

Lecture 14: Bayesian Linear Regression

Professor Ilias Bilonis

**Probabilistic interpretation of least
squares - Estimating the
measurement noise**

Reminder: Generalized linear model and least squares fit

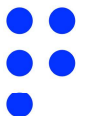
$$\mathbf{x}_{1:n} = (x_1, \dots, x_n) ; \quad \mathbf{y}_{1:n} = (y_1, \dots, y_n) \quad (w_1, \dots, w_m)$$

$$y = w_1 \underbrace{\varphi_1(x)}_{\text{feature, basis function}} + \dots + w_m \underbrace{\varphi_m(x)}_{\text{feature, basis function}} = \sum_{j=1}^m w_j \varphi_j(x) = \underbrace{\varphi(x)^T}_{(\varphi_1(x), \dots, \varphi_m(x))} \underline{w}$$

$$\min L(\underline{w}) = \sum_{i=1}^n (y_i - \varphi(x_i)^T \underline{w})^2 \quad (y_1, \dots, y_n)$$

$$\nabla_{\underline{w}} L(\underline{w}) = 0 \Rightarrow \underbrace{\Phi^T \Phi}_{\text{Design matrix } n \times m} \cdot \underline{w} = \underbrace{\Phi^T}_{\text{Design matrix } n \times m} \cdot \underline{y}$$

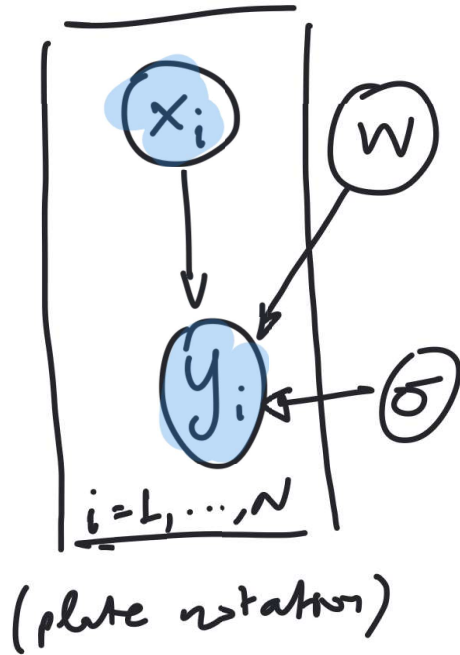
$$\Phi = \begin{pmatrix} \varphi_1(x_1) & \dots & \varphi_m(x_1) \\ \vdots & & \vdots \\ \varphi_1(x_n) & & \varphi_m(x_n) \end{pmatrix}$$



Open questions

- How do I quantify the measurement noise?
- How do we avoid overfitting?
- How do I quantify epistemic uncertainty induced by limited data?
- How do I choose any remaining parameters?
- How do I choose which basis functions to keep?

Probabilistic interpretation



Prior: $\underline{w} \sim p(\underline{w})$
 $\sigma \sim p(\sigma)$

Likelihood:
 $y_i | x_i, \underline{w}, \sigma^2 \sim \mathcal{N}(\underline{\phi}^T(x_i) \underline{w}, \sigma^2)$

$$P(y_{1:N} | x_{1:N}, \underline{w}, \sigma^2) = \prod_{i=1}^N p(y_i | x_i, \underline{w}, \sigma^2)$$

$$= \prod_{i=1}^N (2\pi)^{-1/2} \sigma^{-1} \exp \left\{ -\frac{(y_i - \underline{\phi}^T(x_i) \underline{w})^2}{2\sigma^2} \right\}$$

$$= (2\pi)^{-N/2} \sigma^{-N} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \underline{\phi}^T(x_i) \underline{w})^2 \right\}$$

Posterior \propto Likelihood \cdot Prior

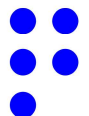
$$P(\underline{w}, \sigma | x_{1:N}, y_{1:N}) \propto \underbrace{P(y_{1:N} | x_{1:N}, \underline{w}, \sigma^2)}_{\text{Likelihood}} \underbrace{p(\underline{w})}_{\text{Prior}} \underbrace{p(\sigma)}_{\text{Prior}}$$

Maximum likelihood estimate of weights yields least squares

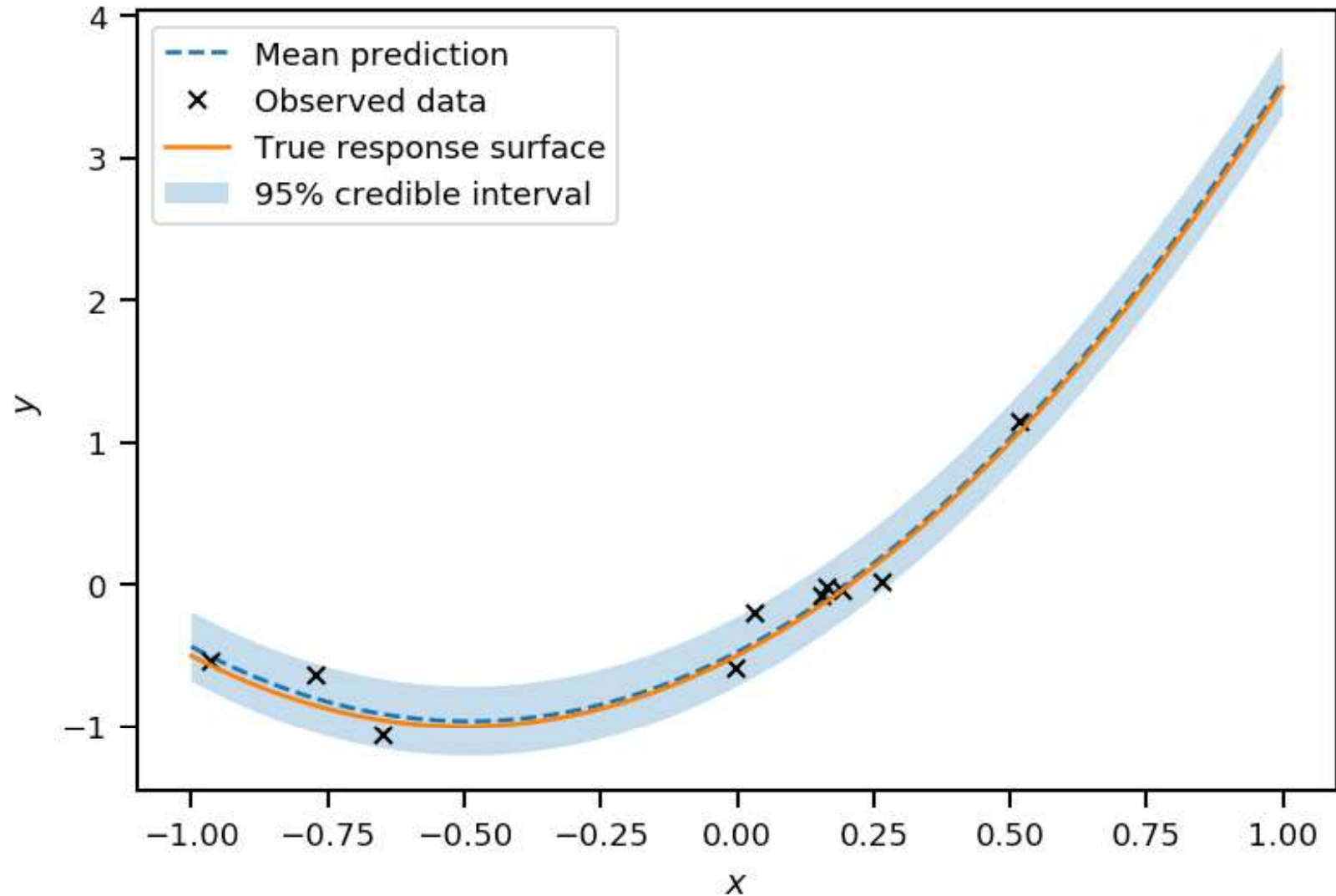
$$\log p(y_{1:n} | x_{1:n}, \underline{w}, \sigma) = -\frac{N}{2} \log 2\pi - N \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \underbrace{\Phi^T(x_i) \underline{w}}_{L(\underline{w})})^2$$

$$\max_{\underline{w}} \log \text{like} \equiv \min_{\underline{w}} L(\underline{w})$$

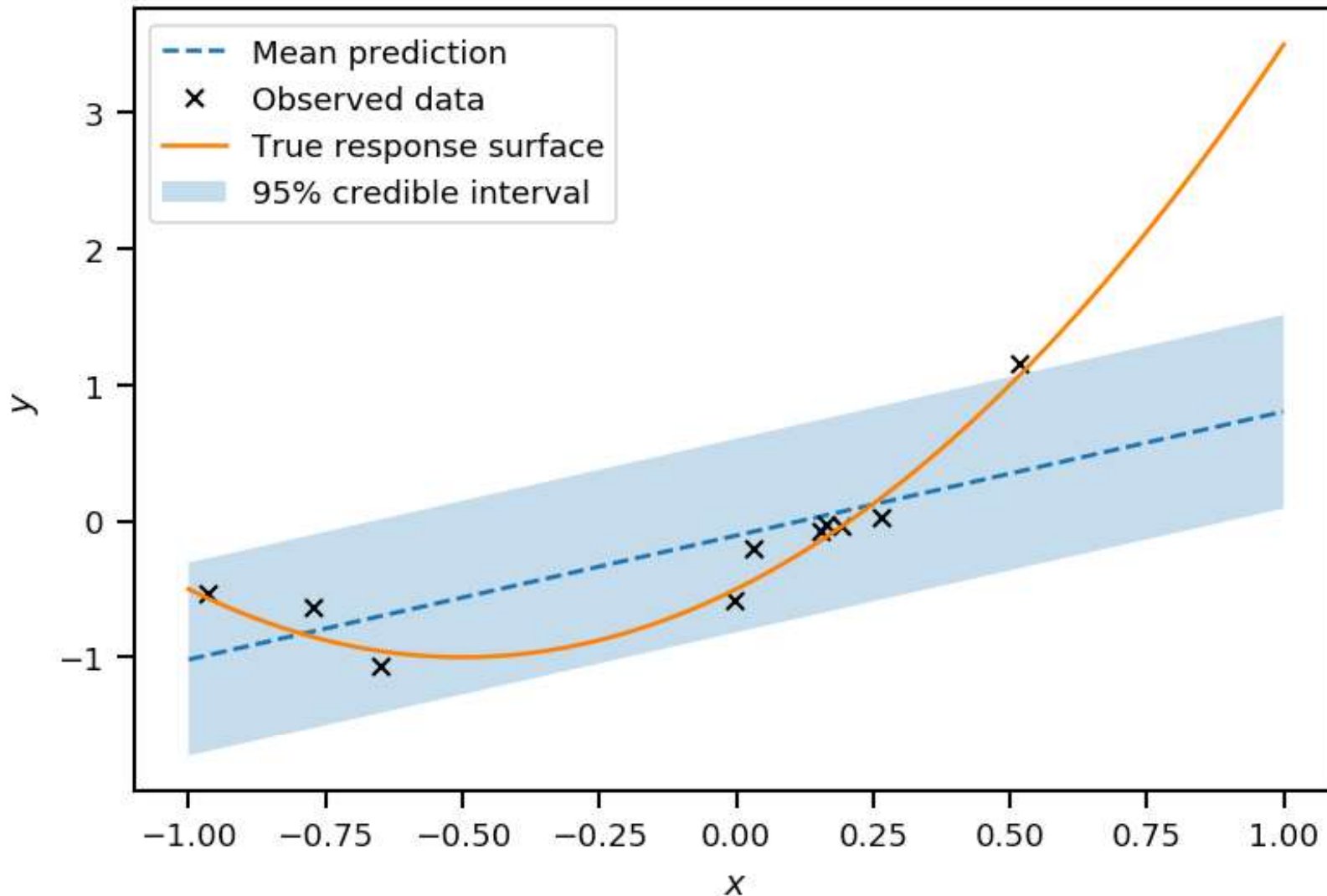
$$\Phi^T \Phi \underline{w} = \Phi^T \underline{y}$$



Example: Degree 2 polynomial



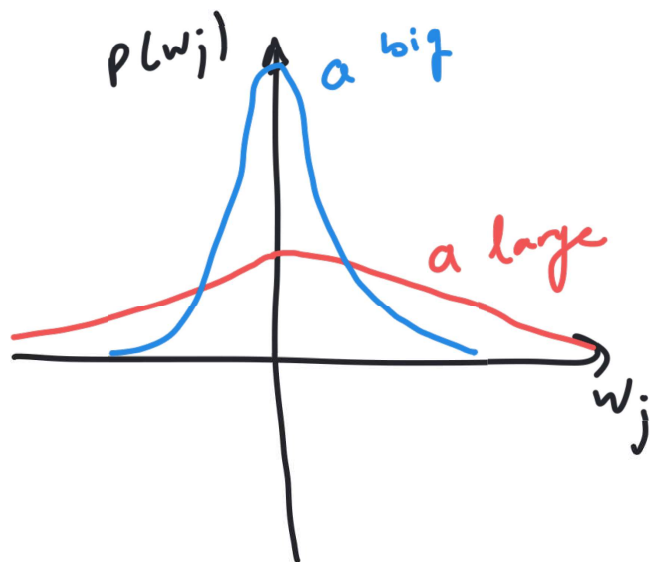
Example: Degree 1 polynomial



Open questions

- How do I quantify the measurement noise?
- How do we avoid overfitting?
- How do I quantify epistemic uncertainty induced by limited data?
- How do I choose any remaining parameters?
- How do I choose which basis functions to keep?

Gaussian prior on weights



$$w_j \sim \mathcal{N}(0, \alpha^{-1})$$

$\alpha = \text{precision} = \frac{1}{\text{variance}}$

$$p(\underline{w}) = \prod_{j=1}^m p(w_j)$$

$$\propto \exp \left\{ -\frac{\alpha}{2} \sum_{j=1}^m w_j^2 \right\}$$

Maximum a posteriori estimate

posterior \propto likelihood \times prior

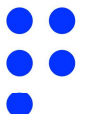
$$p(\underline{w} | x_{1:n}, y_{1:n}, \sigma) \propto p(y_{1:n} | x_{1:n}, \sigma^2) p(\underline{w})$$

$$\begin{aligned} \max_{\underline{w}} \log \text{post} &= \log \text{like} + \log p(\underline{w}) \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \underline{\varphi}^T(x_i) \underline{w})^2 - \frac{\alpha}{2} \sum_{j=1}^m w_j^2 + \text{const.} \end{aligned}$$

$$\min_{\underline{w}} \underbrace{\sum_{i=1}^n (y_i - \underline{\varphi}^T(x_i) \underline{w})^2}_{\text{SF}} + \underbrace{\frac{\alpha}{2} \sum_{j=1}^m w_j^2}_{\text{Regularization term}}$$

SF

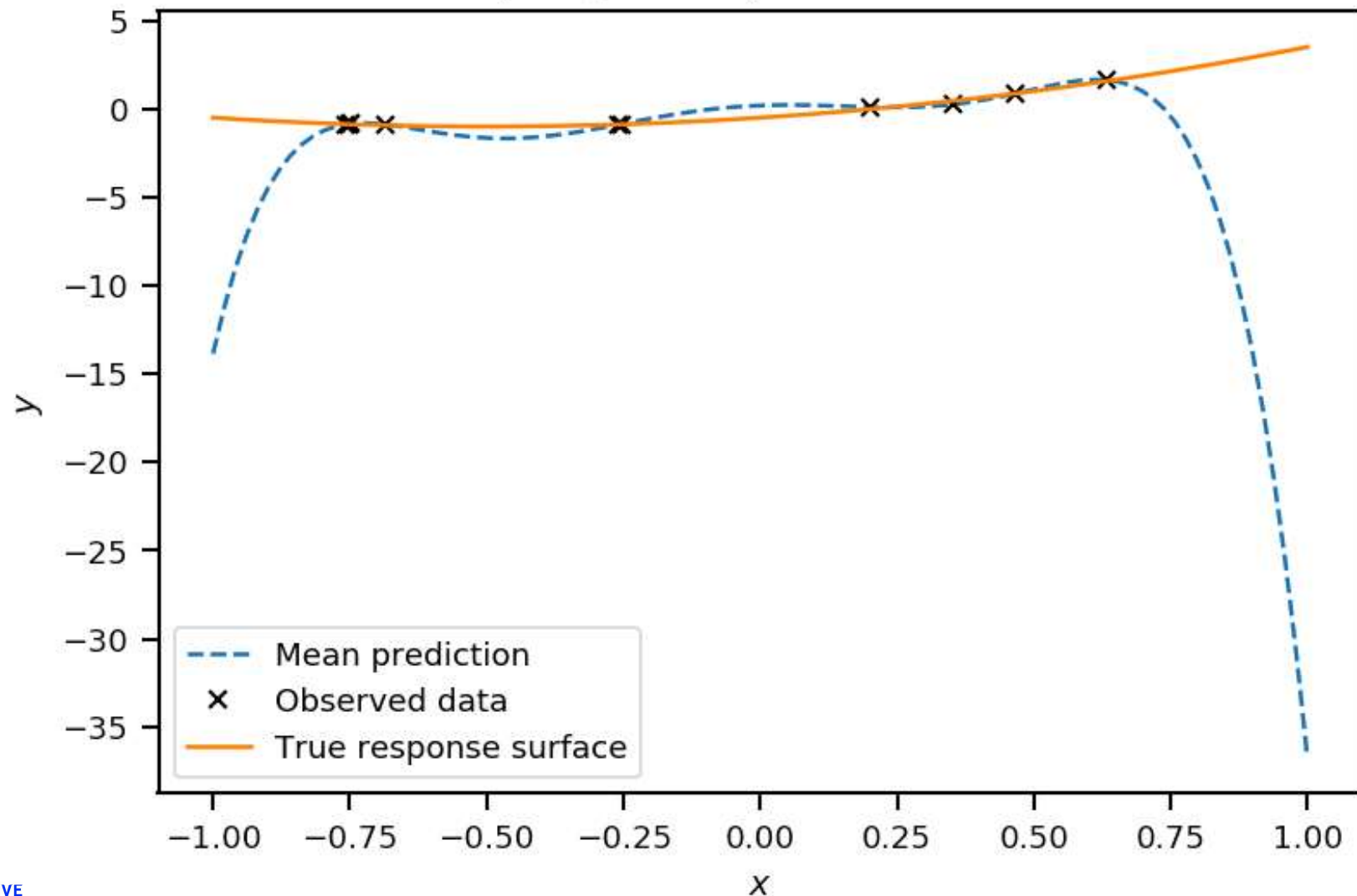
Regularization term



Example: Degree 6 polynomial

$$(\alpha = 0)$$

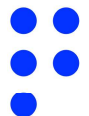
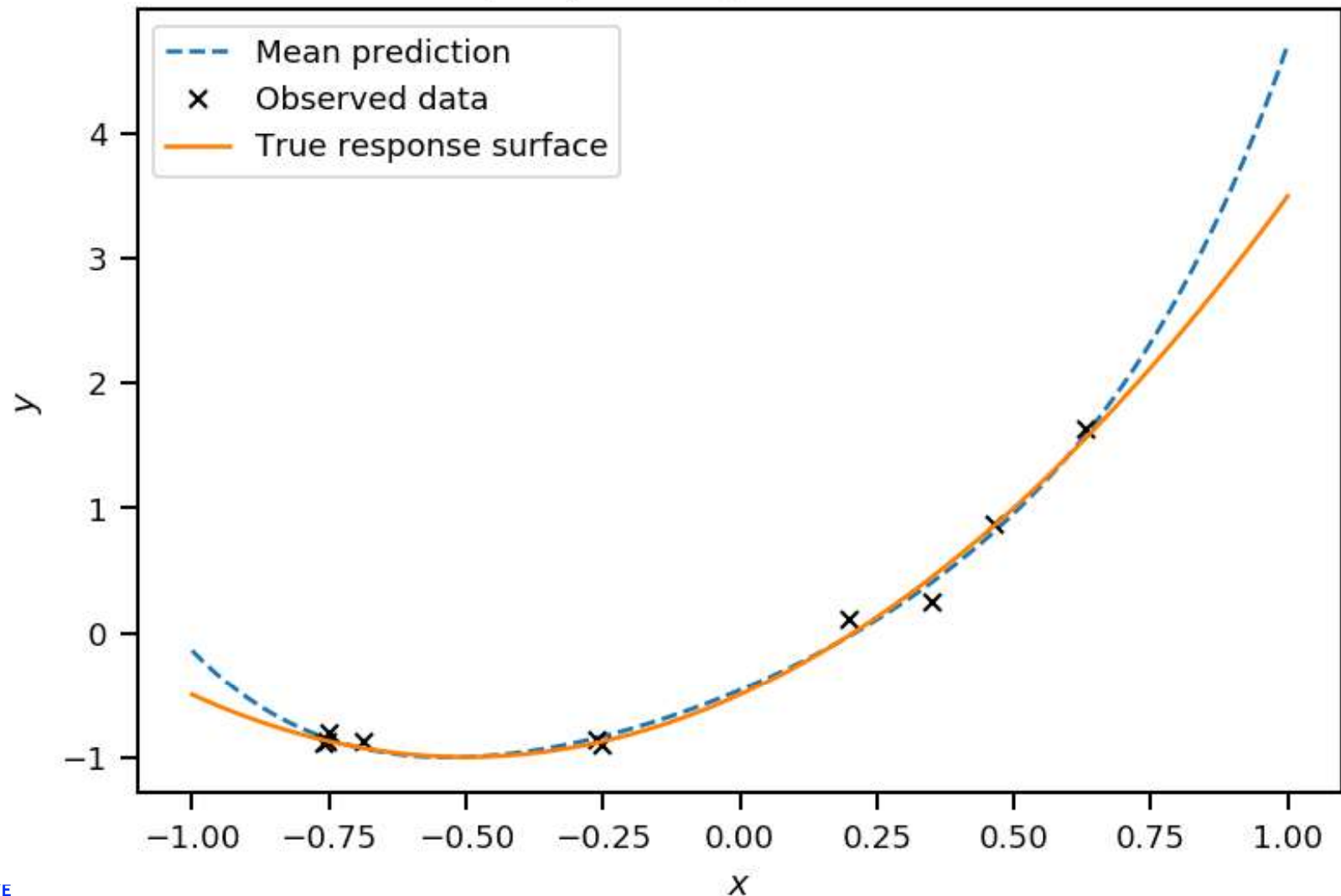
$\rho = 6, \sigma = 0.10, \alpha = 0.00e + 00$



Example: Degree 6 polynomial

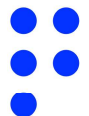
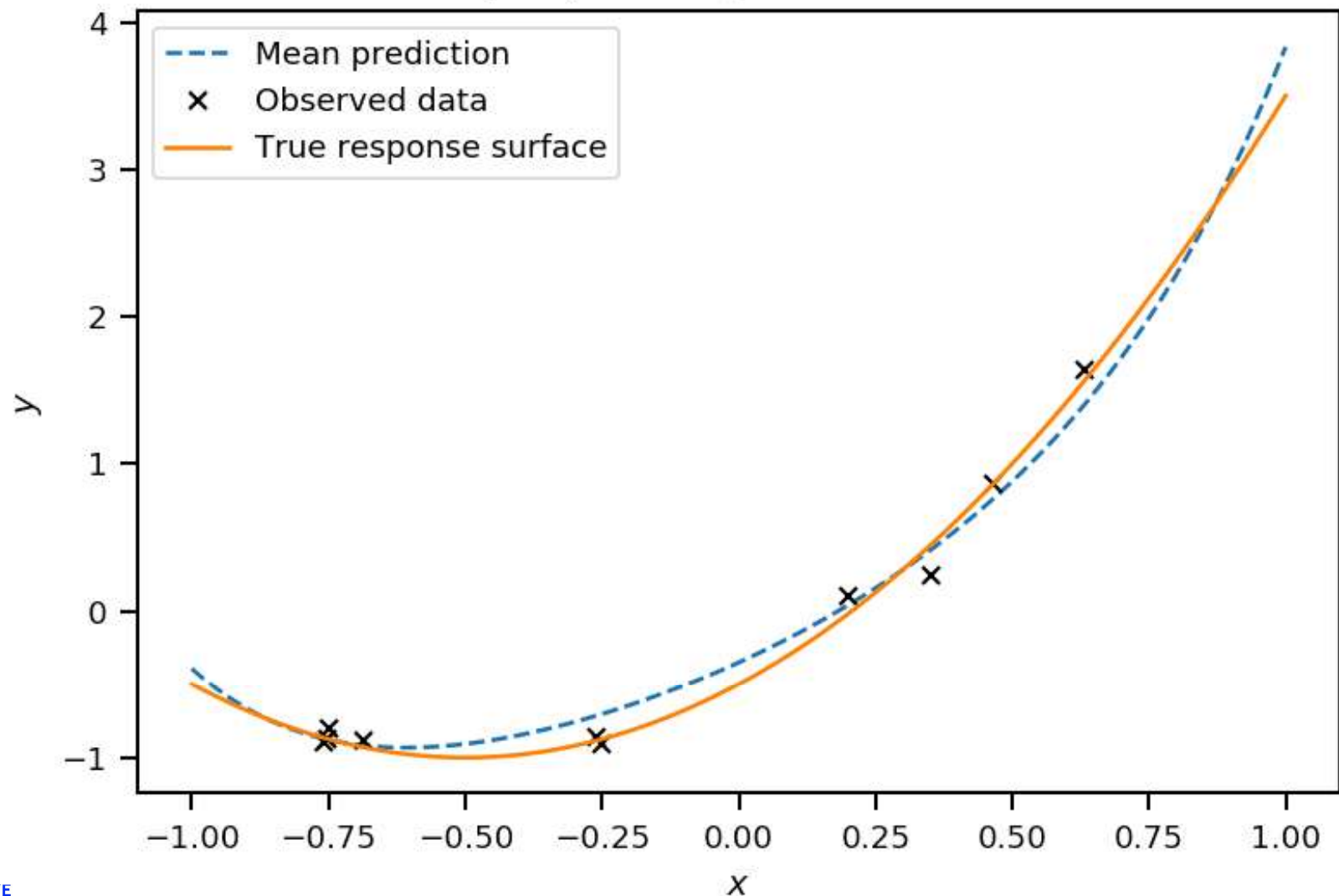
$(\alpha = 1)$

$\rho = 6, \sigma = 0.10, \alpha = 1.00e + 00$



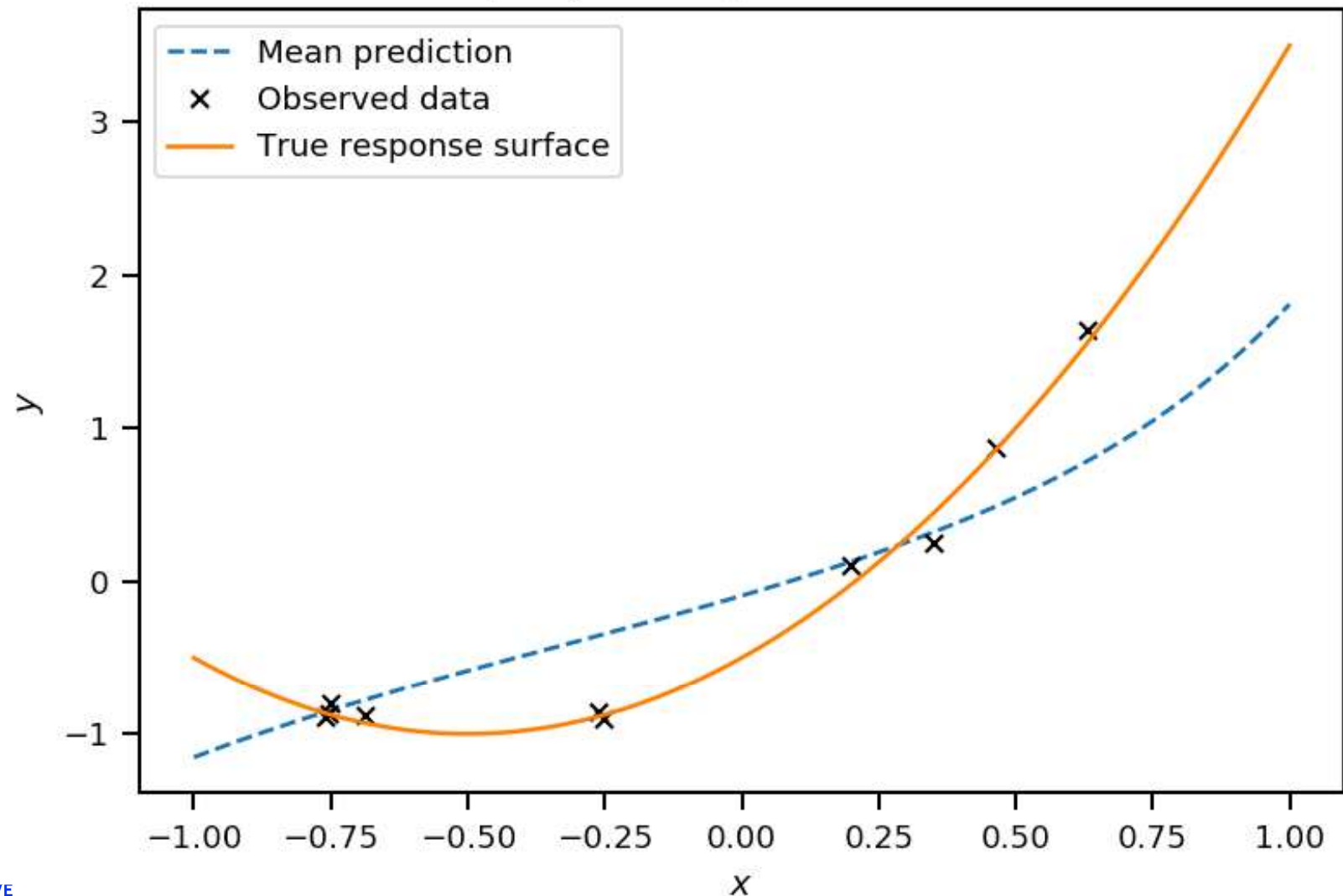
Example: Degree 6 polynomial ($\alpha = 10$)

$\rho = 6, \sigma = 0.10, \alpha = 1.00e + 01$

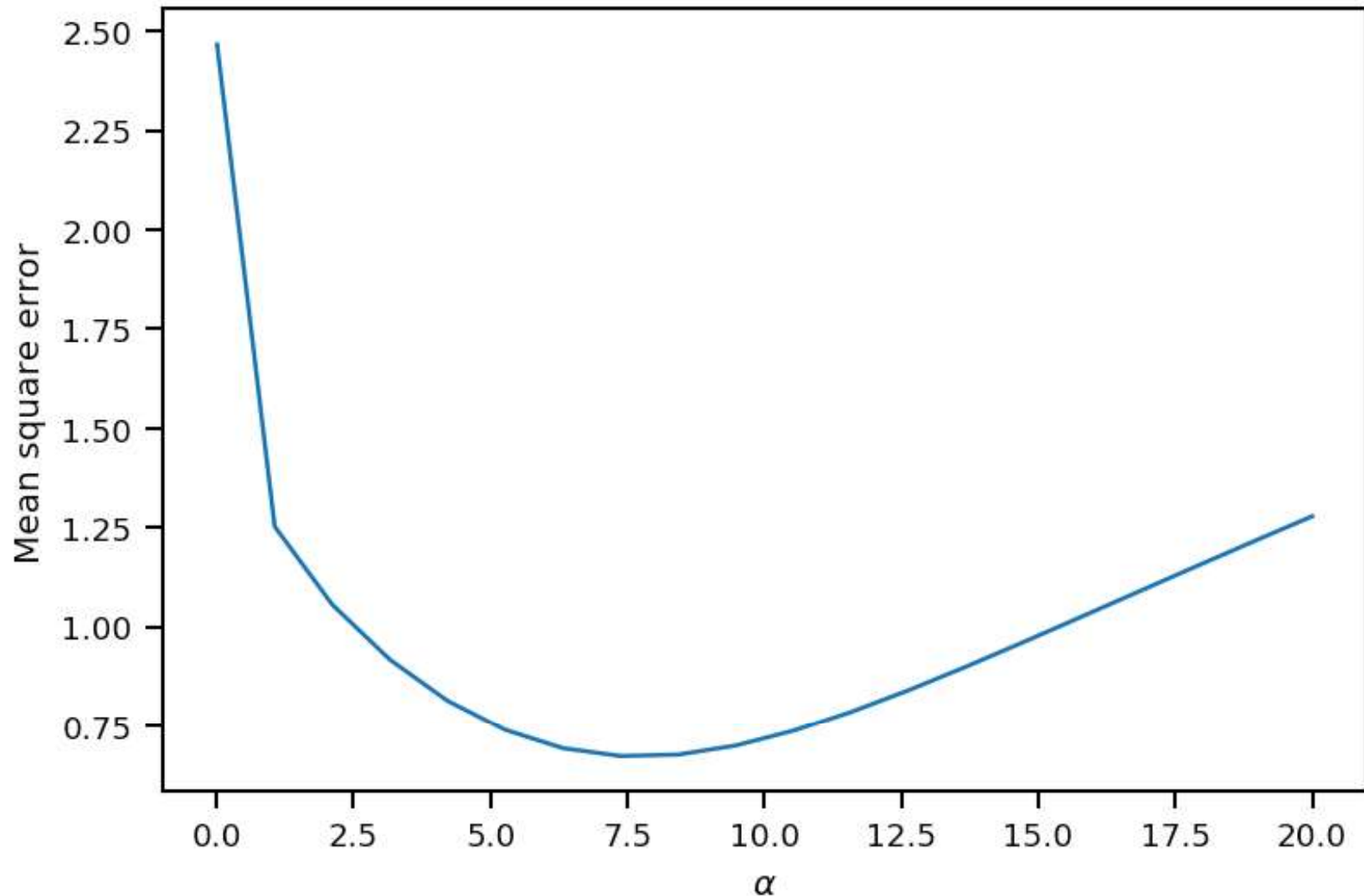


Example: Degree 6 polynomial ($\alpha = 100$)

$\rho = 6, \sigma = 0.10, \alpha = 1.00e + 02$



Mean square error over a validation dataset as a function of α

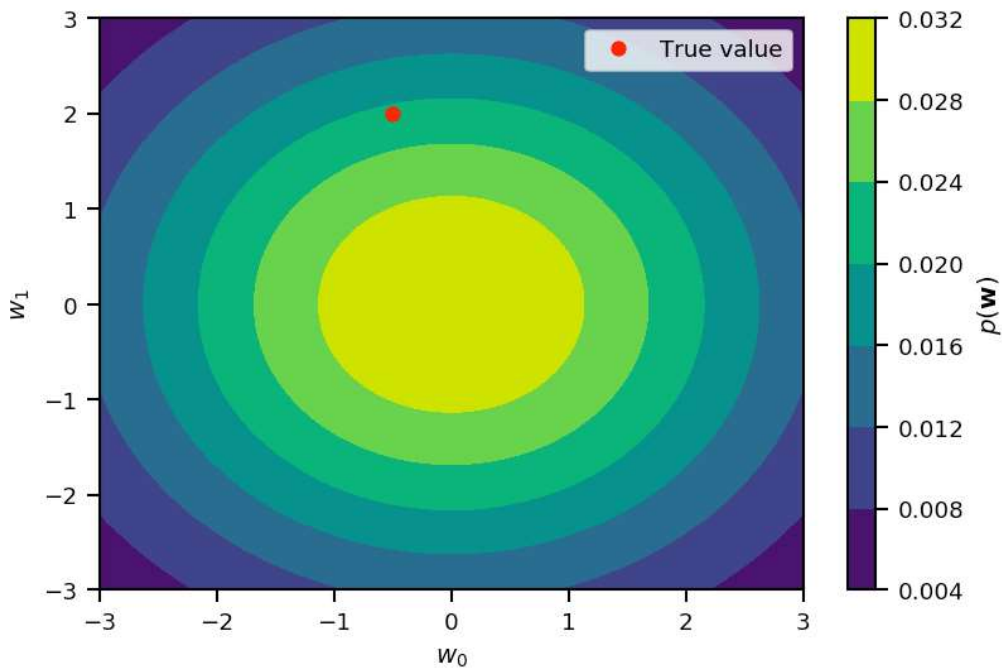


Open questions

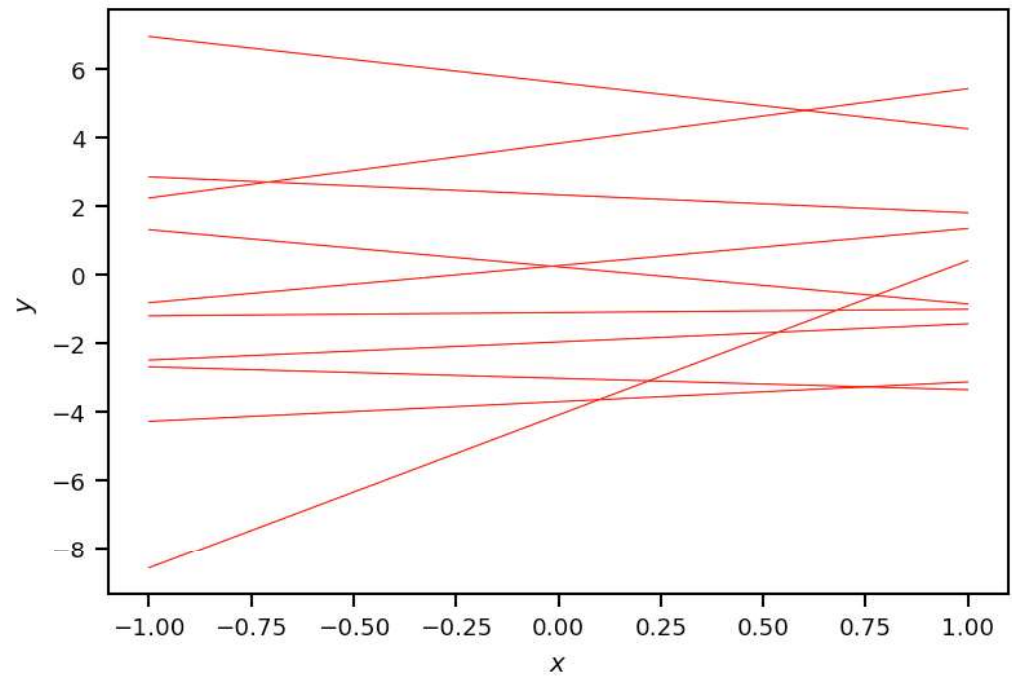
- How do I quantify the measurement noise?
- How do we avoid overfitting?
- How do I quantify epistemic uncertainty induced by limited data?
- How do I choose any remaining parameters?
- How do I choose which basis functions to keep?

Weight prior (linear regression)

Weight space



Model space



Weight posterior

posterior \propto likelihood \times prior

$$p(\underline{w} | x_{1:n}, y_{1:n}, \sigma) \propto p(y_{1:n} | x_{1:n}, \sigma, \underline{w}) \cdot p(\underline{w})$$

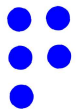
$$= \mathcal{N}(y_{1:n} | \Phi \underline{w}, \sigma^2 \mathbb{I}_N) \times \mathcal{N}(\underline{w} | 0, \alpha^{-1} \mathbb{I}_m)$$

$$\begin{pmatrix} \varphi^T(x_1) \underline{w} \\ \vdots \\ \varphi^T(x_N) \underline{w} \end{pmatrix}$$

$$= \mathcal{N}(\underline{w} | \underline{u}, \underline{\Sigma})$$

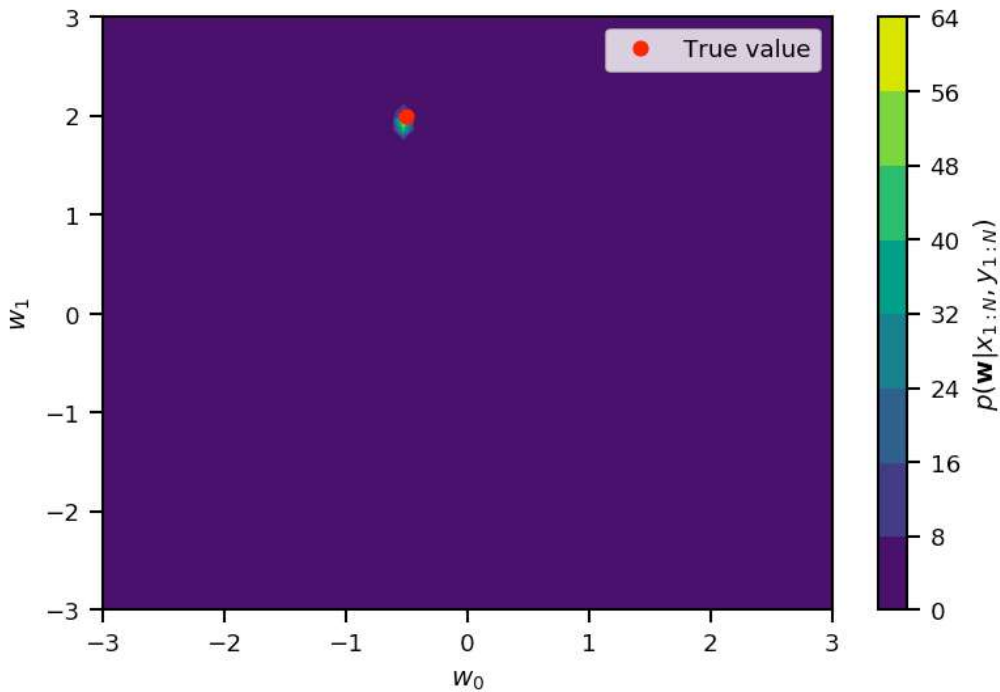
$$\underline{\Sigma} = (\sigma^2 \Phi^T \Phi + \alpha \mathbb{I}_m)^{-1}$$

$$\underline{u} = \sigma^{-2} \sum_{i=1}^N \varphi^T(x_i) y_i$$

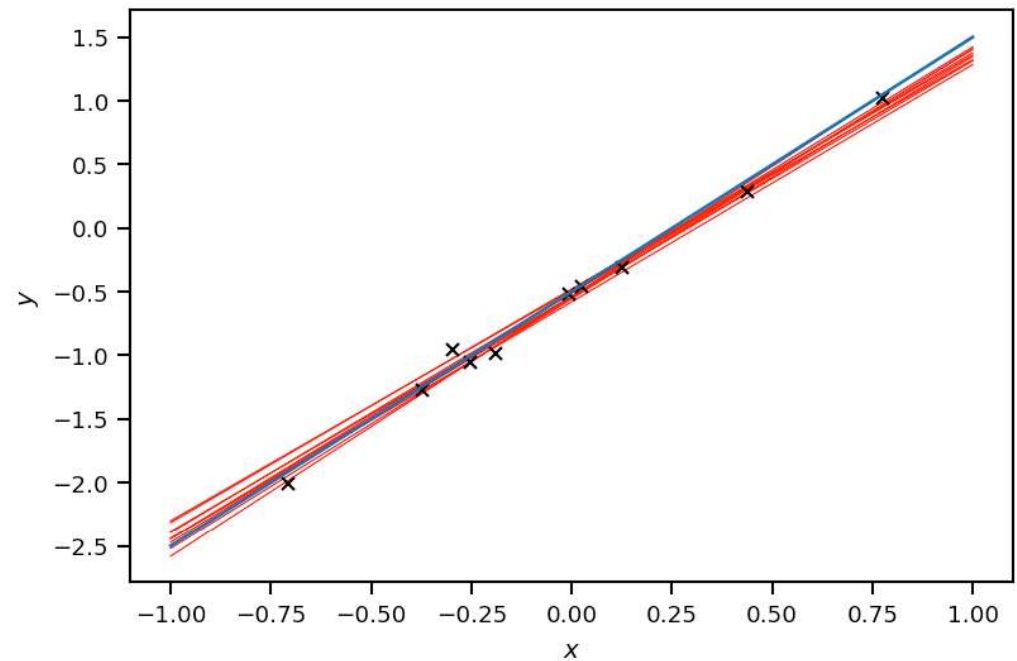


Weight posterior (linear regression)

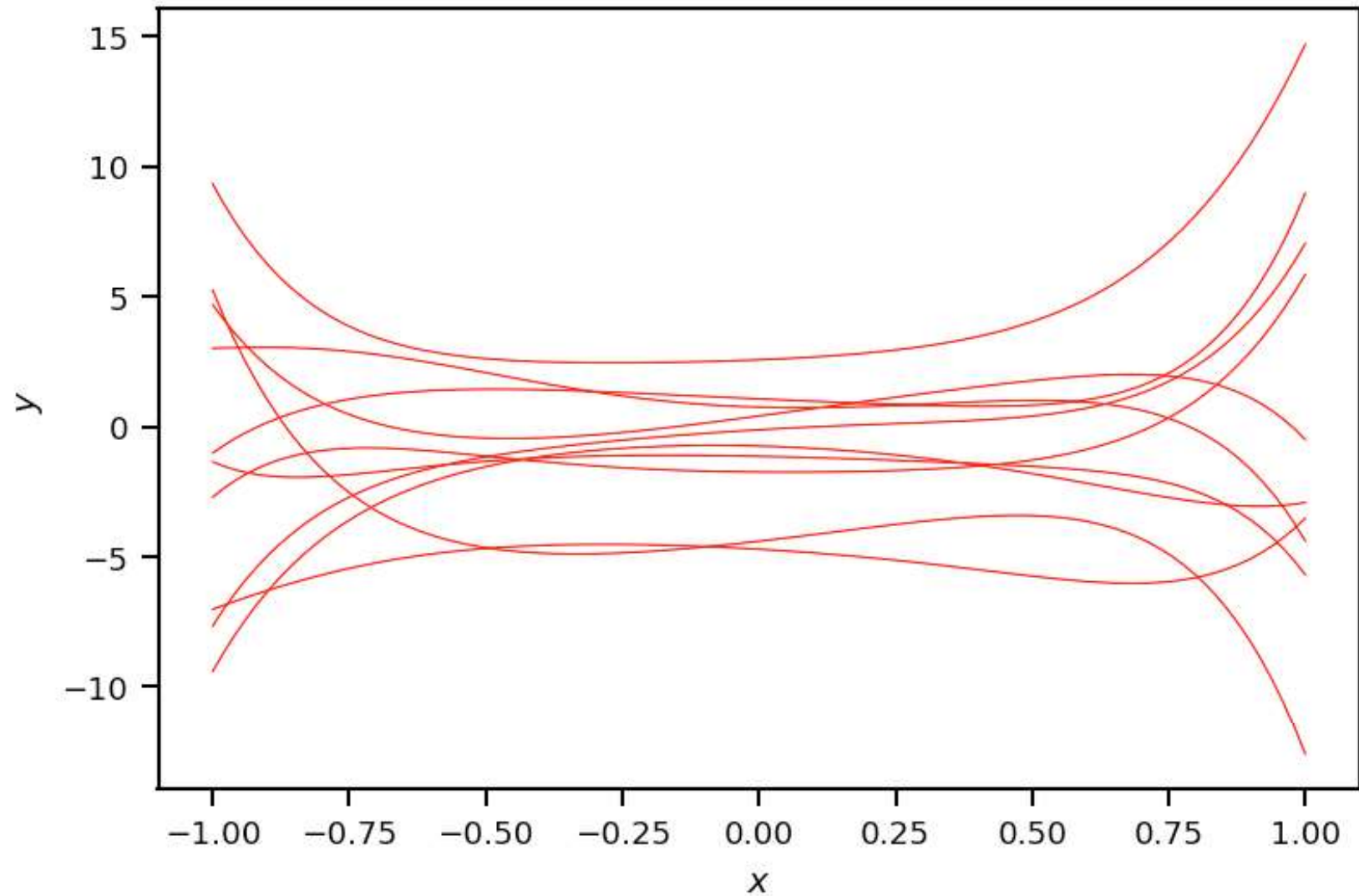
Weight space



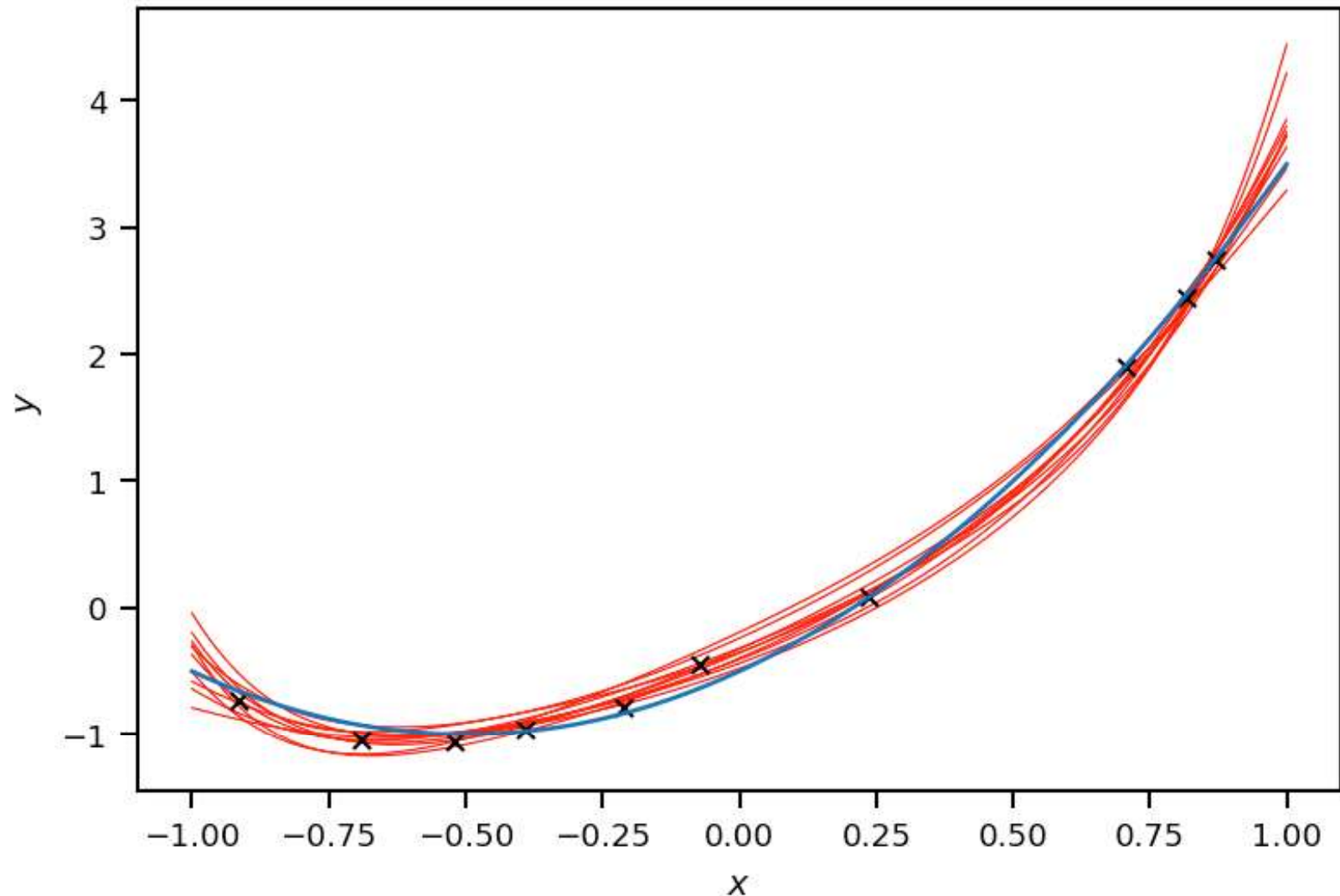
Model space



Example: 7th degree polynomial (prior)



Example: 7th degree polynomial (posterior)



Point-predictive distribution

$$p(y | x, x_{1:n}, y_{1:n}, a, \sigma^2) = ?$$

$$p(\underline{w} | x_{1:n}, y_{1:n}, a, \sigma^2) = \checkmark \quad ; \quad p(y | x, \underline{w}, \sigma^2) = \checkmark$$

Sum Rule : $p(A | I) = \sum_i p(A | B_i, I) p(B_i | I)$
 $p(B_1 \text{ or } B_2 \text{ or } \dots) = 1, \quad p(B_i, B_j) = 0$

$$p(y | \underbrace{x}_{\text{I}}, \underbrace{x_{1:n}, y_{1:n}}_{\text{I}}, a, \sigma^2) = \sum_i p(y | B_i, I) p(B_i | I)$$

$$= \int p(y | \underline{w}, I) p(\underline{w} | I) d\underline{w}$$

$$p(y | \underline{w}, x, x_{1:n}, y_{1:n}, a, \sigma^2) = p(y | x, \underline{w}, \sigma^2)$$

$$= \int p(y | x, \underline{w}, \sigma^2) p(\underline{w} | x_{1:n}, y_{1:n}, a, \sigma^2) d\underline{w}$$

$$= N(y | \underbrace{\varphi^T(x) \underline{\mu}}_{\text{post. mean of weight}}, \underbrace{\varphi^T(x) \left(\sum \varphi(x) + \sigma^2 \right)}_{\text{post. cov. of weight}})$$

post.
mean
of weight.

post.
cov.
of weight

noise
variance

epistemic



Example: Separating epistemic and aleatory uncertainties

