
AIL721: Deep Learning

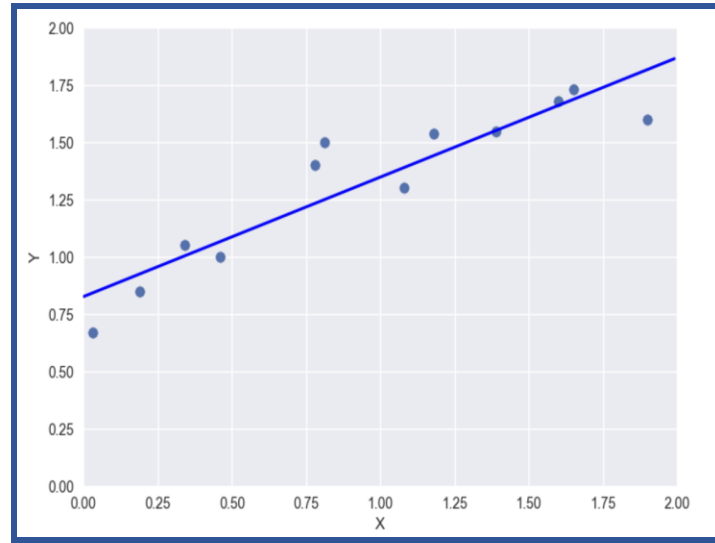
「Instructor: James Arambam」



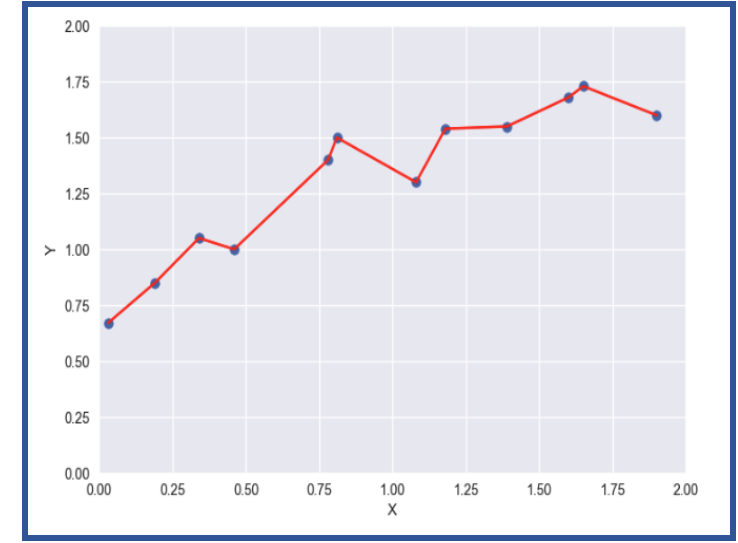
ScAI

Yardi School of Artificial Intelligence
Indian Institute of Technology Delhi

Model Performance



Model 1



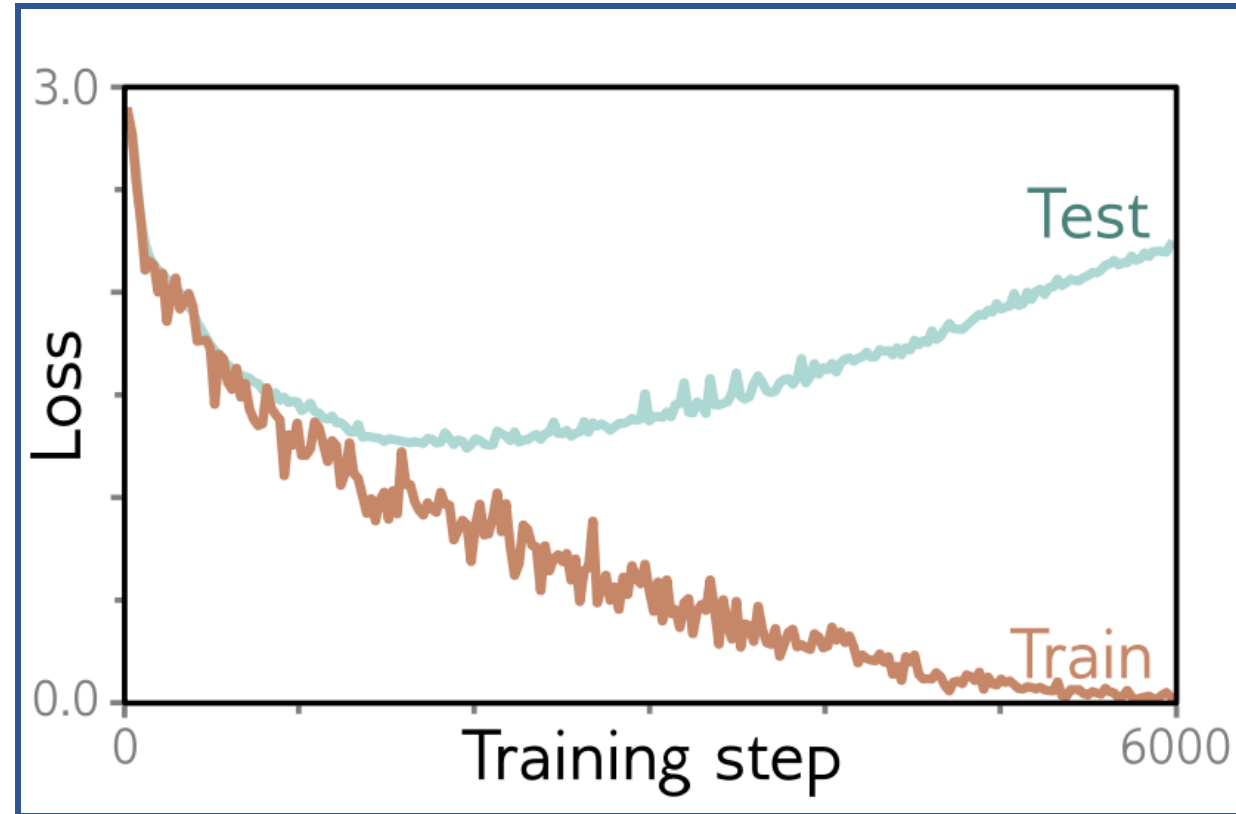
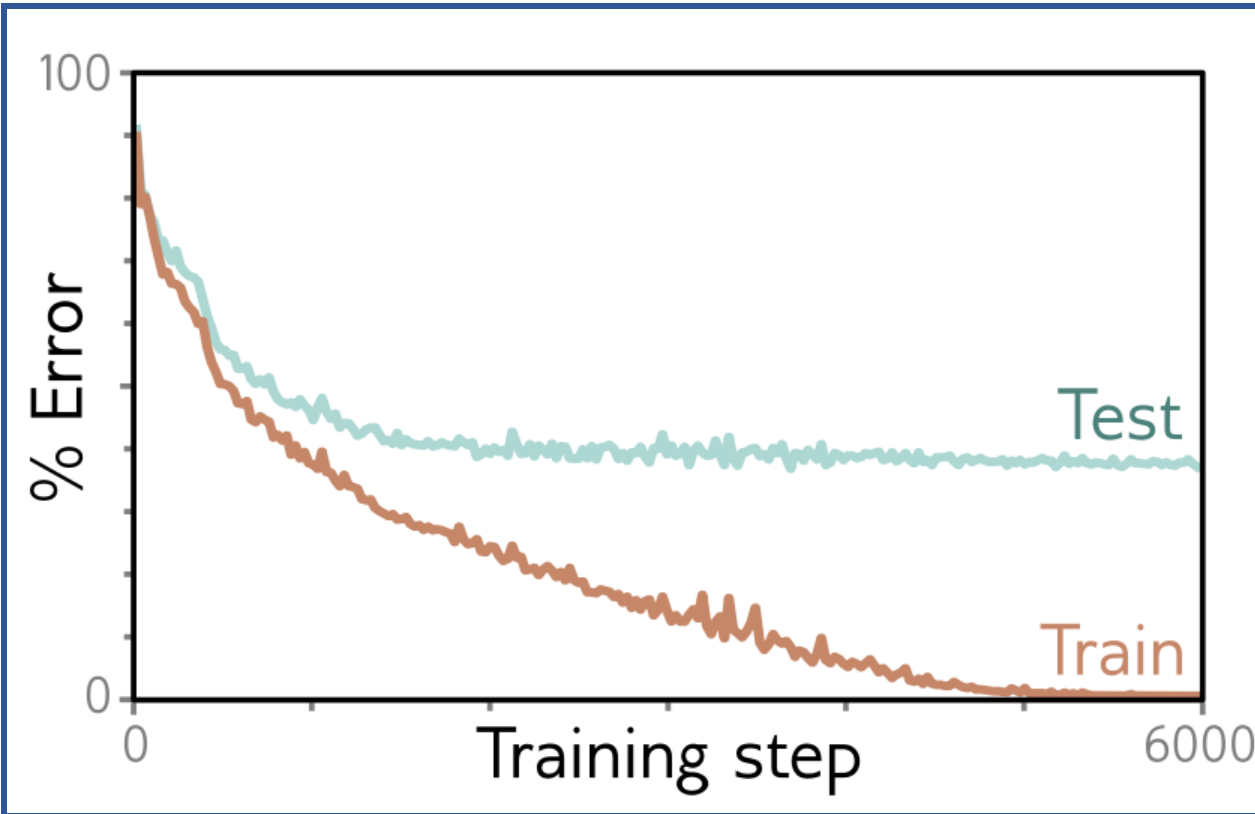
Model 2



Test Loss?

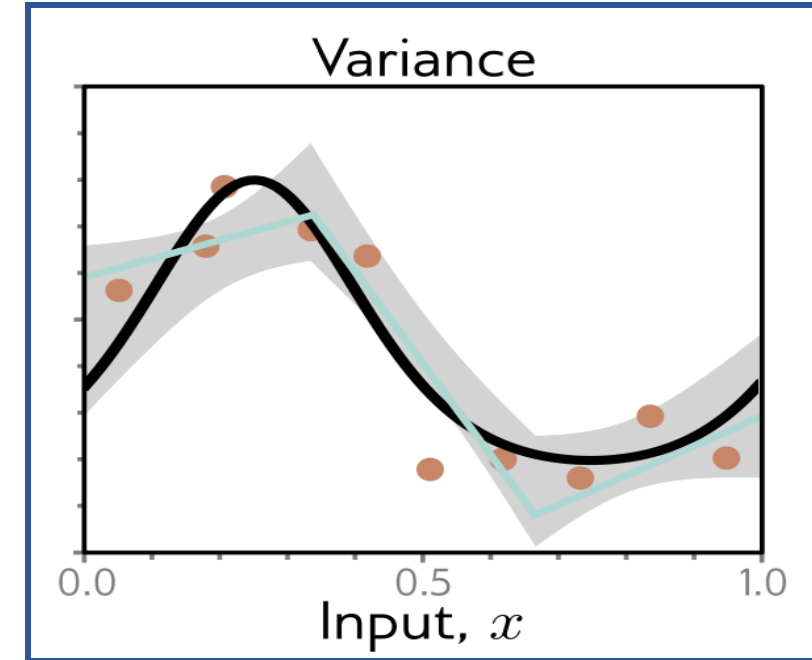
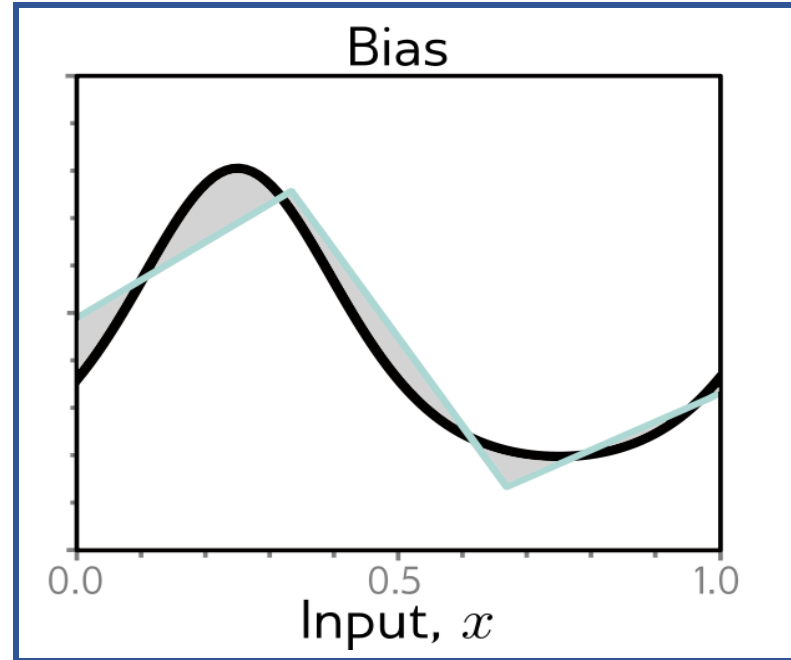
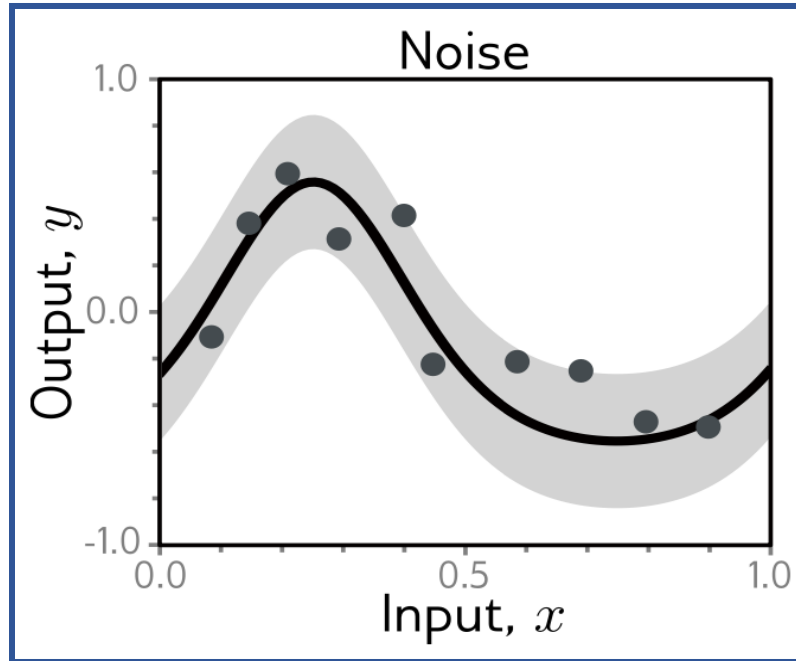
Which model performance is better?

Training Loss?



Model failed to generalize to unseen data

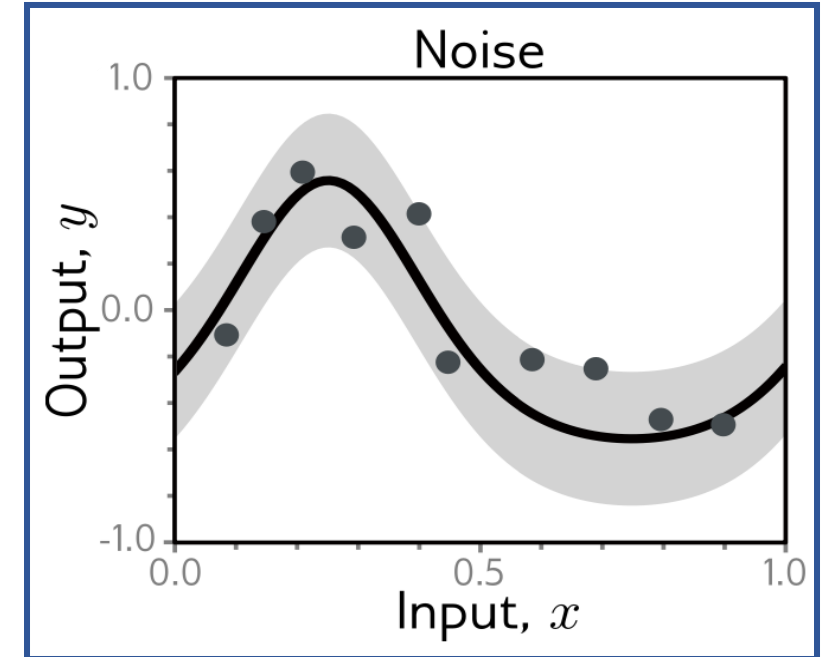
❑ Source of Error



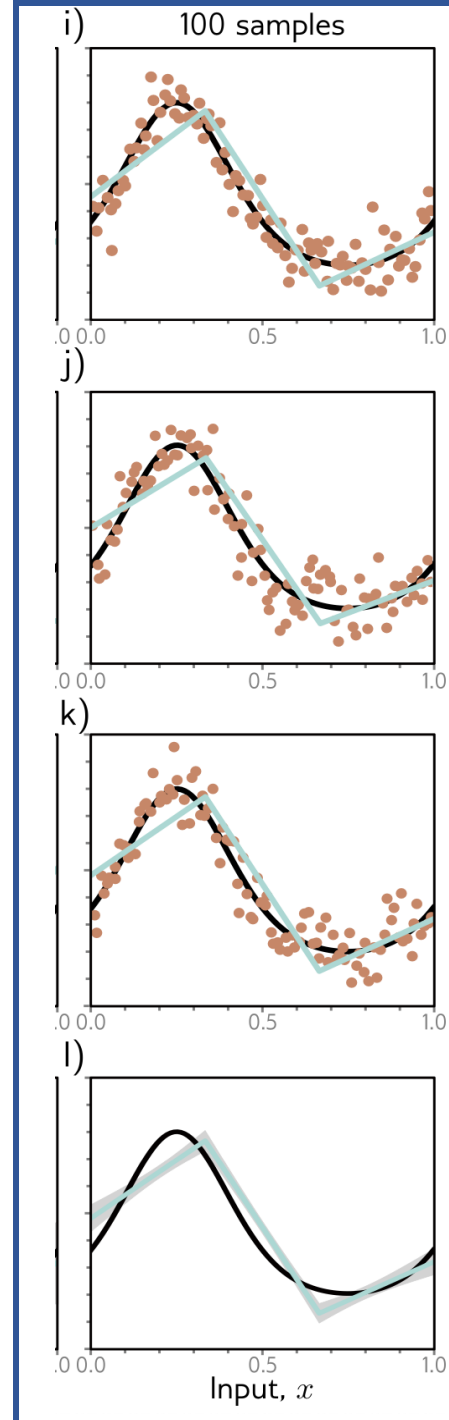
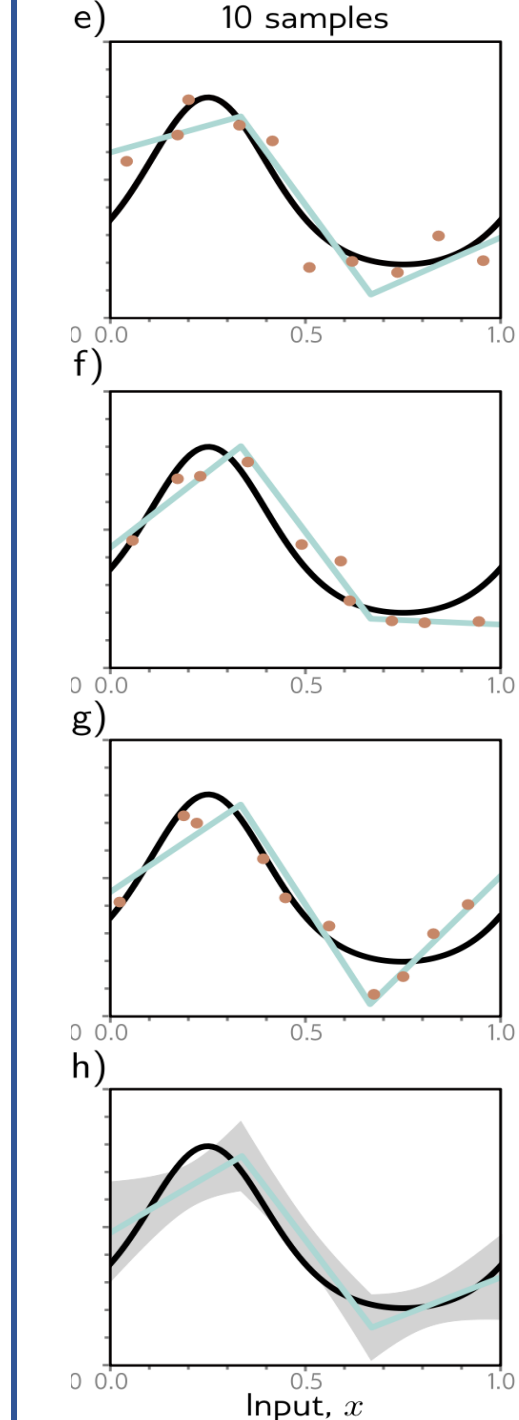
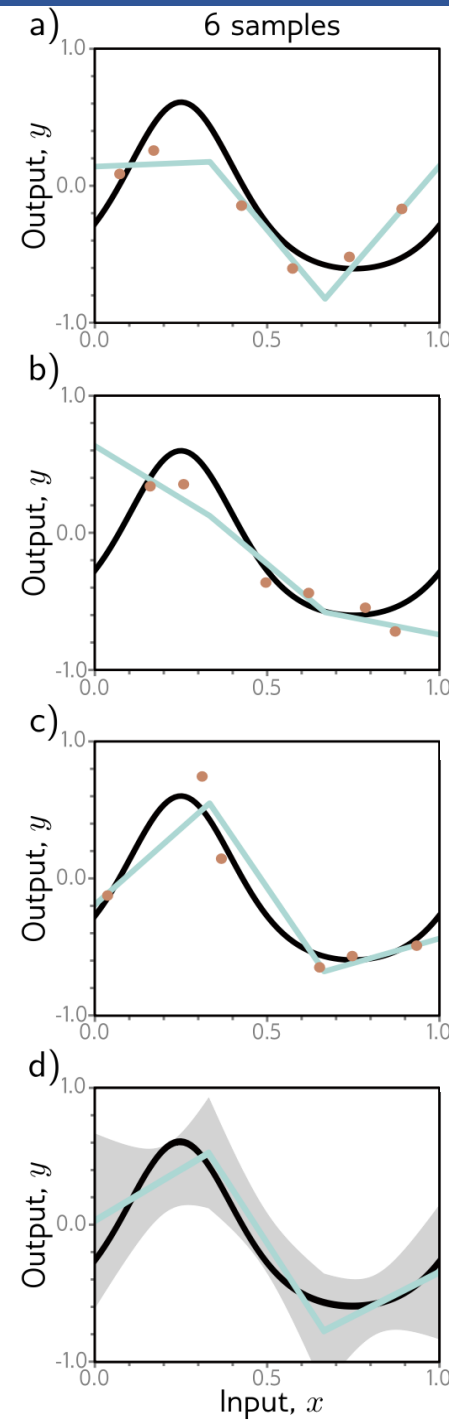
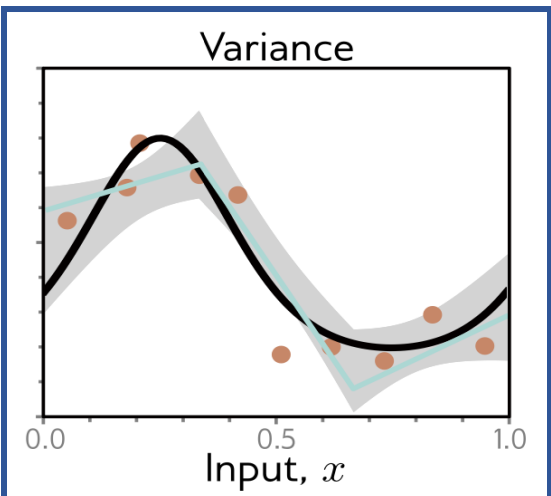
□ Mathematical formulation of Test Error

❑ Reducing Noise

- **Nothing** much we can do there.
- **Fundamental limit** on expected model performance.

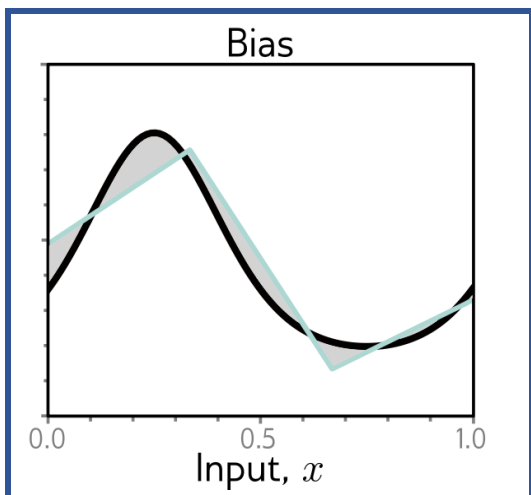


❑ Reducing Variance

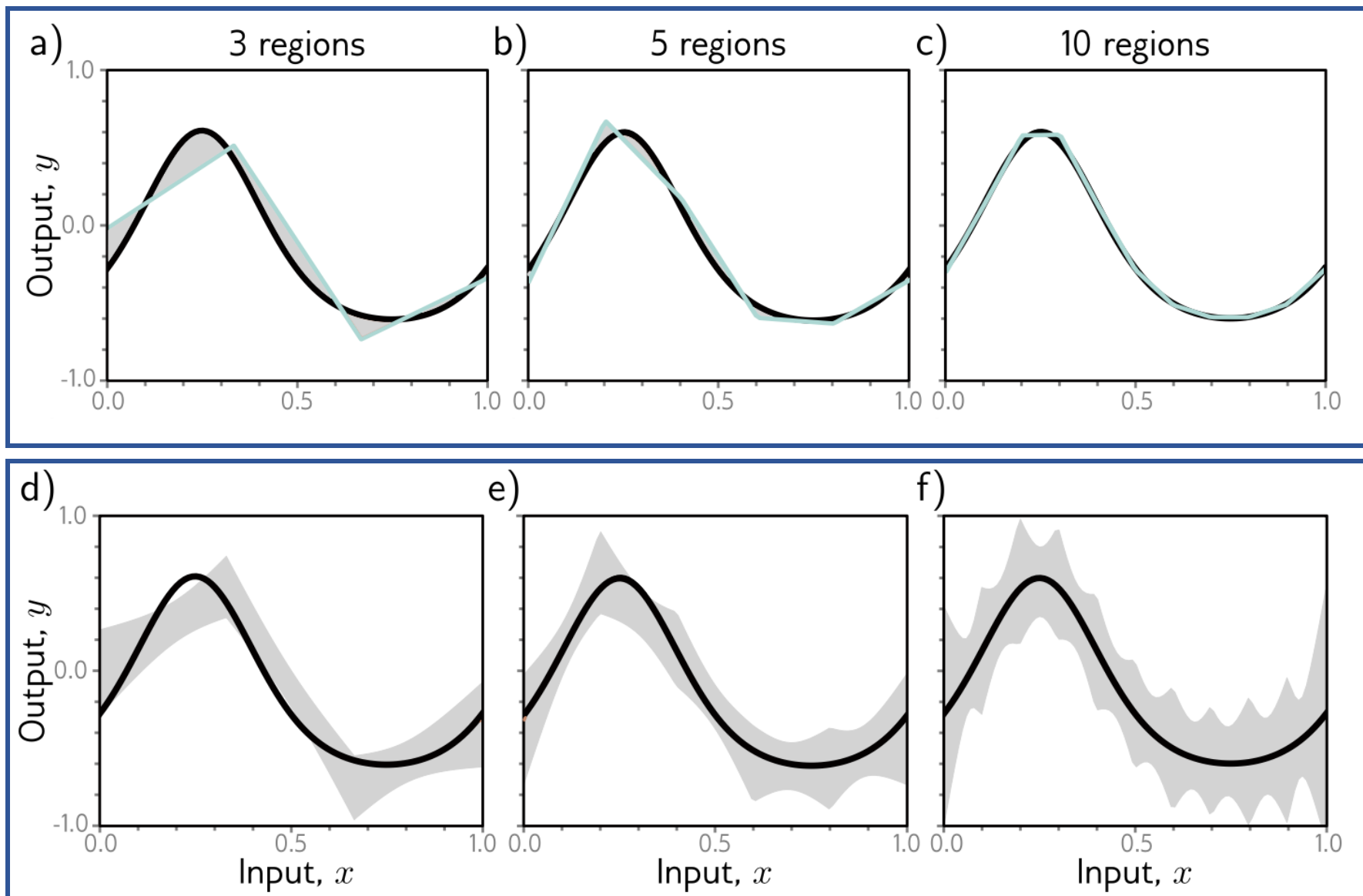


Model Performance

❑ Reducing Bias

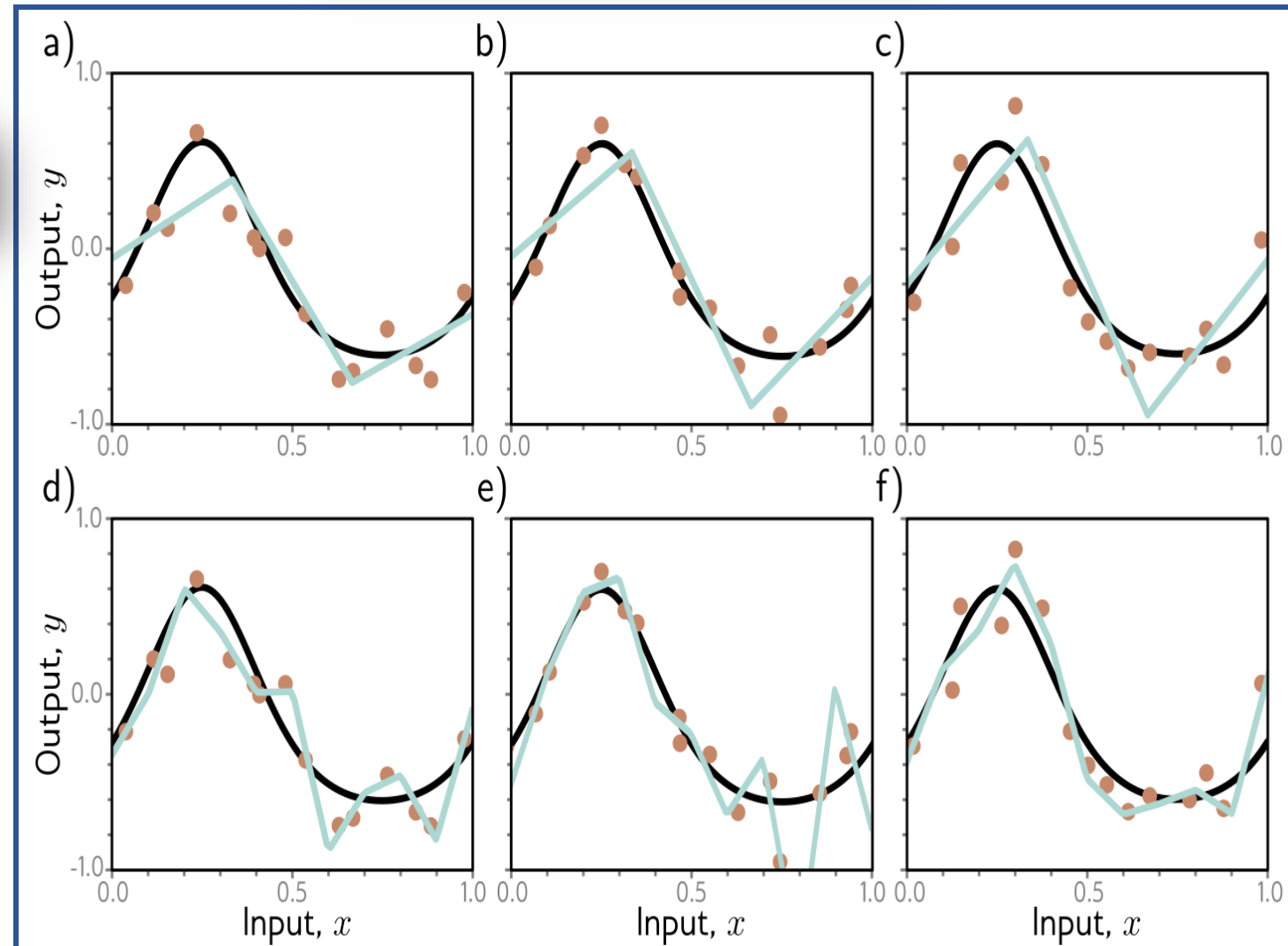
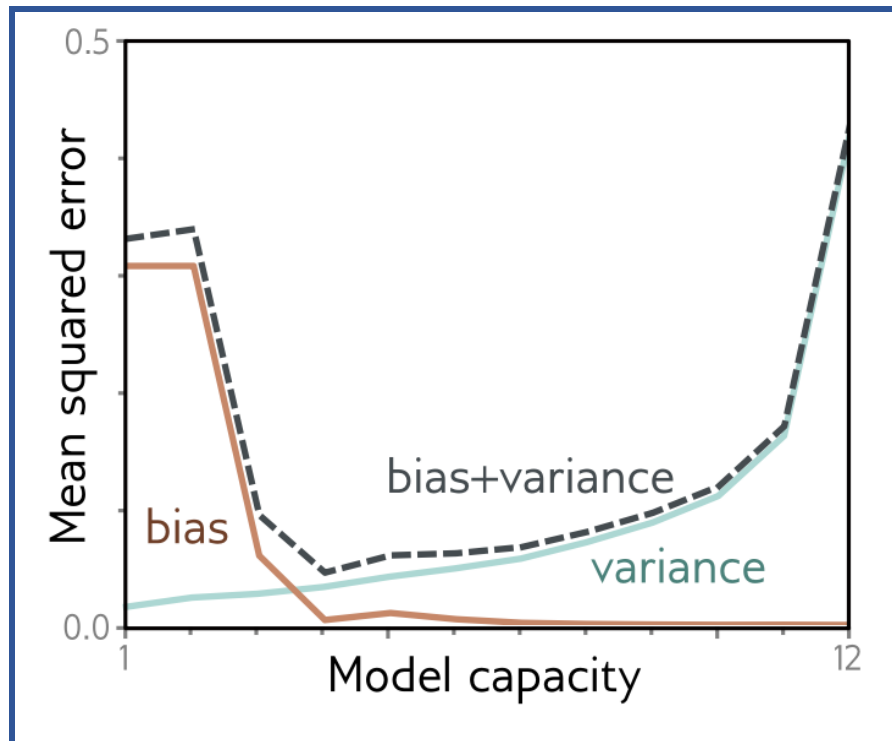


**Bias-Variance
Trade-off**



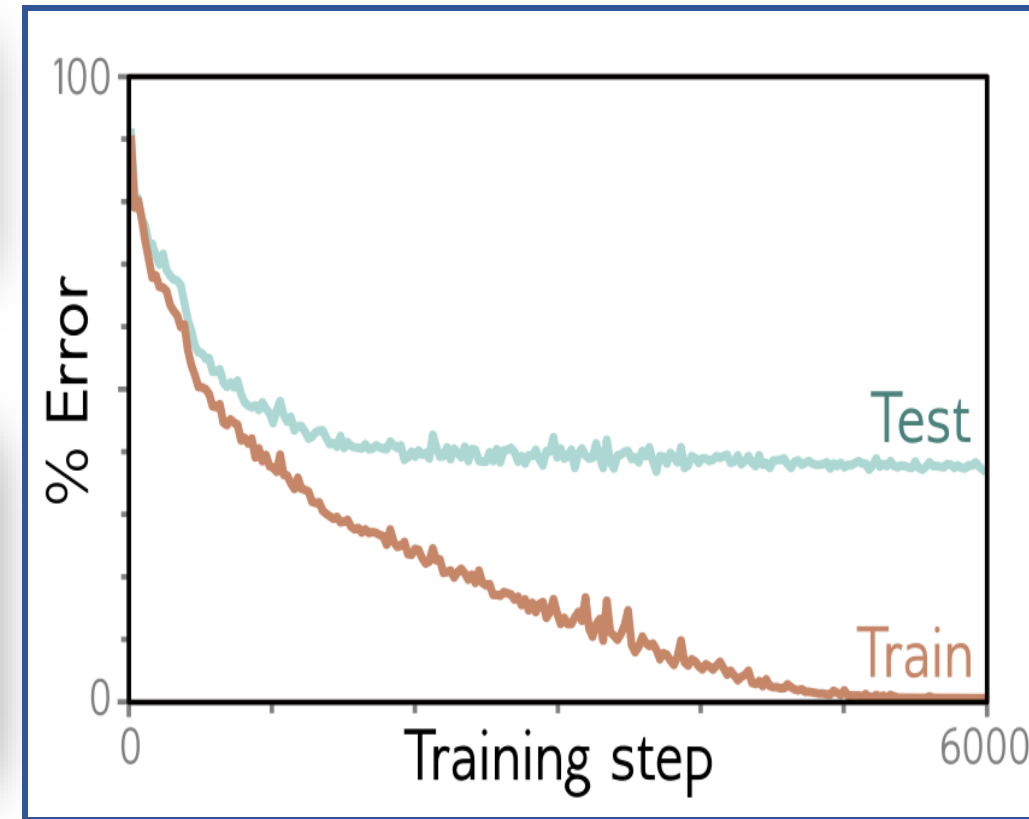
❑ Bias-Variance Trade-Off

What's this phenomenon called?



What is the main challenge?

Reduce gap between **training** and **test performance**.



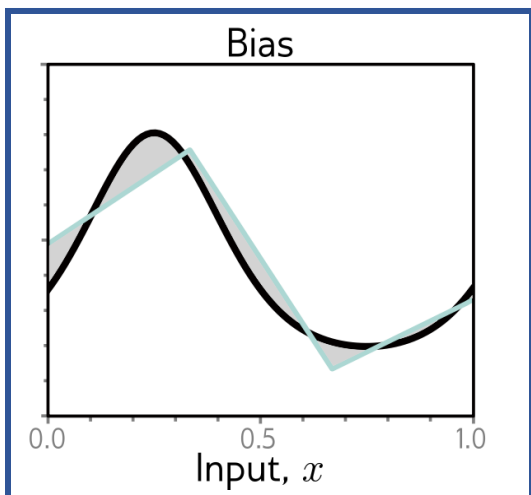
Class Announcement

□ Assignment – 1 will be released today in Piazza.

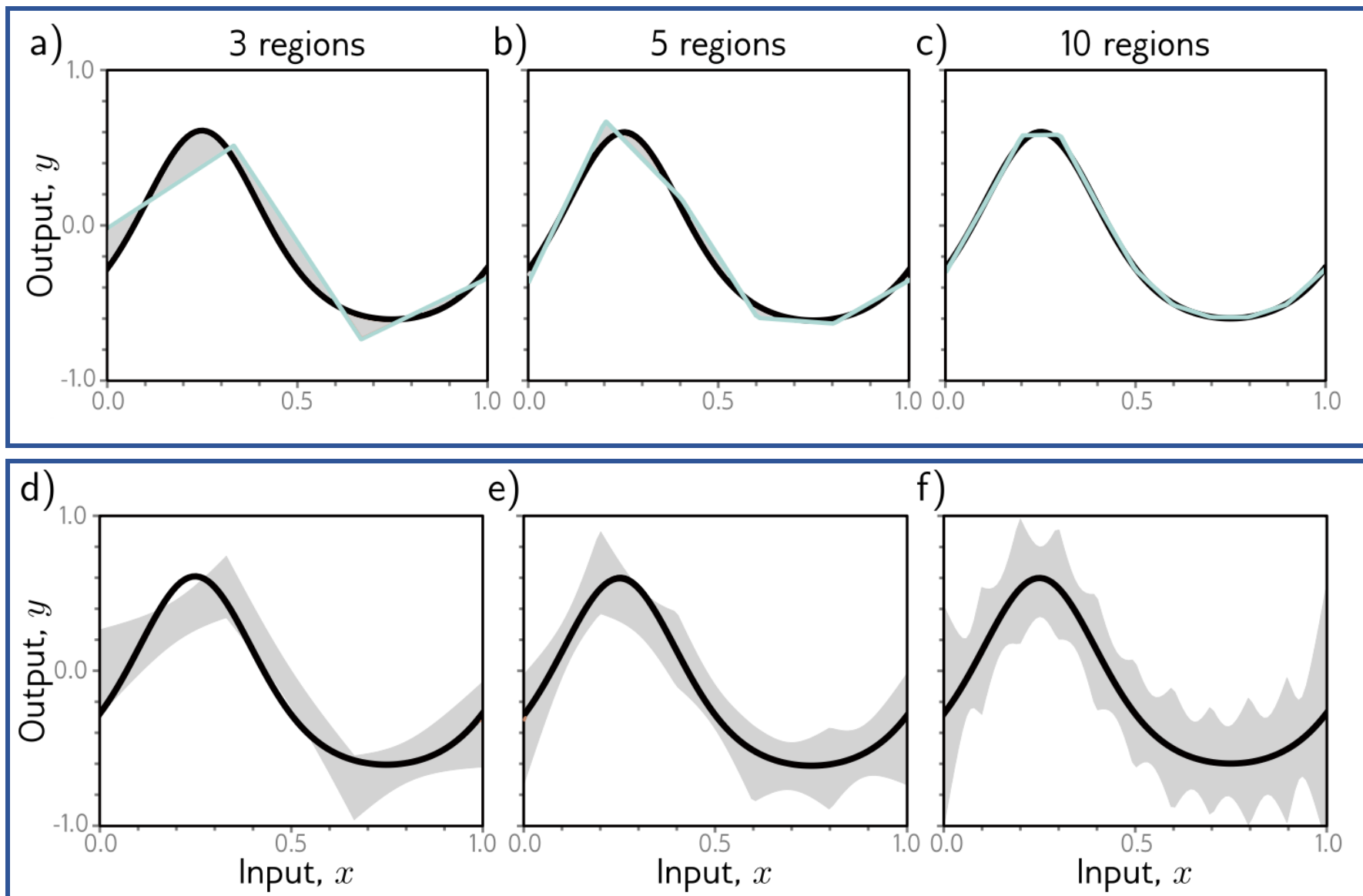
- **Deadline:** 7th Feb 11:59 pm.
- Briefly discuss today if there is time.

Model Performance

❑ Reducing Bias

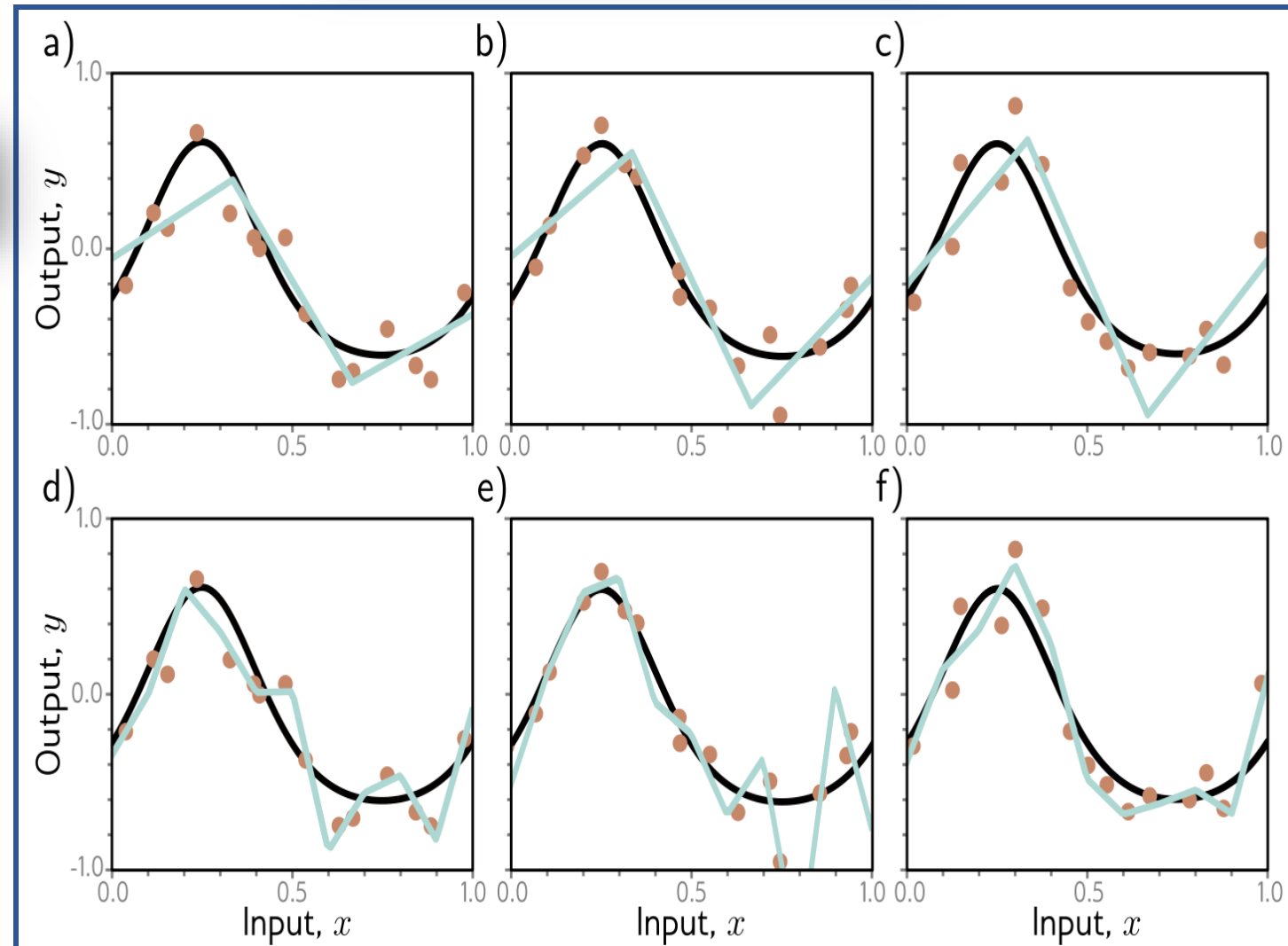
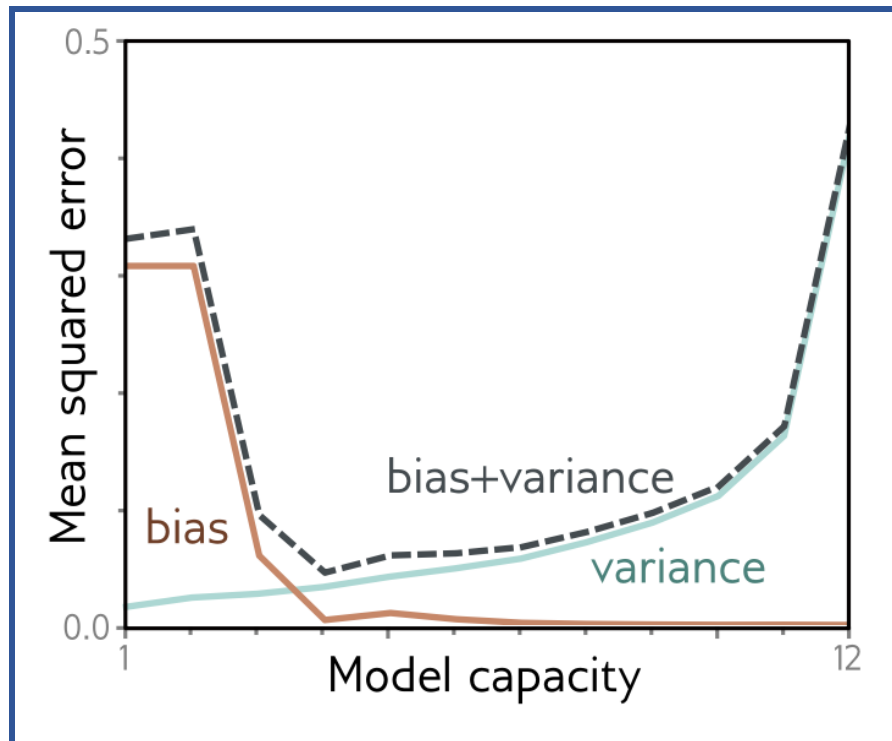


**Bias-Variance
Trade-off**



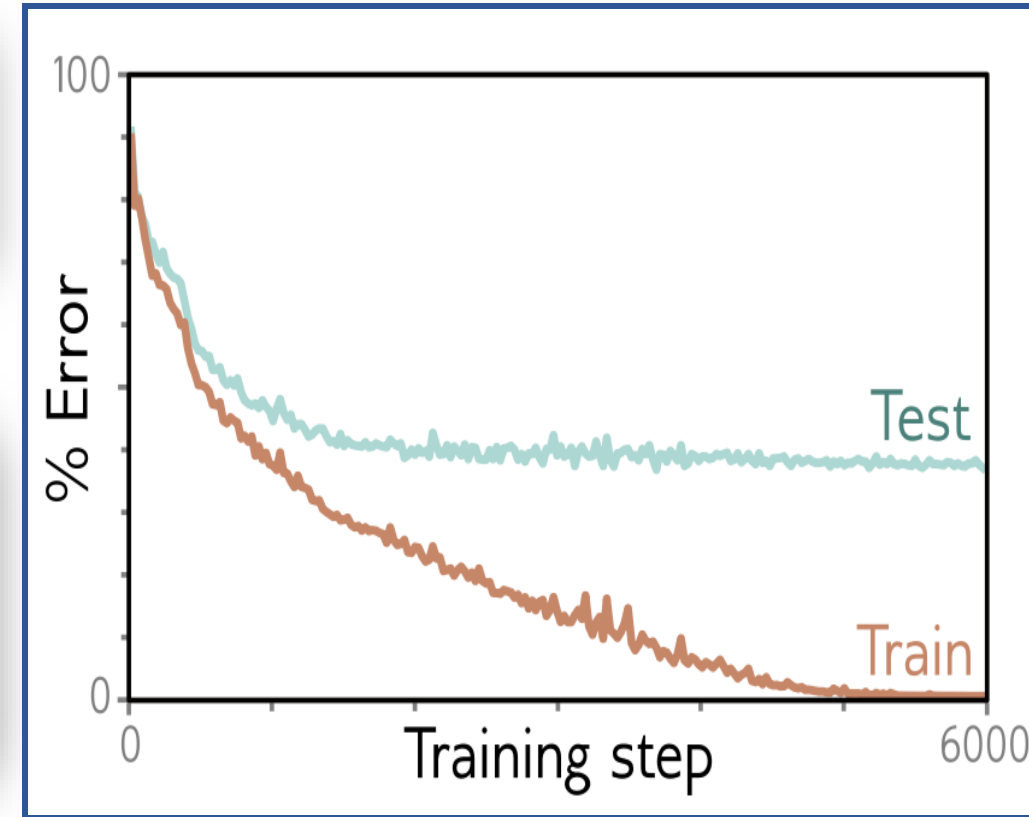
❑ Bias-Variance Trade-Off

What's this phenomenon called?



What is the main challenge?

Reduce gap between **training** and **test performance**.



Regularization

- ❑ **Family of methods** that reduce the **generalization gap** between training and test performance.
- ❑ Involves adding **explicit terms** to the loss function.

□ Explicit Regularization Technique

$$\begin{aligned}\hat{\phi} &= \underset{\phi}{\operatorname{argmin}} [L[\phi]] \\ &= \underset{\phi}{\operatorname{argmin}} \left[\sum_{i=1}^I \ell_i[\mathbf{x}_i, \mathbf{y}_i] \right]\end{aligned}$$

- Additional regularization term:

Why & How?

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \left[\sum_{i=1}^I \ell_i[\mathbf{x}_i, \mathbf{y}_i] + \lambda \cdot g[\phi] \right]$$

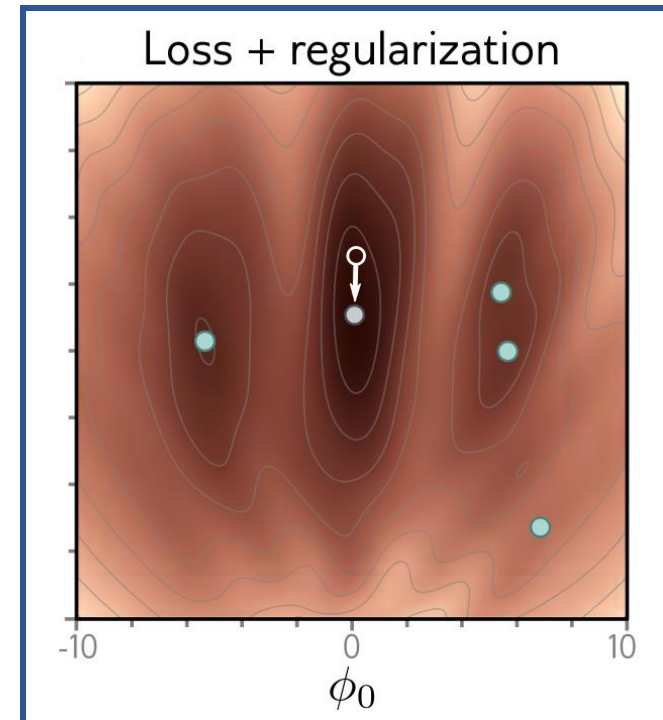
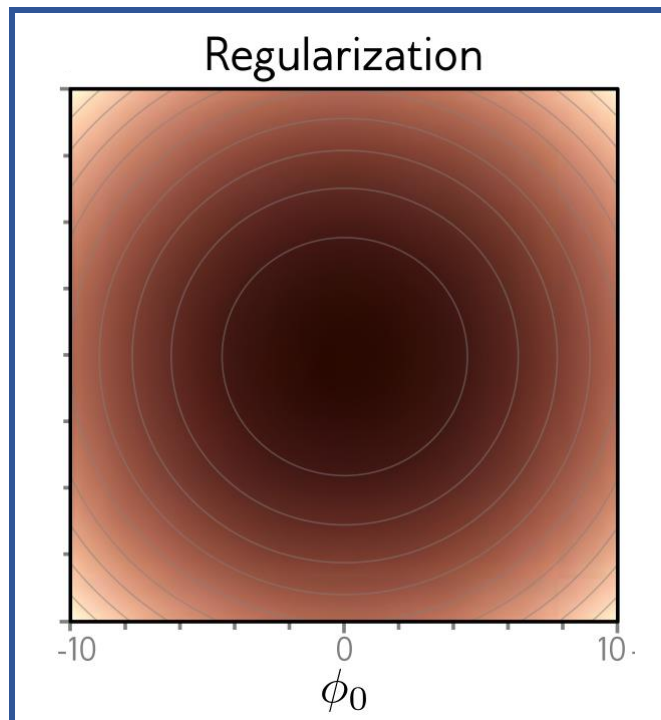
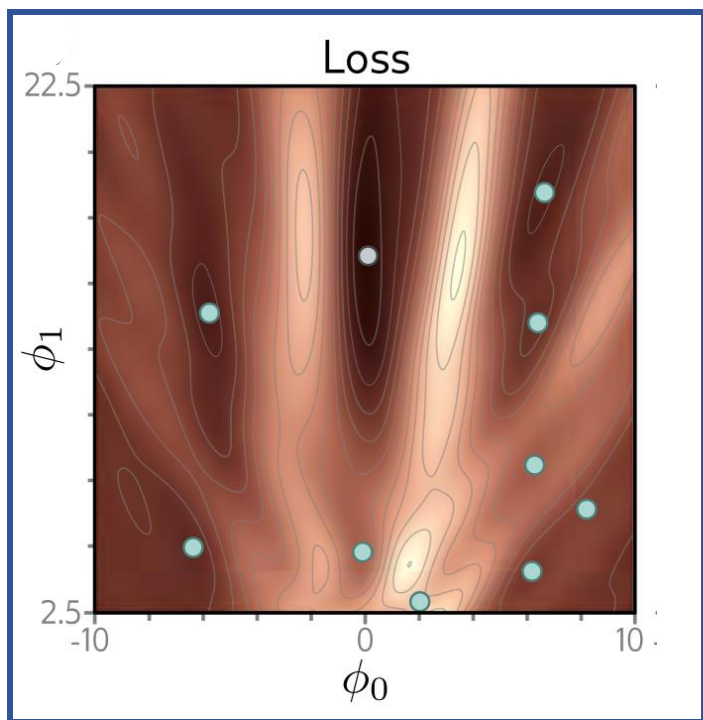
Constrained Optimization

Intuition?

Takes a high value when parameter is not desired.

□ Explicit Regularization

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \left[\sum_{i=1}^I \ell_i[\mathbf{x}_i, \mathbf{y}_i] + \lambda \cdot g[\phi] \right]$$



Regularization

□ Explicit Regularization

- Probabilistic interpretation:

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \left[\sum_{i=1}^I \ell_i[\mathbf{x}_i, \mathbf{y}_i] + \lambda \cdot g[\phi] \right]$$

$$\hat{\phi} = \underset{\phi}{\operatorname{argmax}} \left[\prod_{i=1}^I \operatorname{Pr}(\mathbf{y}_i | \mathbf{x}_i, \phi) \right]$$

Maximum Likelihood Criterion

- What if we have some knowledge about the parameters?

prior $\operatorname{Pr}(\phi)$

$$\hat{\phi} = \underset{\phi}{\operatorname{argmax}} \left[\prod_{i=1}^I \operatorname{Pr}(\mathbf{y}_i | \mathbf{x}_i, \phi) \operatorname{Pr}(\phi) \right]$$

Maximum a posteriori (MAP) Criterion

Negative Log-Likelihood (NLL)

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \left[-\log \left(\prod_i \operatorname{Pr}(y_i | x_i, \phi) \cdot \operatorname{Pr}(\phi) \right) \right]$$

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \left[\sum_i -\log \operatorname{Pr}(y_i | x_i, \phi) - \log \operatorname{Pr}(\phi) \right]$$

$$\lambda \cdot g[\phi] = -\log \operatorname{Pr}(\phi)$$

Regularization

□ L2 Regularization

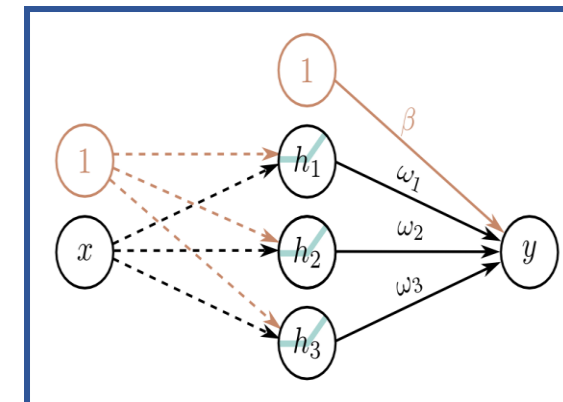
L2 Norm

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \left[\sum_{i=1}^I \ell_i[\mathbf{x}_i, \mathbf{y}_i] + \lambda \sum_j \phi_j^2 \right]$$

- Applied only to weights, not for biases.

What happens to the weights?

- Encourage smaller weights.
- Leads to smoother functions.



$$y = \sum_i w_i \cdot h_i + \beta$$

Regularization

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \left[\sum_{i=1}^I \ell_i[\mathbf{x}_i, \mathbf{y}_i] + \lambda \sum_j \phi_j^2 \right]$$



ScAI

□ L2 Regularization

How avoiding overfitting help us?

- Improves test error.

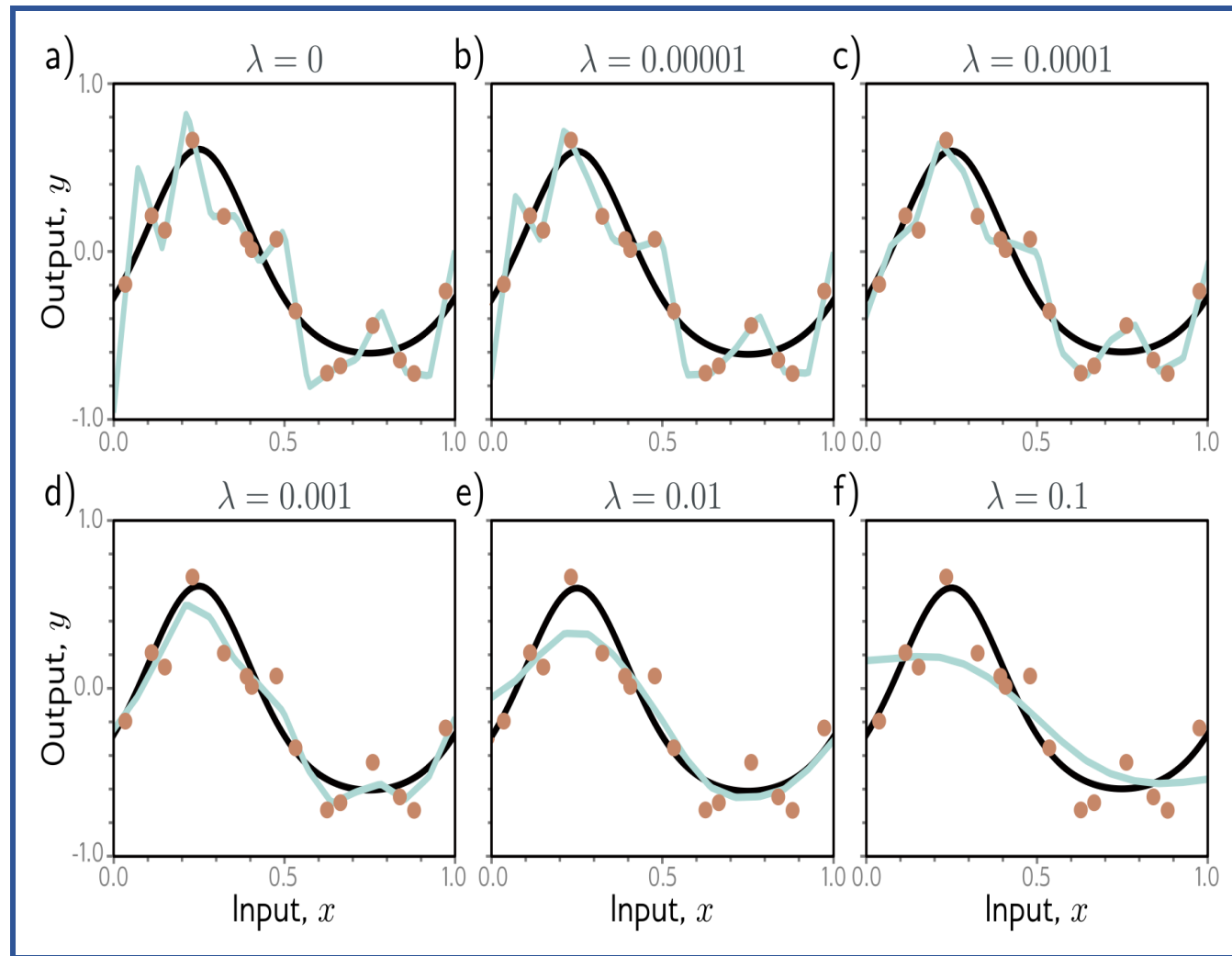
- Noise / Bias / Variance?

Bias Error -> Overfitting -> Variance Error

- Missing data

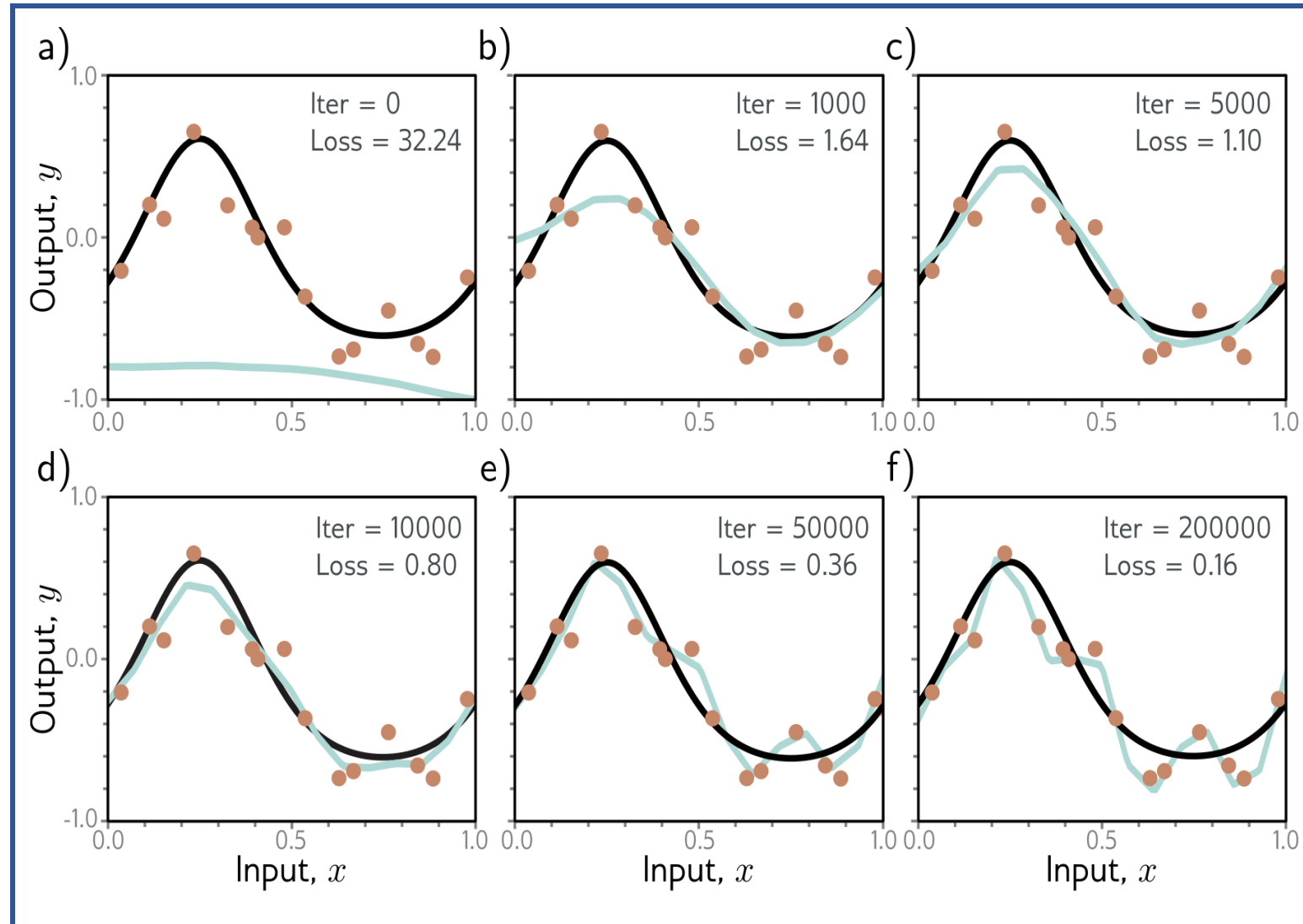
Complex vs Smooth function

How smooth fn help us?

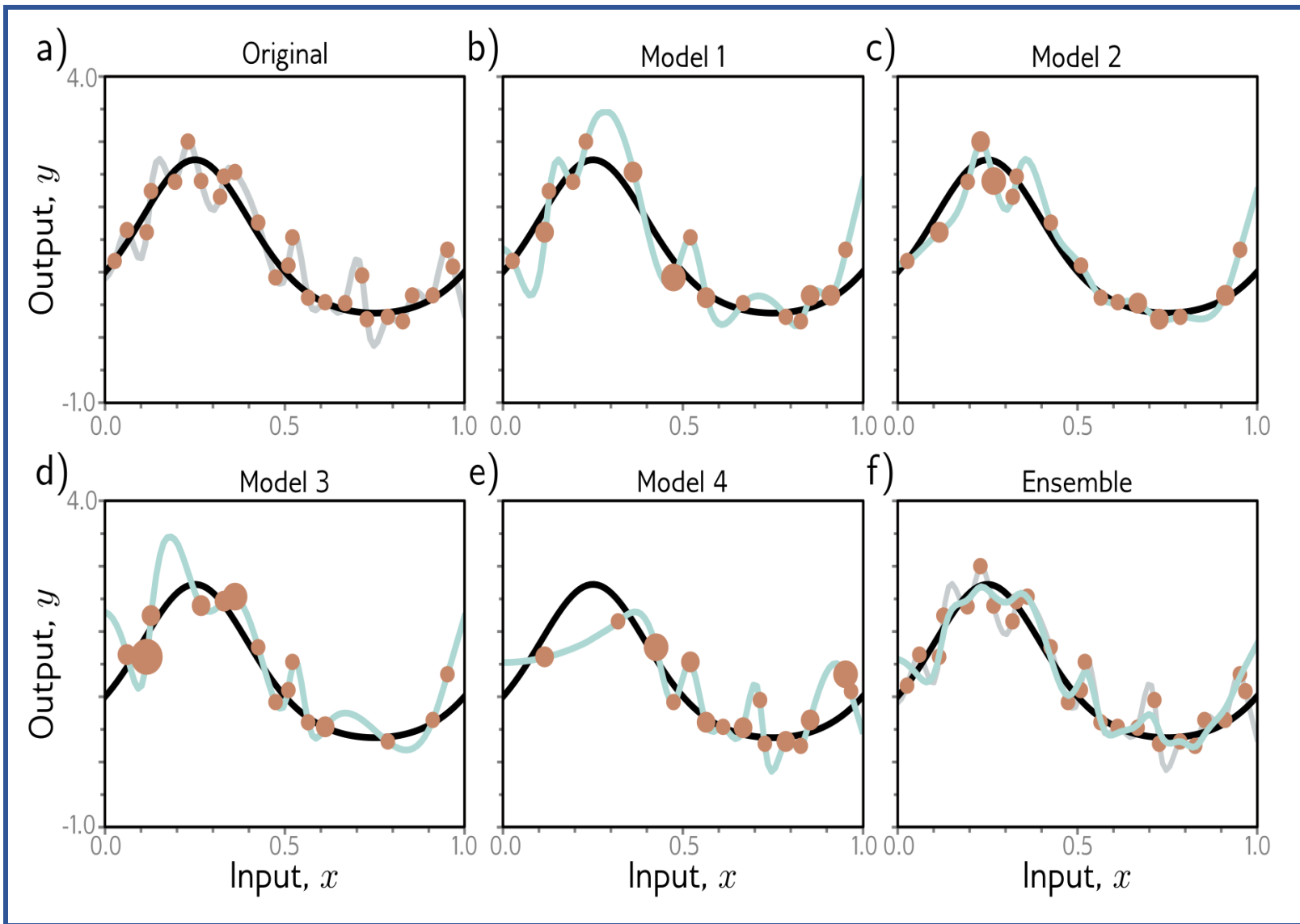


Smoother functions help alleviate overfitting

❑ Implicit Regularization Technique: Early Stopping



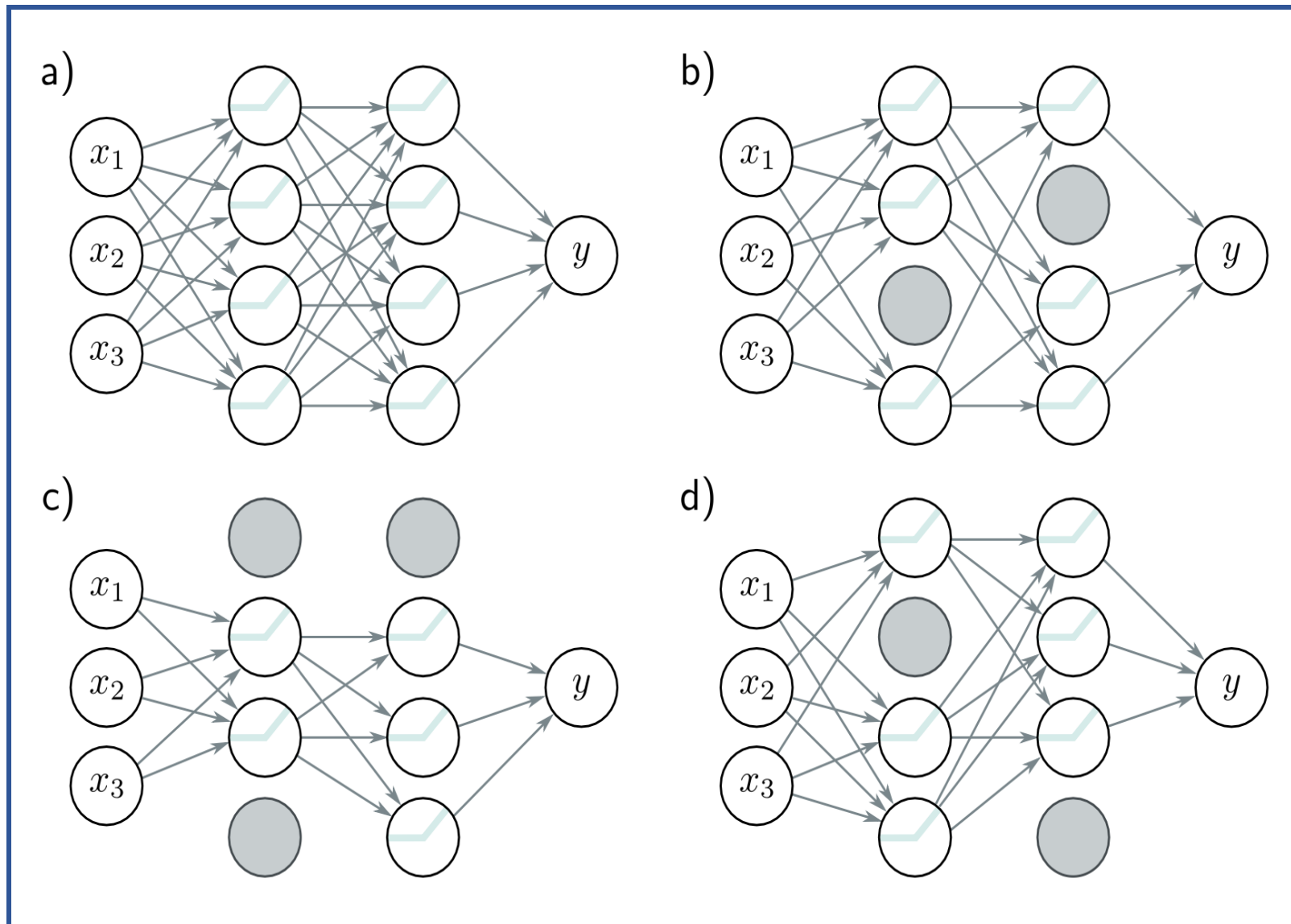
❑ Implicit Regularization Technique: Ensembling



Regularization

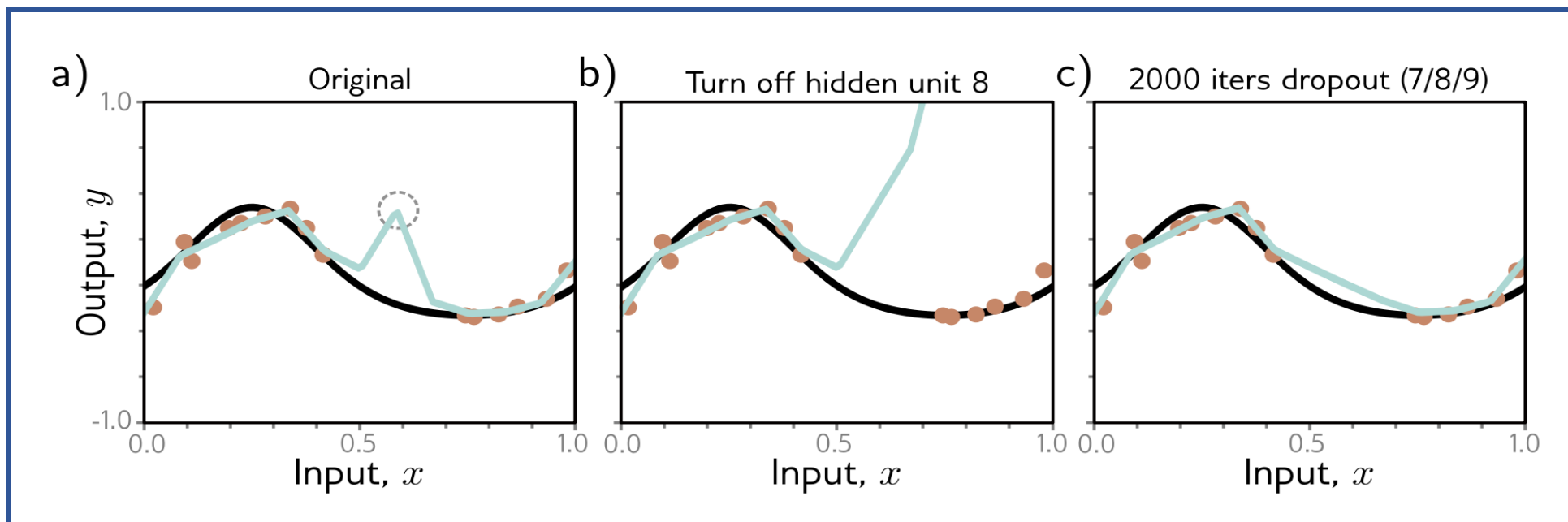
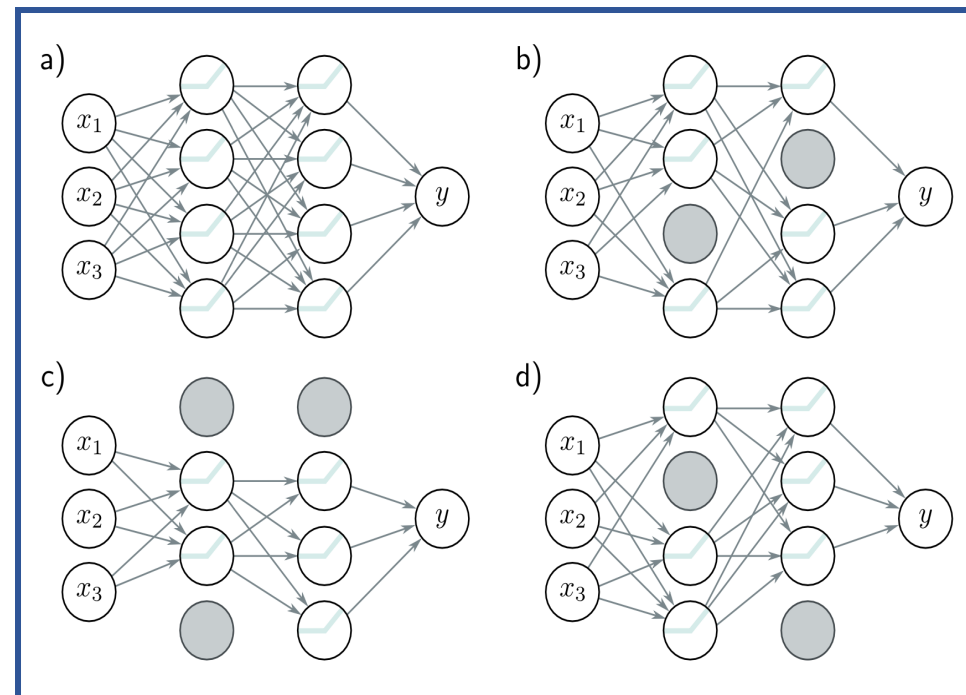


❑ Implicit Regularization Technique: Dropout

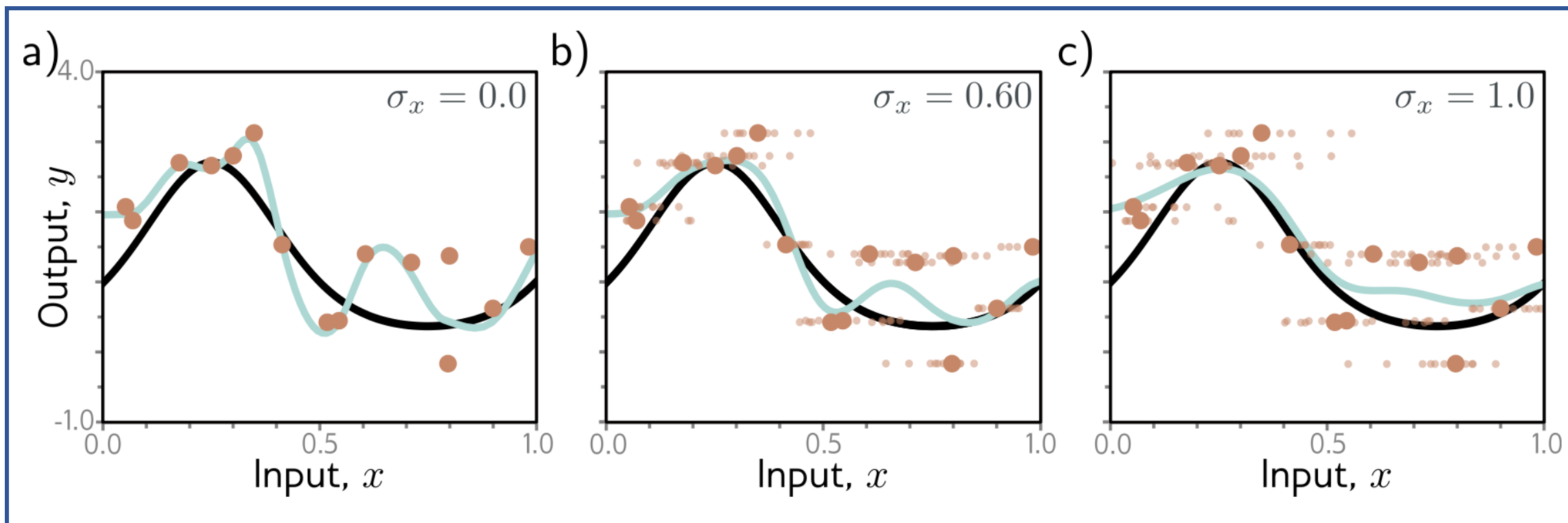


Regularization

❑ Implicit Regularization Technique: Dropout



❑ Implicit Regularization Technique: Applying Noise



END