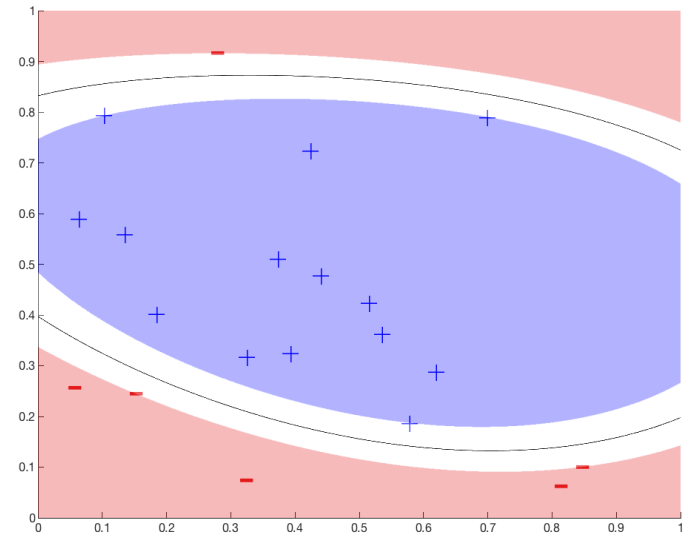
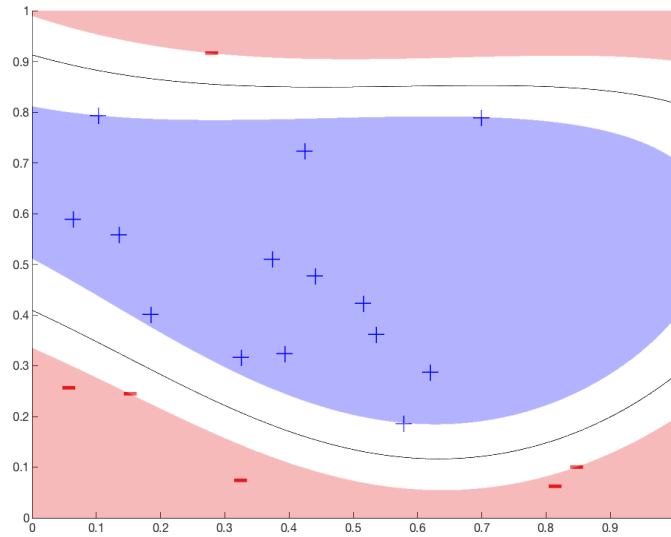


# Gaussian- RBF Kernel





Smaller  $r$



Larger  $r$

# Gaussian-RBF Kernel

# Nonlinear Dual Soft-Margin SVMs

- Decide on a (valid) kernel function  $K_{\Phi}$
- Find a maximal-margin separating hyperplane in the transformed space by solving the QP:

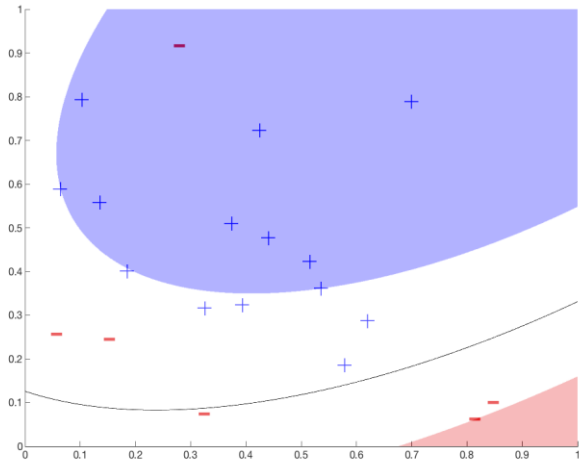
$$\text{minimize } \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} K_{\Phi}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) - \sum_{i=1}^N \alpha^{(i)}$$

$$\text{subject to } \sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0$$

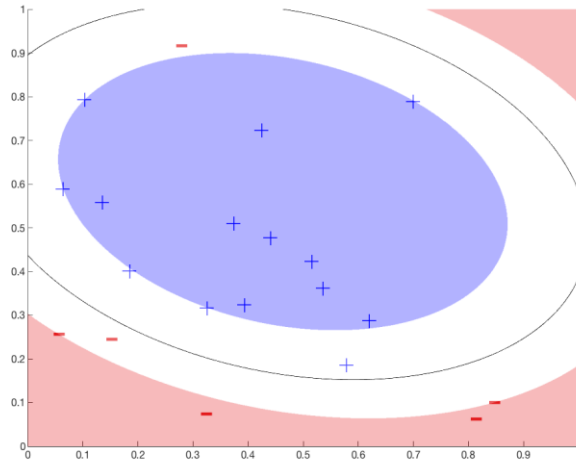
$$0 \leq \alpha^{(i)} \leq C \quad \forall i \in \{1, \dots, N\}$$

- Return the corresponding predictor in the original space:

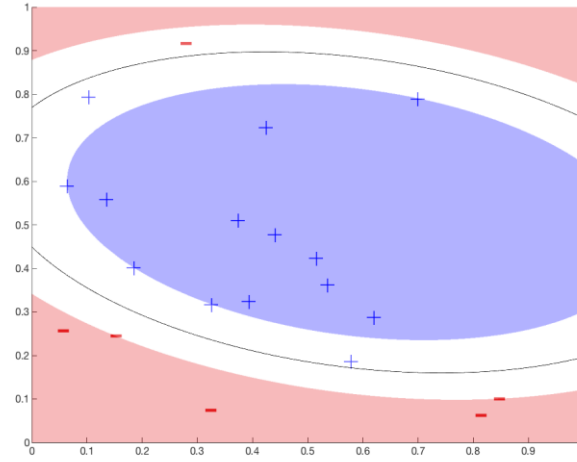
$$h(\mathbf{x}) = \text{sign} \left( \sum_{i: \hat{\alpha}^{(i)} > 0} \hat{\alpha}^{(i)} y^{(i)} K_{\Phi}(\mathbf{x}^{(i)}, \mathbf{x}) + \widetilde{w}_0 \right)$$



Smaller  $C$



Larger  $C$



Hard Margin

## 2<sup>nd</sup>-Degree Polynomial Kernel

$C$  is a trade-off parameter between the size of the margin and the soft training error

# Myth Busters

## MythBusters #1

“The kernel trick means there are no features anymore.”



**Myth: kernels replace feature maps with pure algebraic magic.**



**Reality: a kernel is still an inner product in some transformed feature space  $\Phi(x)$ . The trick is that we compute it implicitly.**

explicit feature map

$\mathbf{x} \rightarrow \Phi(\mathbf{x})$   
then compute  $\Phi(\mathbf{x})^T \Phi(\mathbf{x}')$



kernel trick

$\mathbf{K}(\mathbf{x}, \mathbf{x}')$   
 $= \Phi(\mathbf{x})^T \Phi(\mathbf{x}')$

You keep the transformed-space geometry without writing down all transformed coordinates.

# Valid Kernels

- Any function  $K$  is a valid kernel if and only if:
  - $\exists$  a transformation  $\Phi$  s.t.

$$K(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x})^T \Phi(\mathbf{x}') \quad \forall \mathbf{x}, \mathbf{x}'$$



- the Gram matrix

$$K = \begin{bmatrix} K(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & K(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) & \dots & K(\mathbf{x}^{(1)}, \mathbf{x}^{(N)}) \\ K(\mathbf{x}^{(2)}, \mathbf{x}^{(1)}) & K(\mathbf{x}^{(2)}, \mathbf{x}^{(2)}) & \dots & K(\mathbf{x}^{(2)}, \mathbf{x}^{(N)}) \\ \vdots & \vdots & \ddots & \vdots \\ K(\mathbf{x}^{(N)}, \mathbf{x}^{(1)}) & K(\mathbf{x}^{(N)}, \mathbf{x}^{(2)}) & \dots & K(\mathbf{x}^{(N)}, \mathbf{x}^{(N)}) \end{bmatrix}$$

is symmetric and positive semi-definite  $\forall$  sets

$$\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$$

Recall:  $K$  is positive semi-definite iff  $\mathbf{x}^T K \mathbf{x} \geq 0$  for all  $\mathbf{x}$

## Building New Kernels

- For any valid kernels  $K_1, K_2$  with implied feature transformations  $\Phi_1, \Phi_2$  and non-negative coefficients  $c_1, c_2$ , the following are all valid kernels:

1.  $K(\mathbf{x}, \mathbf{x}') = c_1 K_1(\mathbf{x}, \mathbf{x}') + c_2 K_2(\mathbf{x}, \mathbf{x}')$

$$\Phi(\mathbf{x}) = [\sqrt{c_1} \Phi_1(\mathbf{x}), \sqrt{c_2} \Phi_2(\mathbf{x})]$$

2.  $K(\mathbf{x}, \mathbf{x}') = c_1 K_1(\mathbf{x}, \mathbf{x}') K_2(\mathbf{x}, \mathbf{x}')$

$$\Phi(\mathbf{x}) = \left[ \left\{ \sqrt{c_1} \phi_i(\mathbf{x}) \phi_j(\mathbf{x}) \right\}_{\phi_i(\mathbf{x}) \in \Phi_1(\mathbf{x}), \phi_j(\mathbf{x}) \in \Phi_2(\mathbf{x})} \right]$$

3.  $K(\mathbf{x}, \mathbf{x}') = e^{K_1(\mathbf{x}, \mathbf{x}')}$

Taylor series:  $e^{K_1(\mathbf{x}, \mathbf{x}')} = 1 + K_1(\mathbf{x}, \mathbf{x}') + \frac{K_1(\mathbf{x}, \mathbf{x}')^2}{2!} + \frac{K_1(\mathbf{x}, \mathbf{x}')^3}{3!} + \dots$

# Key Takeaways

- Kernels and the “kernel trick” allow for efficient use of feature transformations for inner product methods
  - Definition of valid kernels
  - Common kernels and combining kernels

# Kernels Everywhere!

- Any method that only depends on the Euclidean distance between data points is an inner product method:

$$\|\mathbf{x} - \mathbf{x}'\|_2 = \sqrt{(\mathbf{x} - \mathbf{x}')^T (\mathbf{x} - \mathbf{x}')} = \sqrt{\mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T \mathbf{x}' + \mathbf{x}'^T \mathbf{x}'}$$

- We can kernelize  $k$ NN!
- We can also kernelize logistic/linear/ridge regression!

# Kernels in Logistic Regression

$$P(Y = 1 | x, \mathbf{w}) = \frac{1}{1 + e^{-(\mathbf{w} \cdot \Phi(\mathbf{x}) + b)}}$$

- Define weights in terms of features:

$$\mathbf{w} = \sum_i \alpha_i \Phi(\mathbf{x}_i)$$

$$\begin{aligned} P(Y = 1 | x, \mathbf{w}) &= \frac{1}{1 + e^{-(\sum_i \alpha_i \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}) + b)}} \\ &= \frac{1}{1 + e^{-(\sum_i \alpha_i K(\mathbf{x}, \mathbf{x}_i) + b)}} \end{aligned}$$

- Derive simple gradient descent rule on  $\alpha_i$

# Ridge regression

$$\min_{\beta} \sum_{i=1}^n (Y_i - X_i\beta)^2 + \lambda \|\beta\|_2^2 \quad \hat{\beta} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$

Recall

$$\mathbf{A} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix} = \begin{bmatrix} X_1^{(1)} & \dots & X_1^{(p)} \\ \vdots & \ddots & \vdots \\ X_n^{(1)} & \dots & X_n^{(p)} \end{bmatrix}$$

Hence  $\mathbf{A}^T \mathbf{A}$  is a  $p \times p$  matrix whose entries denote the (sample) correlation between the features

NOT inner products between the data points – the inner product matrix would be  $\mathbf{A} \mathbf{A}^T$  which is  $n \times n$  (also known as Gram matrix)

Using dual formulation, we will write the solution in terms of  $\mathbf{A} \mathbf{A}^T$

# Ridge regression

$$\min_{\beta} \sum_{i=1}^n (Y_i - X_i\beta)^2 + \lambda \|\beta\|_2^2$$

$$\hat{\beta} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$

## Similarity with SVMs

Primal problem:

$$\begin{aligned} \min_{\beta, z_i} \quad & \sum_{i=1}^n z_i^2 + \lambda \|\beta\|_2^2 \\ \text{s.t.} \quad & z_i = Y_i - X_i\beta \end{aligned}$$

SVM Primal problem:

$$\min_{w, \xi_i} C \sum_{i=1}^n \xi_i + \frac{1}{2} \|w\|_2^2$$

$$\text{s.t. } \xi_i = \max(1 - Y_i X_i \cdot w, 0)$$

Lagrangian:

$$\sum_{i=1}^n z_i^2 + \lambda \|\beta\|_2^2 + \sum_{i=1}^n \alpha_i (z_i - Y_i + X_i\beta)$$

$\alpha_i$  – Lagrange parameter, one per training point

# Ridge regression (dual)

$$\min_{\beta} \sum_{i=1}^n (Y_i - X_i\beta)^2 + \lambda \|\beta\|_2^2 \quad \hat{\beta} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$

Dual problem:

$$\max_{\alpha} \min_{\beta, z_i} \sum_{i=1}^n z_i^2 + \lambda \|\beta\|_2^2 + \sum_{i=1}^n \alpha_i (z_i - Y_i + X_i\beta)$$

$\alpha = \{\alpha_i\}$  for  $i = 1, \dots, n$

Taking derivatives of Lagrangian wrt  $\beta$  and  $z_i$  we get:

$$\beta = -\frac{1}{2\lambda} \mathbf{A}^T \alpha \quad z_i = -\frac{\alpha_i}{2}$$

$$\text{Dual problem: } \max_{\alpha} -\frac{\alpha^T \alpha}{4} - \frac{1}{4\lambda} \alpha^T \mathbf{A} \mathbf{A}^T \alpha - \alpha^T \mathbf{Y}$$

n-dimensional optimization problem

# Ridge regression (dual)

$$\min_{\beta} \sum_{i=1}^n (Y_i - X_i\beta)^2 + \lambda \|\beta\|_2^2 \quad \hat{\beta} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$
$$= \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{Y}$$

Dual problem:

$$\max_{\alpha} -\frac{\alpha^T \alpha}{4} - \frac{1}{4\lambda} \alpha^T \mathbf{A} \mathbf{A}^T \alpha - \alpha^T \mathbf{Y} \quad \Rightarrow \hat{\alpha} = -\left( \frac{\mathbf{A} \mathbf{A}^T}{\lambda} + \mathbf{I} \right)^{-1} \mathbf{Y}$$

can get back  $\hat{\beta} = -\frac{1}{2\lambda} \mathbf{A}^T \hat{\alpha} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{Y}$

Weighted average of training points

Weight of each training point (but typically not sparse)

# Kernelized ridge regression

$$\hat{\beta} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$

Using dual, can re-write solution as:

$$\hat{\beta} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{Y}$$

How does this help?

- Only need to invert  $n \times n$  matrix (instead of  $p \times p$  or  $m \times m$ )
- More importantly, kernel trick!

$\mathbf{A} \mathbf{A}^T$  involves only inner products between the training points  
BUT still have an extra  $\mathbf{A}^T$

Recall the predicted label is  $\hat{f}_n(X) = \mathbf{X} \hat{\beta}$

$$= \mathbf{X} \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{Y}$$

$\mathbf{X} \mathbf{A}^T$  contains inner products between test point  $\mathbf{X}$  and training points!

# Kernelized ridge regression

$$\hat{\beta} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$

$$\hat{f}_n(X) = \mathbf{X} \hat{\beta}$$

Using dual, can re-write solution as:

$$\hat{\beta} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{Y}$$

How does this help?

- Only need to invert  $n \times n$  matrix (instead of  $p \times p$  or  $m \times m$ )
- More importantly, kernel trick!

$$\hat{f}_n(X) = \mathbf{K}_X (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{Y}$$

$$\mathbf{K}_X(i) = \phi(X) \cdot \phi(X_i)$$

$$\mathbf{K}(i, j) = \phi(X_i) \cdot \phi(X_j)$$

Work with kernels, never need to write out the high-dim vectors

# Kernelized ridge regression

$$\hat{f}_n(X) = \mathbf{K}_X (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{Y}$$

$$\mathbf{K}_X(i) = \boldsymbol{\phi}(X) \cdot \boldsymbol{\phi}(X_i)$$

$$\mathbf{K}(i, j) = \boldsymbol{\phi}(X_i) \cdot \boldsymbol{\phi}(X_j)$$

Work with kernels, never need to write out the high-dim vectors

Examples of kernels:

Polynomials of degree exactly d  $K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d$

Polynomials of degree up to d  $K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v} + 1)^d$

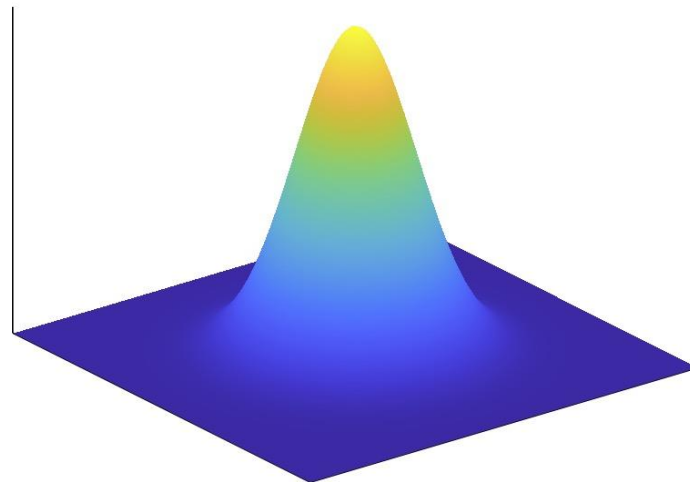
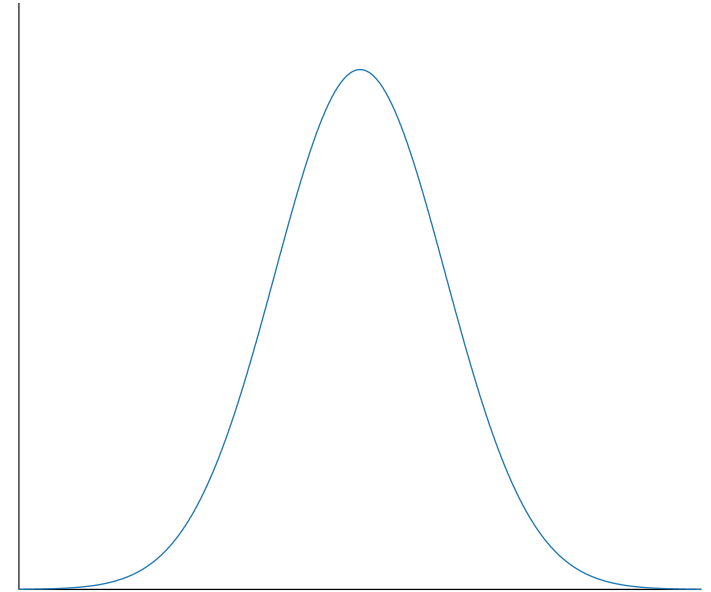
Gaussian/Radial kernels  $K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma^2}\right)$

Ridge Regression with (implicit) nonlinear features  $\boldsymbol{\phi}(X)$ !  $f(X) = \boldsymbol{\phi}(X)\boldsymbol{\beta}$

# Gaussians

(Univariate) Gaussians:

$$x \sim \mathcal{N}(x; \mu = 0, \sigma^2 = 1)$$



- Multivariate Gaussians:

$$\mathbf{x} = [x_1, \dots, x_D]^T$$

$$\sim \mathcal{N}(\mathbf{x}; \boldsymbol{\mu} = \mathbf{0}_D, \boldsymbol{\Sigma} = I_D)$$

# Some fun facts about Gaussians

- Closure under linear transformations:

$$\text{If } \mathbf{x} \sim \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

$$\text{then } A\mathbf{x} + b \sim \mathcal{N}(A\boldsymbol{\mu} + b, A\boldsymbol{\Sigma}A^T)$$

- Closure under addition

$$\text{If } \mathbf{x} \sim \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \text{ and } \mathbf{y} \sim \mathcal{N}(\mathbf{y}; \mathbf{m}, S),$$

$$\text{then } \mathbf{x} + \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu} + \mathbf{m}, \boldsymbol{\Sigma} + S)$$

- Closure under conditioning:

$$\text{If } \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}; \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right),$$

$$\text{then } x_1 | x_2 = c \sim \mathcal{N}(x_1; \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(c - \mu_2), \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})$$

Some old  
friends

Gaussian process =

Bayesian linear regression + Kernels

Some old  
friends

Gaussian process =

Bayesian linear regression + Kernels

## Recall: MAP for Linear Regression

- If we assume a linear model with additive Gaussian noise  $\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon}$  where  $\boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \rightarrow \mathbf{y} \sim N(X\mathbf{w}, \sigma^2 I_N)$  and independent identical Gaussian priors on the weights...

$$\mathbf{w} \sim N\left(\mathbf{0}, \frac{\sigma^2}{\lambda} I_{D+1}\right) \rightarrow p(\mathbf{w}) \propto \exp\left(-\frac{1}{2\sigma^2} (\lambda \mathbf{w}^T \mathbf{w})\right)$$

- ... then, the MAP of  $\mathbf{w}$  is the ridge regression solution!

$$\begin{aligned}\mathbf{w}_{MAP} &= \underset{\mathbf{w}}{\operatorname{argmin}} (X\mathbf{w} - \mathbf{y})^T (X\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w} \\ &= (X^T X + \lambda I_{D+1})^{-1} X^T \mathbf{y}\end{aligned}$$

# Bayesian Linear Regression

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \mathbf{w} \sim N(\mathbf{0}_{D+1}, \Sigma)$$

then,

$$\mathbf{y} \sim N((X\mathbf{0}_{D+1} + \mathbf{0}_N) \equiv \mathbf{0}_N, X\Sigma X^T + \sigma^2 I_N)$$

# Bayesian Linear Regression

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \mathbf{w} \sim N(\mathbf{0}_{D+1}, \Sigma)$$

then,

$$\begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix} \sim N \left( \begin{bmatrix} \mathbf{0}_{D+1} \\ \mathbf{0}_N \end{bmatrix}, \begin{bmatrix} \Sigma & ??? \\ ??? & X\Sigma X^T + \sigma^2 I_N \end{bmatrix} \right)$$

- Covariance between  $\mathbf{y}$  and  $\mathbf{w}$ :

$$\text{Cov}(\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon}, \mathbf{w}) = \text{Cov}(X\mathbf{w}, \mathbf{w}) = X\text{Cov}(\mathbf{w}, \mathbf{w}) = X\Sigma$$

# Bayesian Linear Regression

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \mathbf{w} \sim N(\mathbf{0}_{D+1}, \Sigma)$$

then,

$$\begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix} \sim N \left( \begin{bmatrix} \mathbf{0}_{D+1} \\ \mathbf{0}_N \end{bmatrix}, \begin{bmatrix} \Sigma & X\Sigma \\ X\Sigma & X\Sigma X^T + \sigma^2 I_N \end{bmatrix} \right)$$

- Covariance between  $\mathbf{y}$  and  $\mathbf{w}$ :

$$\text{Cov}(\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon}, \mathbf{w}) = \text{Cov}(X\mathbf{w}, \mathbf{w}) = X\text{Cov}(\mathbf{w}, \mathbf{w}) = X\Sigma$$

# Bayesian Linear Regression

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \mathbf{w} \sim N(\mathbf{0}_{D+1}, \Sigma)$$

then,

$$\mathbf{w} \mid \mathbf{y} \sim N(\boldsymbol{\mu}_{POST}, \Sigma_{POST})$$

where

$$\begin{aligned}\boldsymbol{\mu}_{POST} &= \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} \mathbf{y}, \\ \Sigma_{POST} &= \Sigma - \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} X \Sigma\end{aligned}$$

# Bayesian Linear Regression

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \mathbf{w} \sim N(\mathbf{0}_{D+1}, \Sigma)$$

then given a new test data point  $\mathbf{x}'$ , the prediction is

$$y' \mid \mathbf{y} = \mathbf{x}'^T \mathbf{w} \mid \mathbf{y} \sim N(\mathbf{x}'^T \boldsymbol{\mu}_{POST}, \mathbf{x}'^T \Sigma_{POST} \mathbf{x}')$$

where

$$\begin{aligned} \boldsymbol{\mu}_{POST} &= \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} \mathbf{y}, \\ \Sigma_{POST} &= \Sigma - \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} X \Sigma \end{aligned}$$

# Bayesian Linear Regression

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \mathbf{w} \sim N(\mathbf{0}_{D+1}, \Sigma)$$

then given a new test data point  $\mathbf{x}'$ , the prediction is

$$y' \mid \mathbf{y} = \mathbf{x}'^T \mathbf{w} \mid \mathbf{y} \sim N(\boldsymbol{\mu}_{PRED}, \Sigma_{PRED})$$

where

$$\boldsymbol{\mu}_{PRED} = \mathbf{x}'^T \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} \mathbf{y},$$

$$\Sigma_{PRED} = \mathbf{x}'^T \Sigma \mathbf{x}' - \mathbf{x}'^T \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} X \Sigma \mathbf{x}'$$

Some old  
friends

Gaussian process =

Bayesian linear regression + Kernels

Some old  
friends

Gaussian process =

Bayesian linear regression + **Kernels**

# Bayesian Linear Regression...

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \mathbf{w} \sim N(\mathbf{0}_{D+1}, \Sigma)$$

then given a new test data point  $\mathbf{x}'$ , the prediction is

$$y' \mid \mathbf{y} = \mathbf{x}'^T \mathbf{w} \mid \mathbf{y} \sim N(\boldsymbol{\mu}_{PRED}, \Sigma_{PRED})$$

where

$$\boldsymbol{\mu}_{PRED} = \mathbf{x}'^T \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} \mathbf{y},$$

$$\Sigma_{PRED} = \mathbf{x}'^T \Sigma \mathbf{x}' - \mathbf{x}'^T \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} X \Sigma \mathbf{x}'$$

# Bayesian Linear Regression can be kernelized!

$$\Phi = \begin{bmatrix} 1 & \phi(\mathbf{x}^{(1)})^T \\ 1 & \phi(\mathbf{x}^{(2)})^T \\ \vdots & \vdots \\ 1 & \phi(\mathbf{x}^{(N)})^T \end{bmatrix}$$

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = \Phi\boldsymbol{\omega} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \boldsymbol{\omega} \sim N(\mathbf{0}_{D'+1}, \Sigma)$$

then given a new test data point  $\mathbf{x}'$ , the prediction is

$$y' | \mathbf{y} = \phi(\mathbf{x}')^T \boldsymbol{\omega} | \mathbf{y} \sim N(\boldsymbol{\mu}_{PRED}, \Sigma_{PRED})$$

where

$$\boldsymbol{\mu}_{PRED} = \phi(\mathbf{x}')^T \Sigma \Phi^T (\Phi \Sigma \Phi^T + \sigma^2 I_N)^{-1} \mathbf{y},$$

$$\Sigma_{PRED}$$

$$= \phi(\mathbf{x}')^T \Sigma \phi(\mathbf{x}') - \phi(\mathbf{x}')^T \Sigma \Phi^T (\Phi \Sigma \Phi^T + \sigma^2 I_N)^{-1} \Phi \Sigma \phi(\mathbf{x}')$$

# Bayesian Linear Regression can be kernelized!

$$\Phi = \begin{bmatrix} 1 & \phi(\mathbf{x}^{(1)})^T \\ 1 & \phi(\mathbf{x}^{(2)})^T \\ \vdots & \vdots \\ 1 & \phi(\mathbf{x}^{(N)})^T \end{bmatrix}$$

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = \Phi\boldsymbol{\omega} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \boldsymbol{\omega} \sim N(\mathbf{0}_{D'+1}, \Sigma)$$

then given a new test data point  $\mathbf{x}'$ , the prediction is

$$y' | \mathbf{y} = \phi(\mathbf{x}')^T \boldsymbol{\omega} | \mathbf{y} \sim N(\boldsymbol{\mu}_{PRED}, \Sigma_{PRED})$$

where

$$\boldsymbol{\mu}_{PRED} = \phi(\mathbf{x}')^T \Sigma \Phi^T (\Phi \Sigma \Phi^T + \sigma^2 I_N)^{-1} \mathbf{y},$$

$$\Sigma_{PRED}$$

$$= \phi(\mathbf{x}')^T \Sigma \phi(\mathbf{x}') - \phi(\mathbf{x}')^T \Sigma \Phi^T (\Phi \Sigma \Phi^T + \sigma^2 I_N)^{-1} \Phi \Sigma \phi(\mathbf{x}')$$

- Define the kernel function to be

$$K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \Sigma \phi(\mathbf{x}')$$

# Bayesian Linear Regression can be kernelized!

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = \Phi\boldsymbol{\omega} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \boldsymbol{\omega} \sim N(\mathbf{0}_{D'+1}, \Sigma)$$

then given a new test data point  $\mathbf{x}'$ , the prediction is

$$y' | \mathbf{y} = \phi(\mathbf{x}')^T \boldsymbol{\omega} | \mathbf{y} \sim N(\boldsymbol{\mu}_{PRED}, \Sigma_{PRED})$$

where

$$\boldsymbol{\mu}_{PRED} = K(\mathbf{x}', X)(K(X, X) + \sigma^2 I_N)^{-1} \mathbf{y},$$

$$\Sigma_{PRED} = K(\mathbf{x}', \mathbf{x}') - K(\mathbf{x}', X)(K(X, X) + \sigma^2 I_N)^{-1} K(X, \mathbf{x}')$$

- Define the kernel function to be

$$K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \Sigma \phi(\mathbf{x}')$$

Wait, what happened to the weights?

Recall: models as functions!

- Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\mathbf{y} = \Phi\boldsymbol{\omega} + \boldsymbol{\epsilon} \text{ where } \boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 I_N) \text{ and } \boldsymbol{\omega} \sim N(\mathbf{0}_{D'+1}, \Sigma)$$

then given a new test data point  $\mathbf{x}'$ , the prediction is

$$y' | \mathbf{y} = \phi(\mathbf{x}')^T \boldsymbol{\omega} | \mathbf{y} \sim N(\boldsymbol{\mu}_{PRED}, \Sigma_{PRED})$$

where

$$\boldsymbol{\mu}_{PRED} = K(\mathbf{x}', X)(K(X, X) + \sigma^2 I_N)^{-1} \mathbf{y},$$

$$\Sigma_{PRED} = K(\mathbf{x}', \mathbf{x}') - K(\mathbf{x}', X)(K(X, X) + \sigma^2 I_N)^{-1} K(X, \mathbf{x}')$$

- Define the kernel function to be

$$K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \Sigma \phi(\mathbf{x}')$$

Some old  
friends

Gaussian process =

Bayesian linear regression + Kernels

A new  
perspective

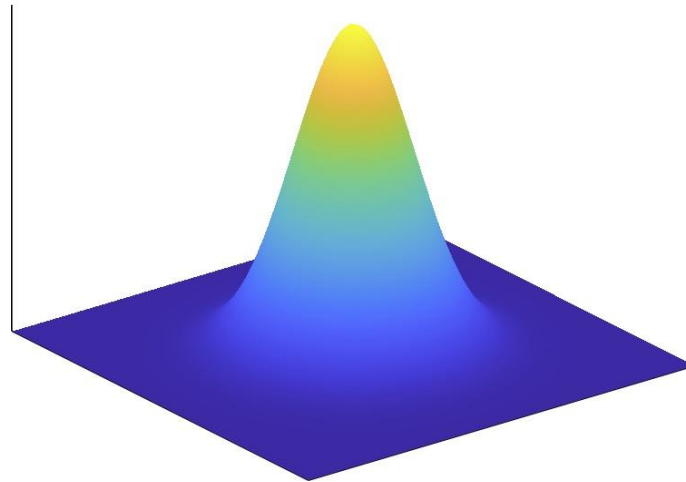
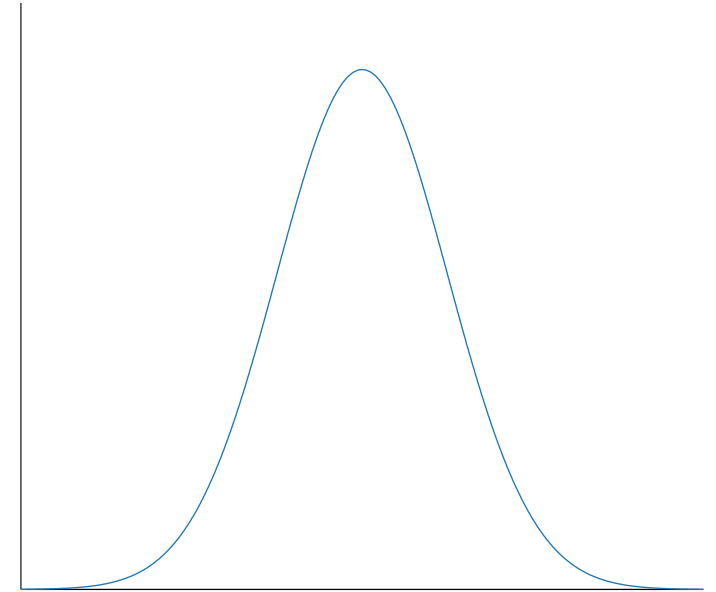
Gaussian process =

The extension of a Gaussian  
distribution to functions

# Gaussians

- (Univariate) Gaussians:

$$x \sim \mathcal{N}(x; \mu = 0, \sigma^2 = 1)$$



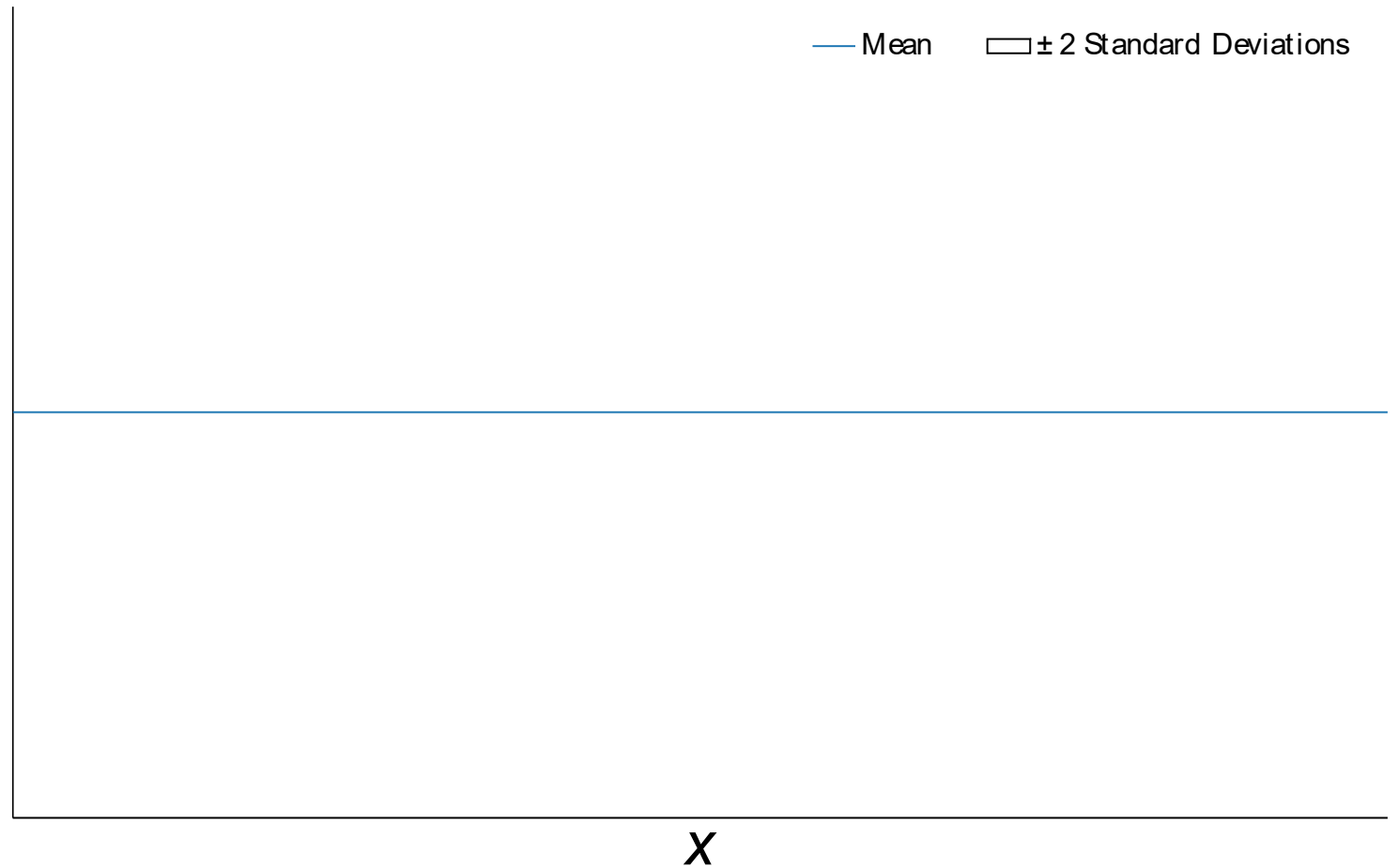
- Multivariate Gaussians:

$$\mathbf{x} = [x_1, \dots, x_D]^T$$

$$\sim \mathcal{N}(\mathbf{x}; \boldsymbol{\mu} = \mathbf{0}_D, \boldsymbol{\Sigma} = I_D)$$

# Gaussian Process (GP)

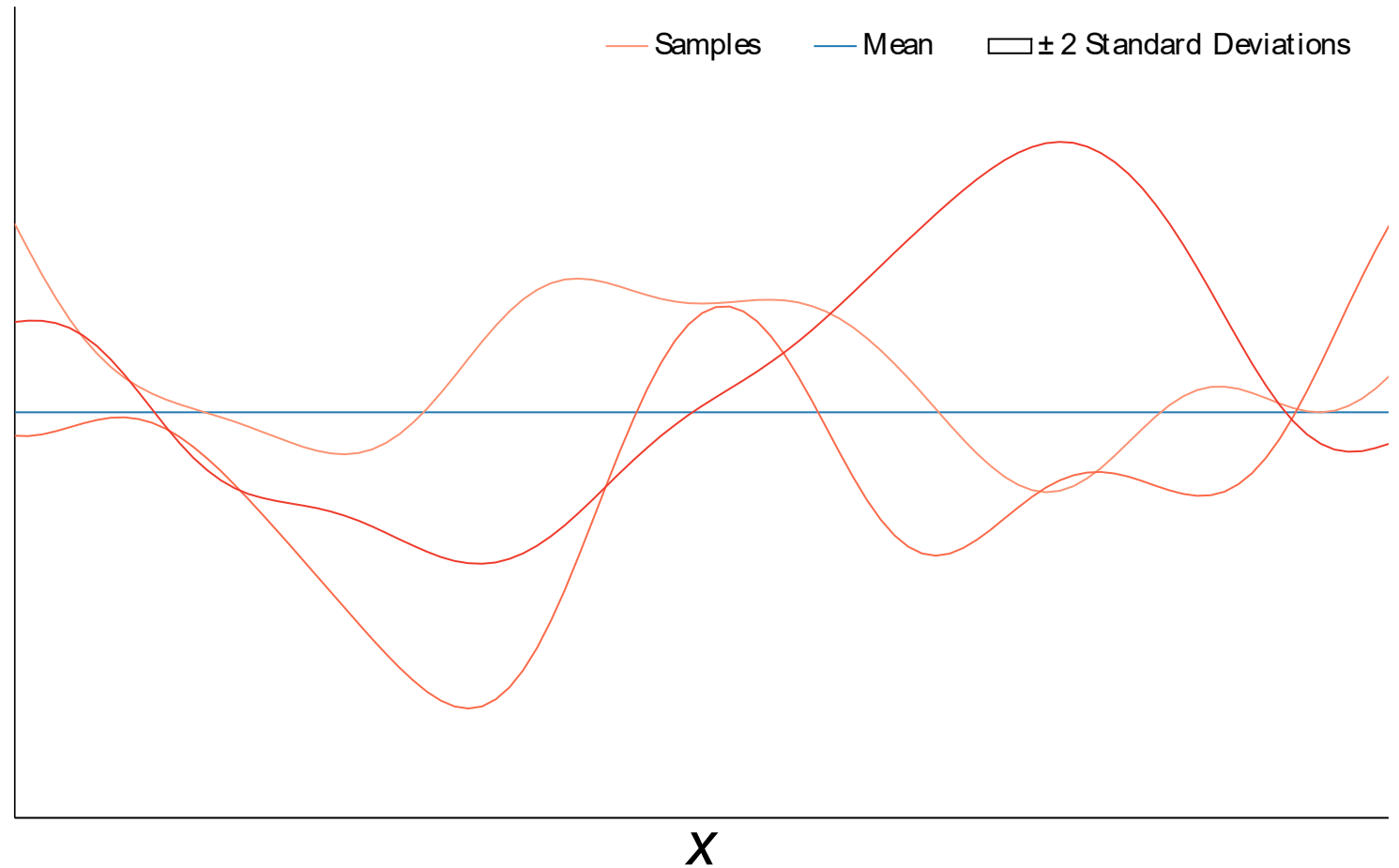
$$f: \mathbb{R}^p \mapsto \mathbb{R} \sim \mathcal{GP}(f; \mu(x), \Sigma(x, x'))$$



$$f \sim \mathcal{GP}(\mu, \Sigma) \rightarrow f(x) \sim \mathcal{N}(\mu(x), \Sigma(x, x))$$

# Gaussian Process (GP)

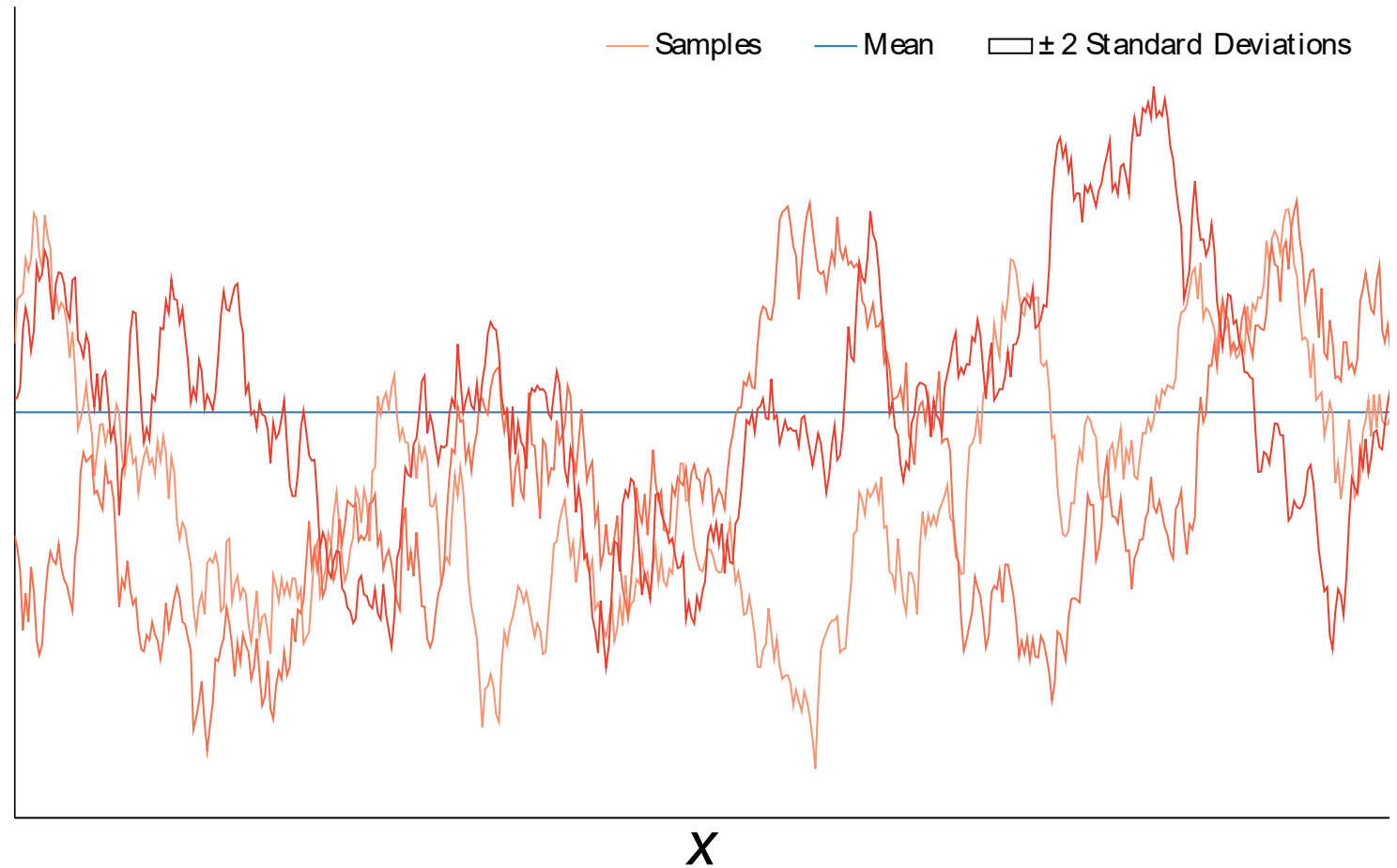
$$f: \mathbb{R}^p \mapsto \mathbb{R} \sim \mathcal{GP}(f; \mu(x) = 0, \Sigma(x, x') = \exp(-(x - x')^2))$$



$$f \sim \mathcal{GP}(\mu, \Sigma) \rightarrow f(x) \sim \mathcal{N}(\mu(x), \Sigma(x, x))$$

# Gaussian Process (GP)

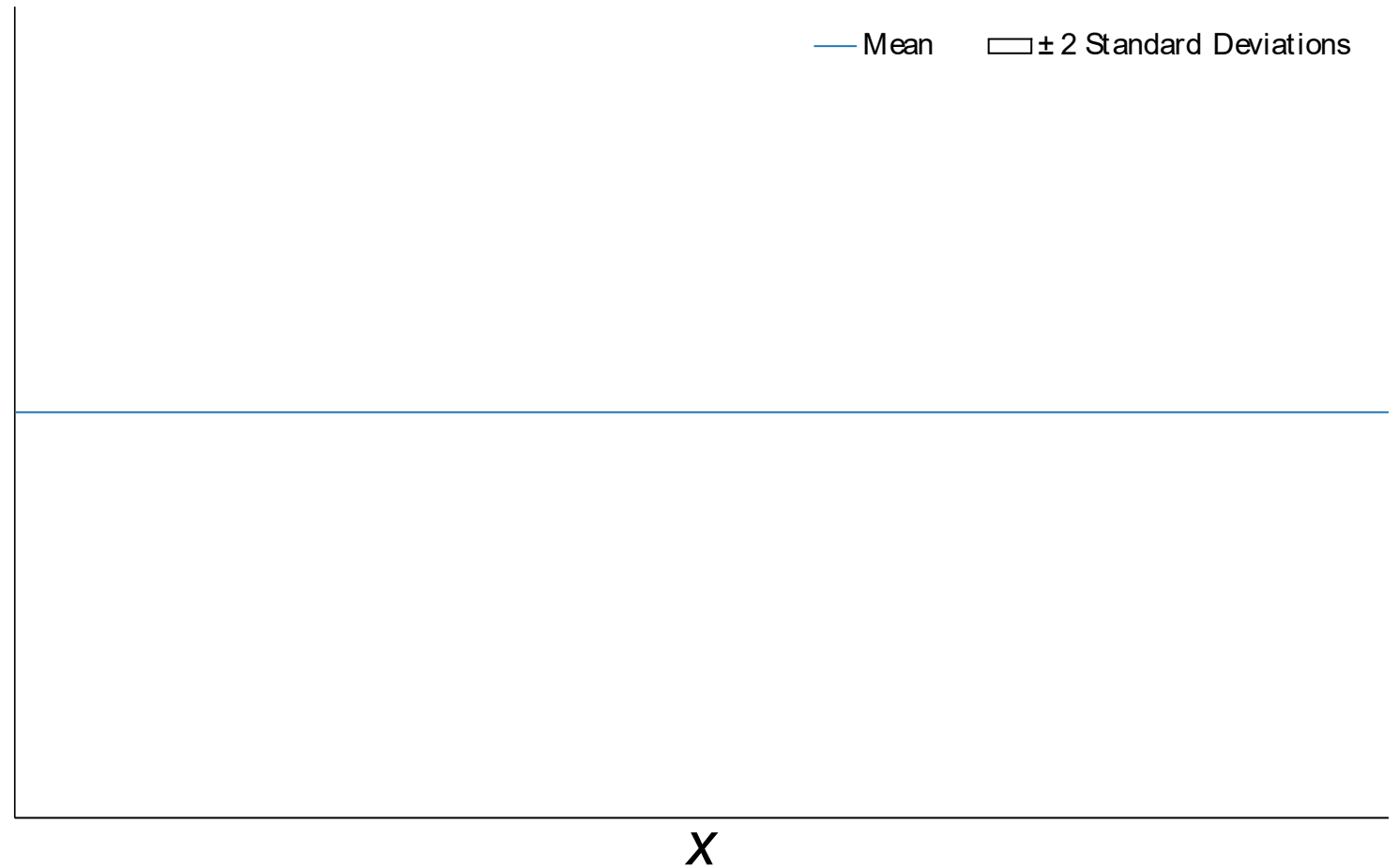
$$f: \mathbb{R}^p \mapsto \mathbb{R} \sim \mathcal{GP}(f; \mu(x) = 0, \Sigma(x, x') = \exp(-|x - x'|))$$



$$f \sim \mathcal{GP}(\mu, \Sigma) \rightarrow f(x) \sim \mathcal{N}(\mu(x), \Sigma(x, x))$$

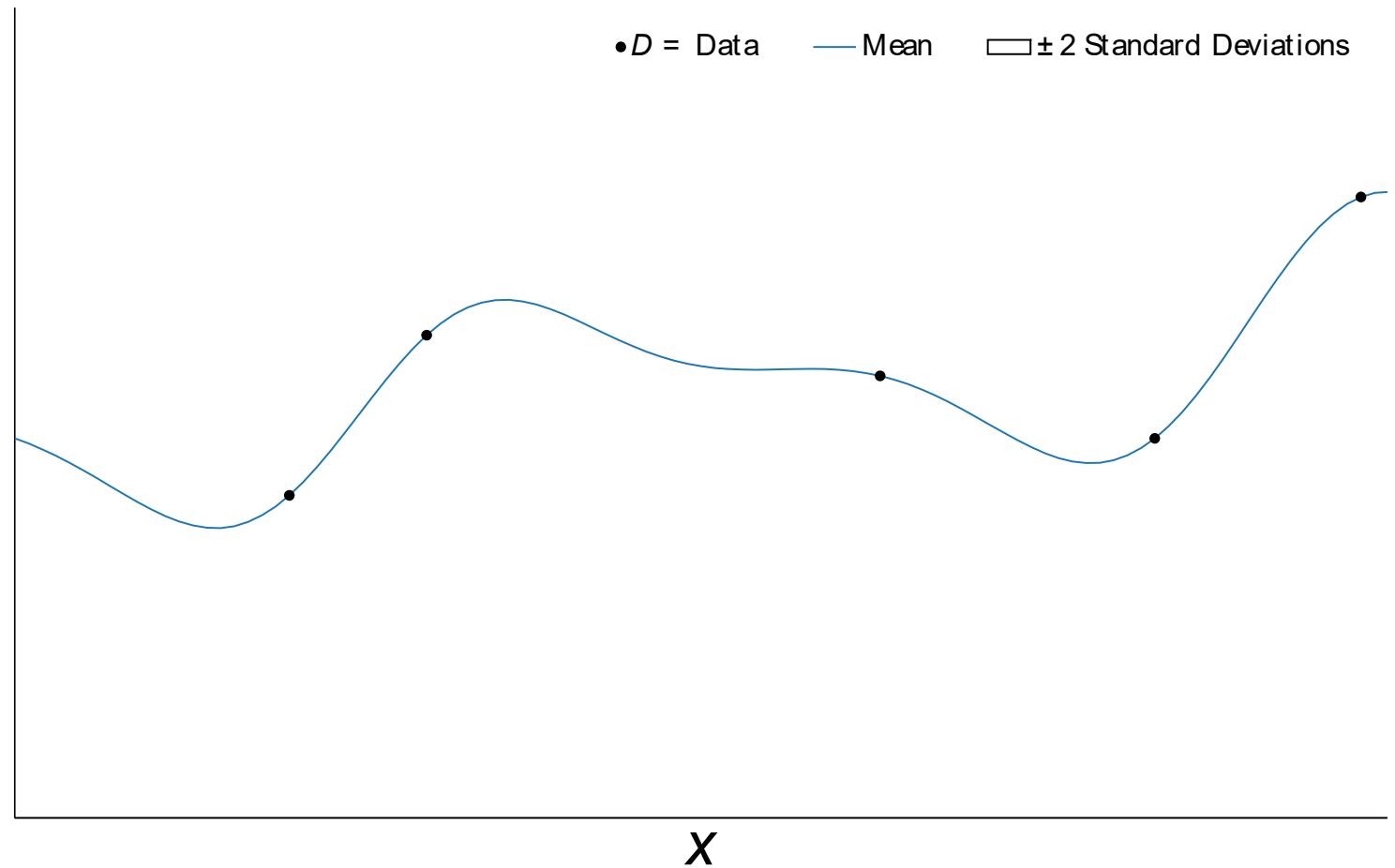
# GP Prior

$$f: \mathbb{R}^p \mapsto \mathbb{R} \sim \mathcal{GP}(f; \mu(x) = 0, \Sigma(x, x') = \exp(-(x - x')^2))$$



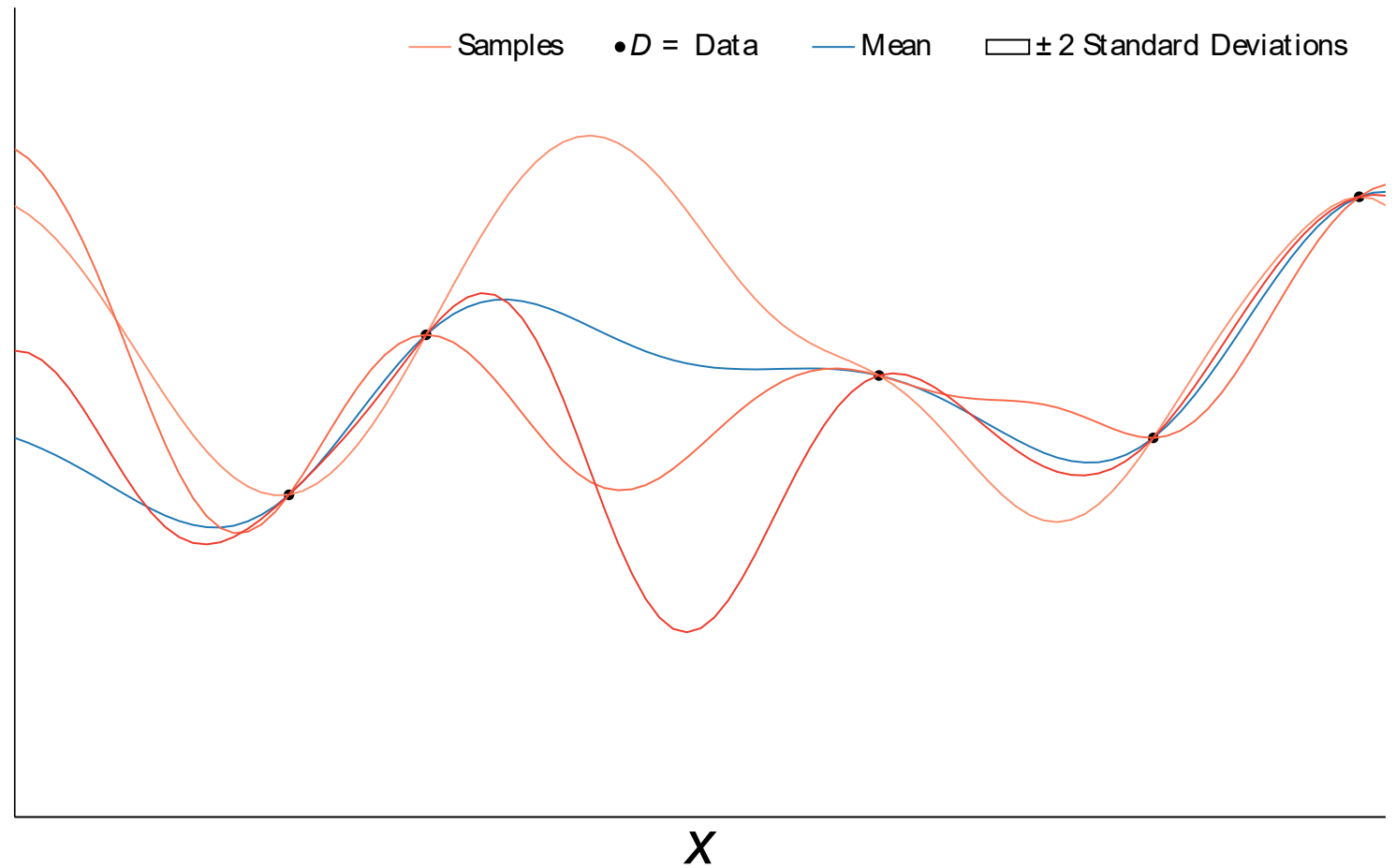
# GP Posterior

$$f | \mathcal{D} \sim \mathcal{GP}(f; \mu_{\mathcal{D}}, \Sigma_{\mathcal{D}})$$



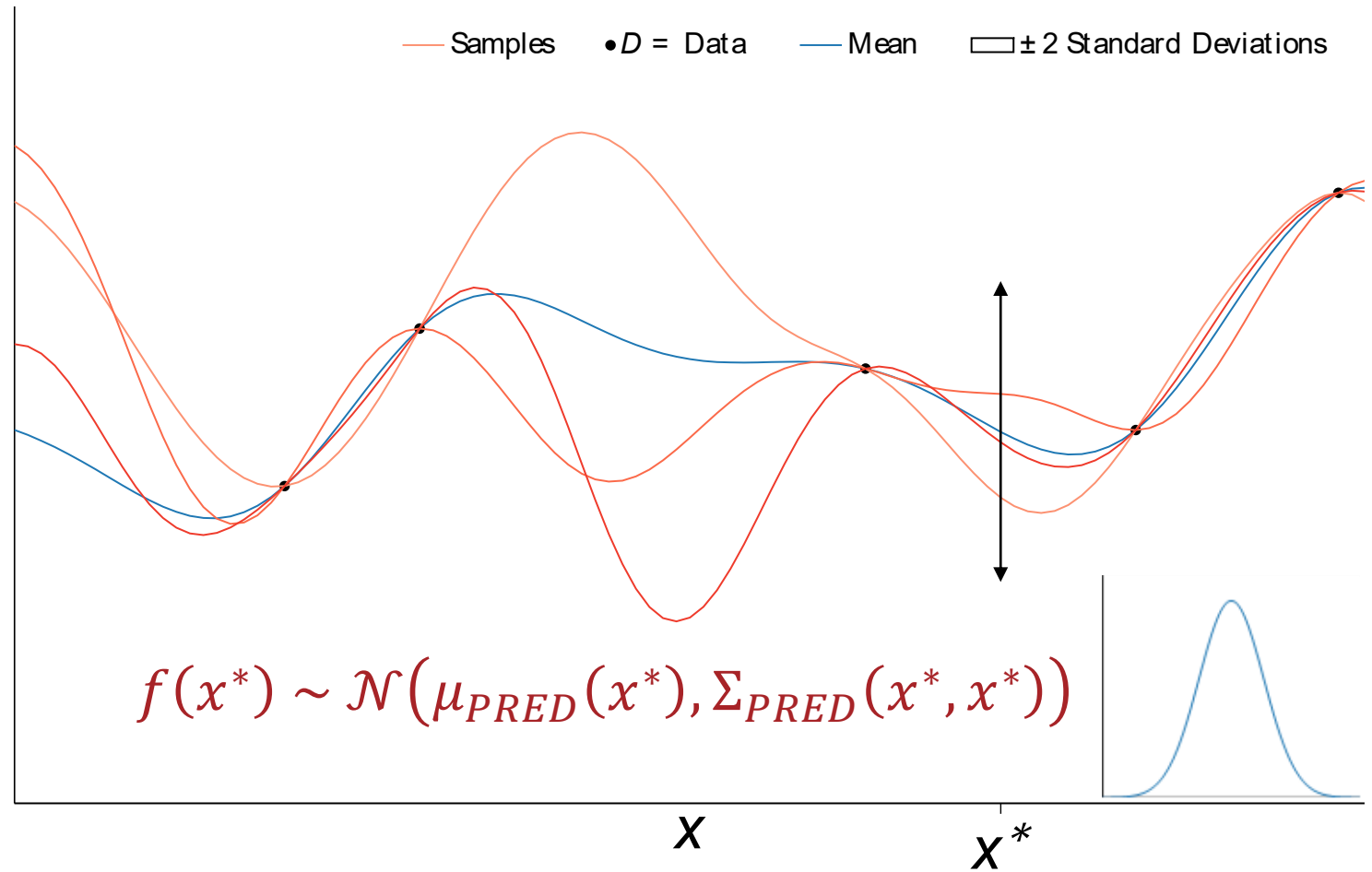
# GP Posterior

$$f | \mathcal{D} \sim \mathcal{GP}(f; \mu_{\mathcal{D}}, \Sigma_{\mathcal{D}})$$



# GP Posterior

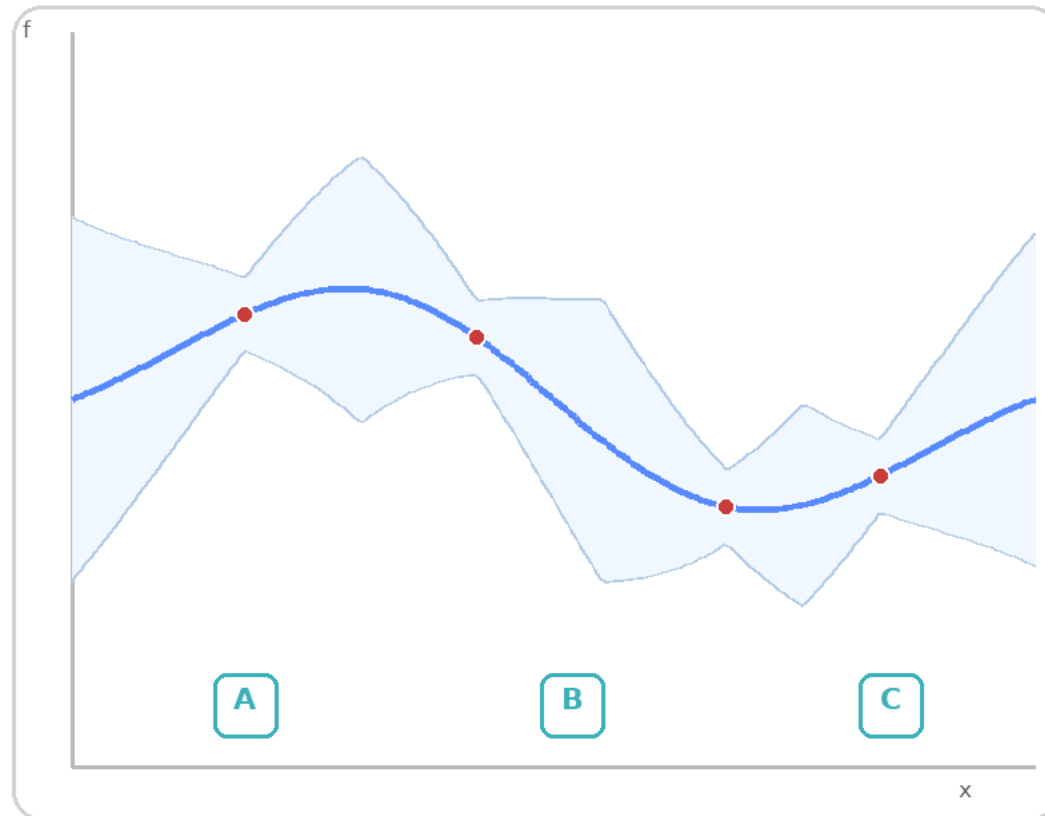
$$f | \mathcal{D} \sim \mathcal{GP}(f; \mu_{\mathcal{D}}, \Sigma_{\mathcal{D}})$$



# Quick Game

## Quick game: where is GP uncertainty smallest?

Look at the posterior picture intuition: where should the shaded confidence band be narrowest?



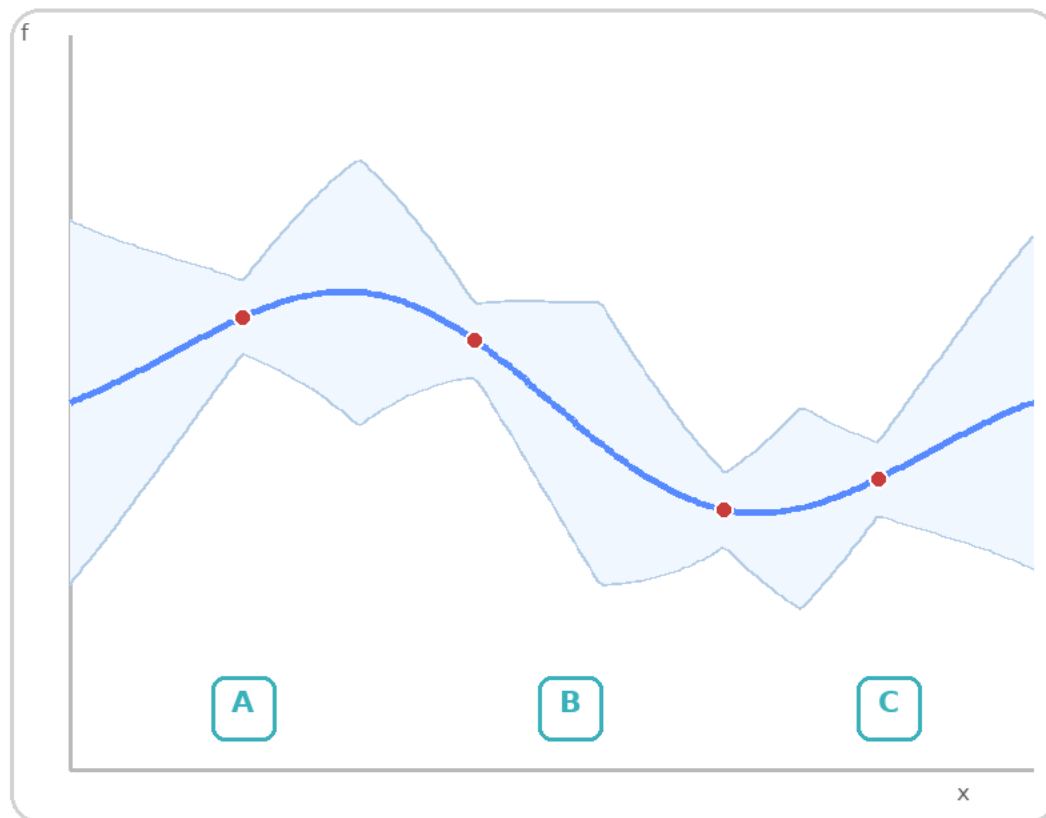
### Vote

- A. near an observed training point**
- B. midway between observed points**
- C. far away from observed points**

# Quick Game

## Quick game: where is GP uncertainty smallest?

Look at the posterior picture intuition: where should the shaded confidence band be narrowest?



### Vote

**A. near an observed training point**

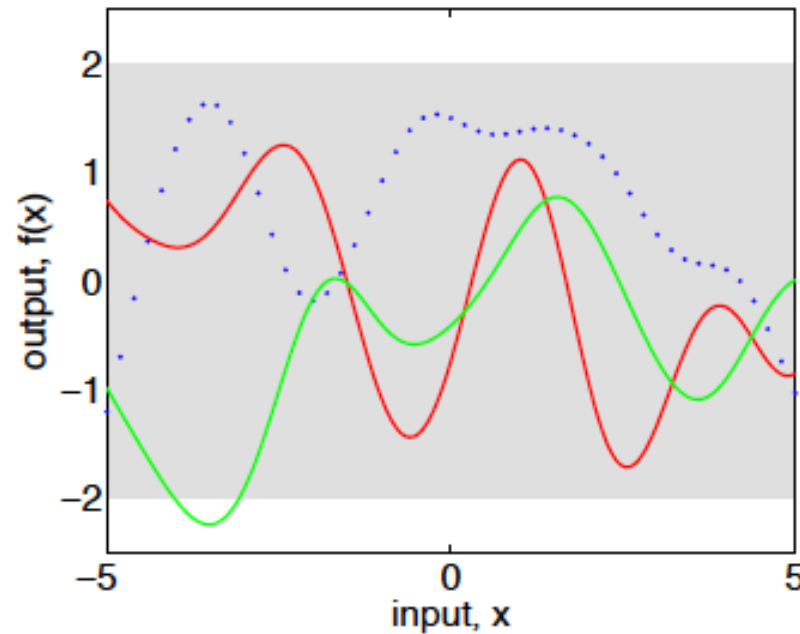
**B. midway between observed points**

**C. far away from observed points**

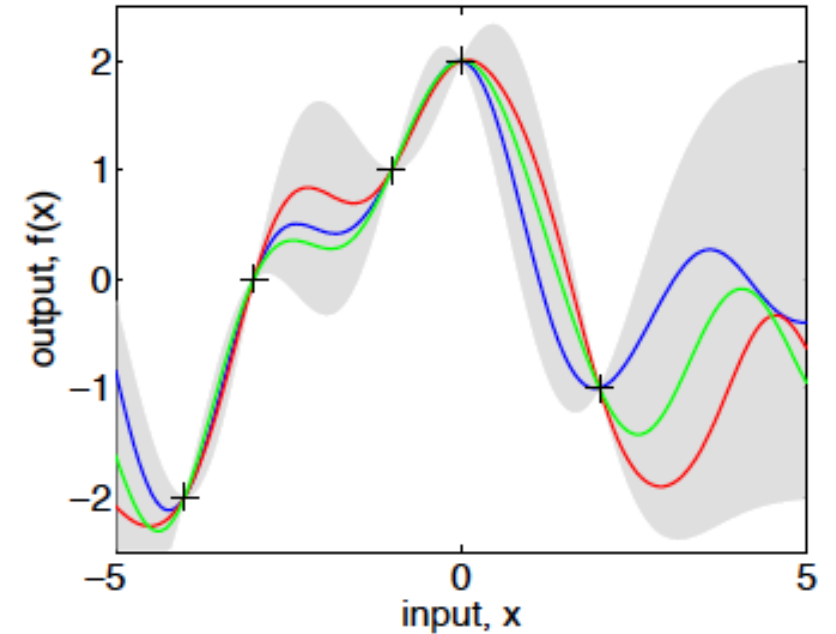
Reveal: A. In the GP posterior picture, uncertainty shrinks near observed points and grows away from them.

# Gaussian process regression

3 functions drawn at random from GP **prior**



3 functions drawn at random from **posterior** based on five noiseless observations



Shaded region represents pointwise mean  $\pm 2$  standard deviation

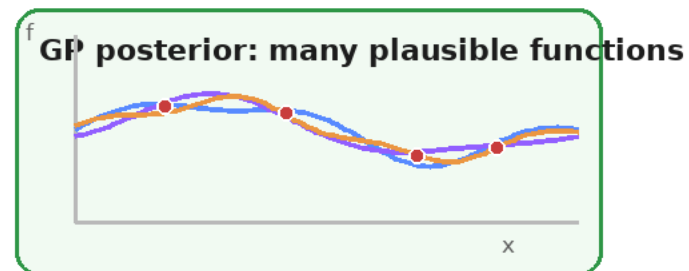
# Myth Busters

## MythBusters #2

“A Gaussian process is just one very smooth fitted curve.”

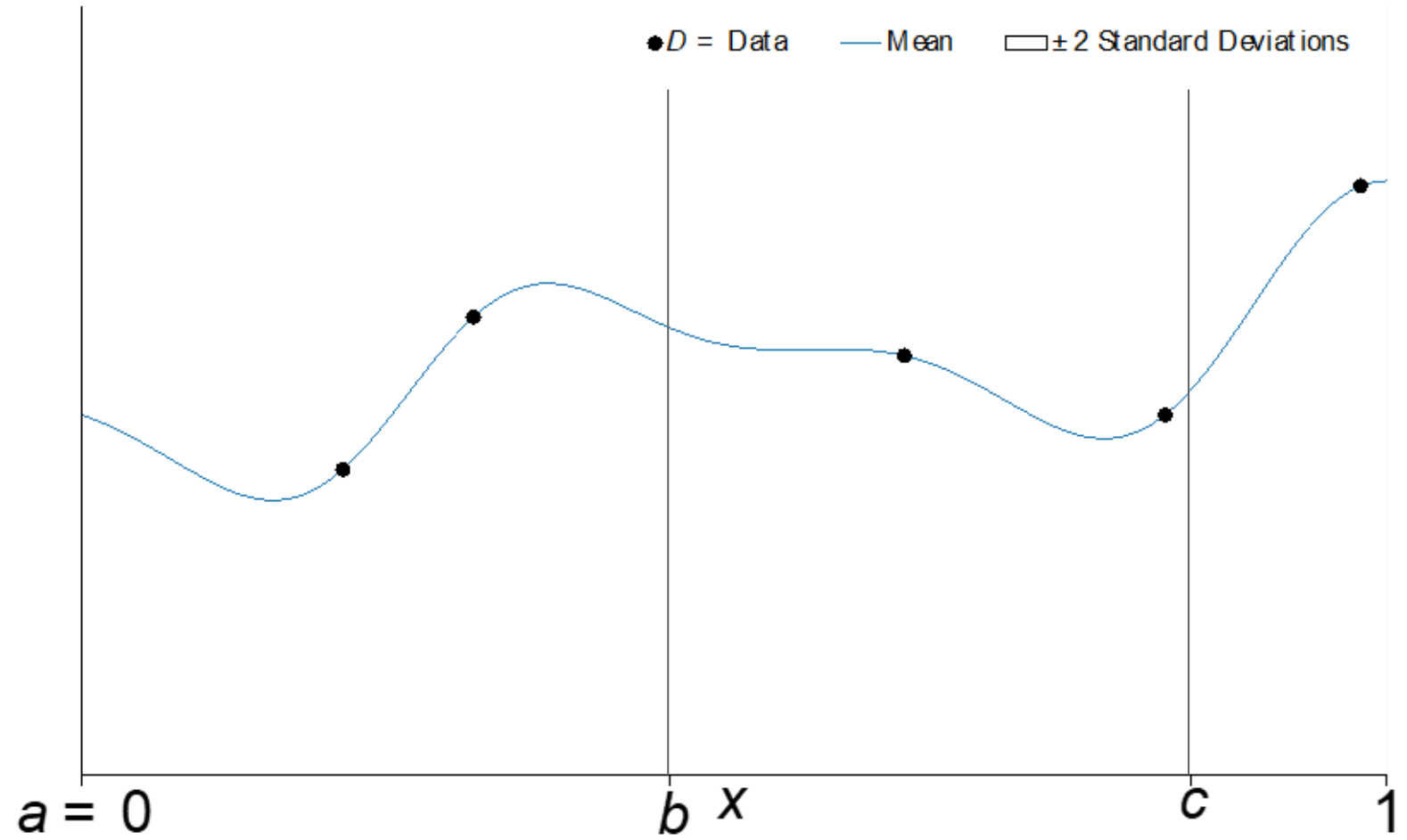
❌ **Myth: a GP is just a single nonlinear regressor with uncertainty added later.**

✅ **Reality: the lecture frames a GP as a distribution over functions. Data turns a prior over functions into a posterior over functions.**



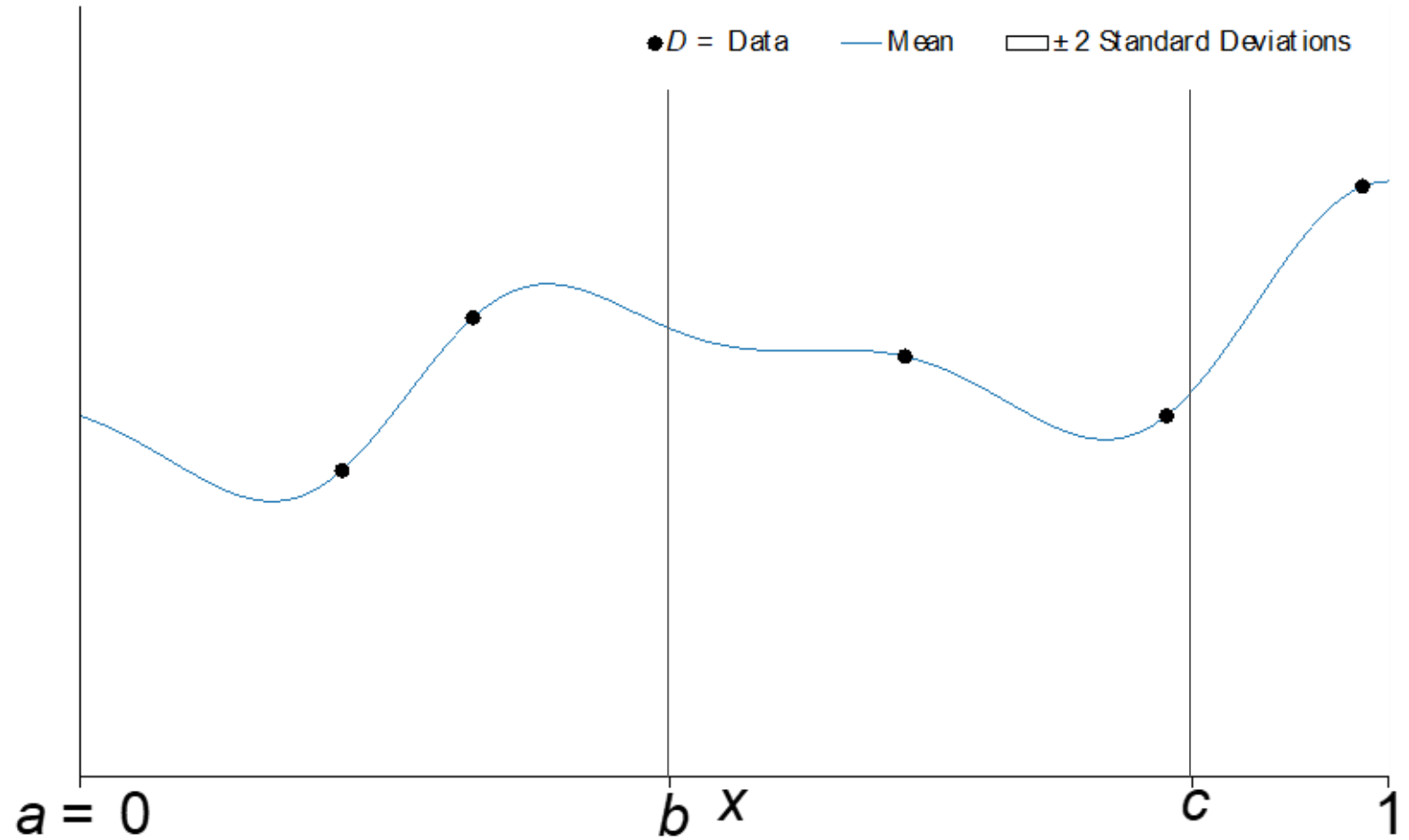
The mean is just one summary; the GP itself is the full distribution over functions.

# Active Learning



Suppose you can add one data point to your training data.

Which value of  $x$  would you add and why?



# Quick Game

## Two truths and a lie (kernels & GPs)

Vote for the lie: A, B, or C.

A. The Gaussian-RBF kernel corresponds to an infinite-dimensional feature map.

B. Gaussian process regression can be viewed as Bayesian linear regression written in kernel form.

C. In the GP posterior picture, predictive variance is largest exactly at the observed noiseless training points.

# Quick Game

## Two truths and a lie (kernels & GPs)

Vote for the lie: A, B, or C.

A. The Gaussian-RBF kernel corresponds to an infinite-dimensional feature map.

B. Gaussian process regression can be viewed as Bayesian linear regression written in kernel form.

C. In the GP posterior picture, predictive variance is largest exactly at the observed noiseless training points.

**Reveal: C is the lie.**

# Final Poll

Poll:



Lecture 25

## Final poll (pick one)

Which statement is **FALSE** according to the lecture?

- A. The kernel trick lets us compute inner products in transformed feature spaces implicitly.
- B. A Gaussian process can be interpreted as a distribution over functions.
- C. In a GP posterior, uncertainty is usually smallest near observed noiseless training points.
- D. A smaller Gaussian-RBF radius  $r$  makes the kernel more local.
- E. Gaussian process regression cannot be related to Bayesian linear regression.