

# CM20315 - Machine Learning

Prof. Simon Prince

9. Regularization



# Regularization

- Why is there a generalization gap between training and test data?
  - Overfitting (model describes statistical peculiarities)
  - Model unconstrained in areas where there are no training examples
- **Regularization** = methods to reduce the generalization gap
- Technically means adding terms to loss function
- But colloquially means any method (hack) to reduce gap

# Regularization

- Explicit regularization
- Implicit regularization
- Early stopping
- Ensembling
- Dropout
- Adding noise
- Bayesian approaches
- Transfer learning, multi-task learning, self-supervised learning
- Data augmentation

# Explicit regularization

- Standard loss function:

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

# Explicit regularization

- Standard loss function:

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

- Regularization adds an extra term

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

# Explicit regularization

- Standard loss function:

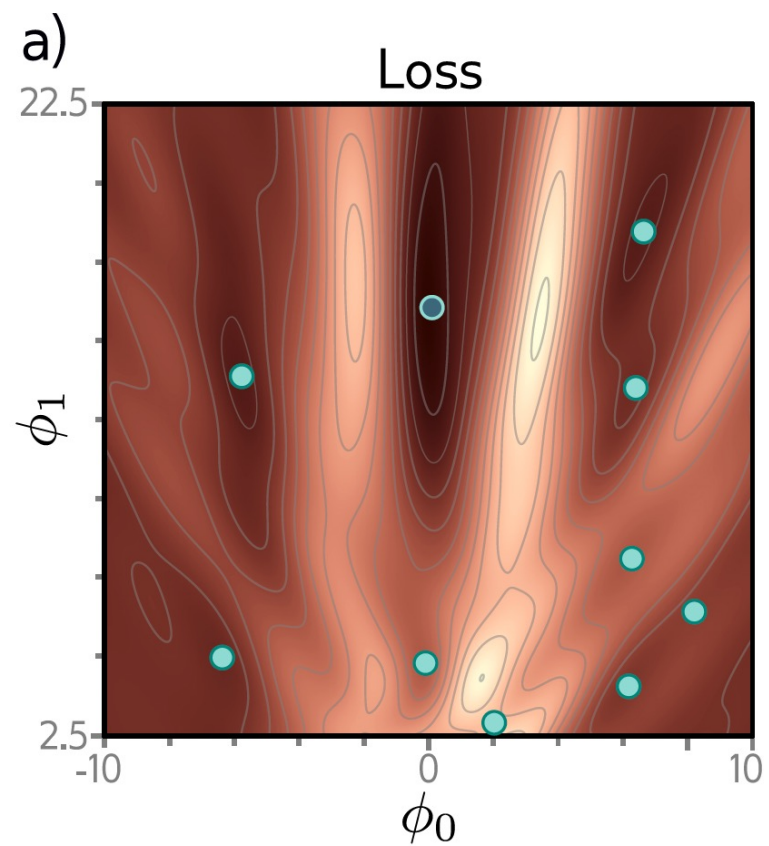
$$\begin{aligned}\hat{\phi} &= \operatorname{argmin}_{\phi} [L[\phi]] \\ &= \operatorname{argmin}_{\phi} \left[ \sum_{i=1}^I \ell_i[\mathbf{x}_i, \mathbf{y}_i] \right]\end{aligned}$$

- Regularization adds an extra term

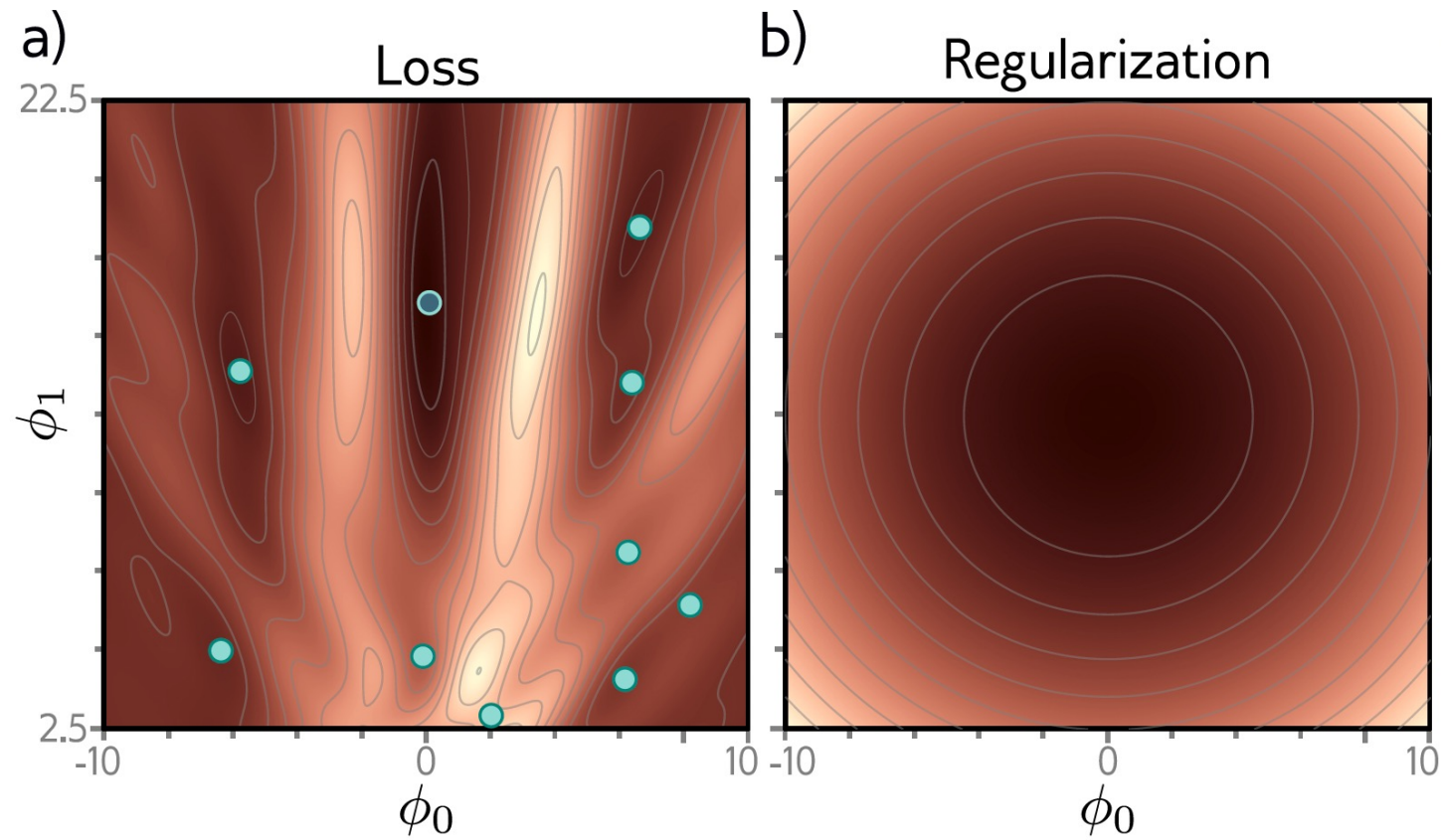
$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[ \sum_{i=1}^I \ell_i[\mathbf{x}_i, \mathbf{y}_i] + \lambda \cdot g[\phi] \right]$$

- Favors some parameters, disfavors others.
- $\lambda > 0$  controls the strength

# Explicit regularization

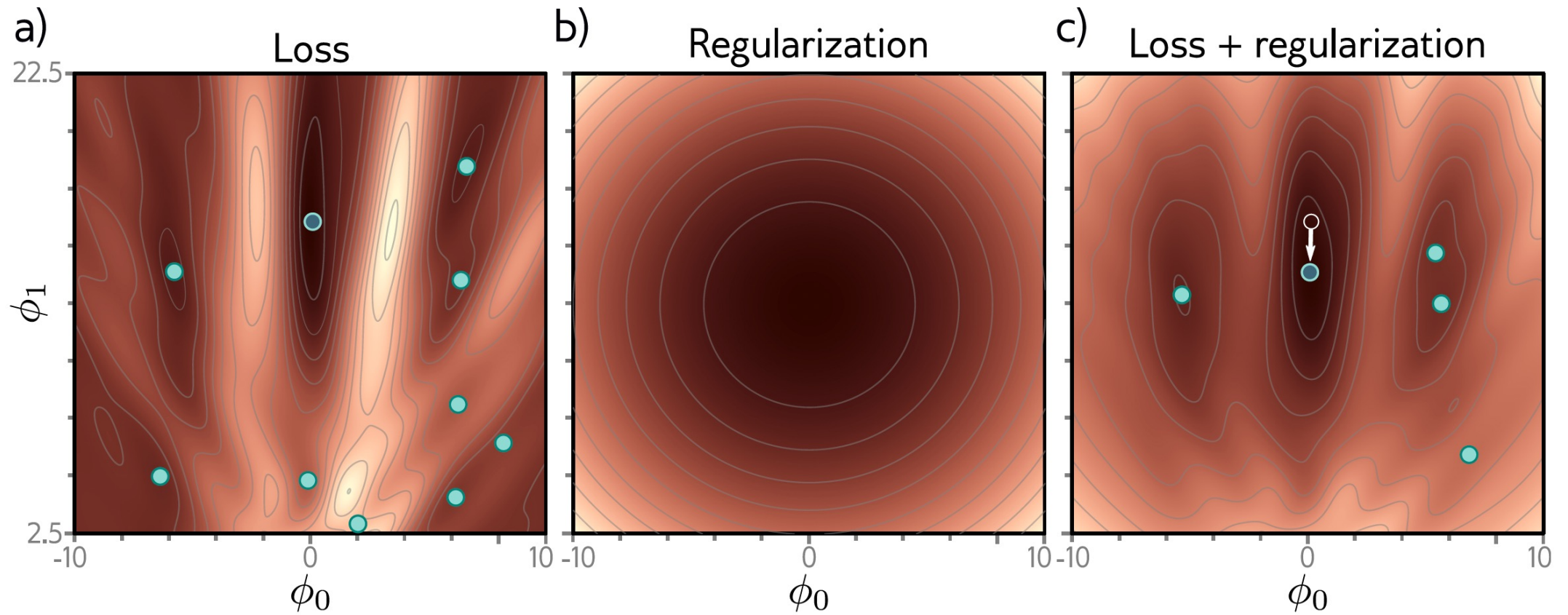


# Explicit regularization





# Explicit regularization



# Probabilistic interpretation

- Maximum likelihood:

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

- Regularization is equivalent to adding a **prior** over parameters

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

... what you know about parameters *before* seeing the data

# Equivalence

- Explicit regularization:

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

- Probabilistic interpretation:

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

# Equivalence

- Explicit regularization:

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

- Probabilistic interpretation:

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

- Mapping:

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

# L2 Regularization

- Can only use very general terms
- Most common is **L2 regularization**
- Favors smaller parameters

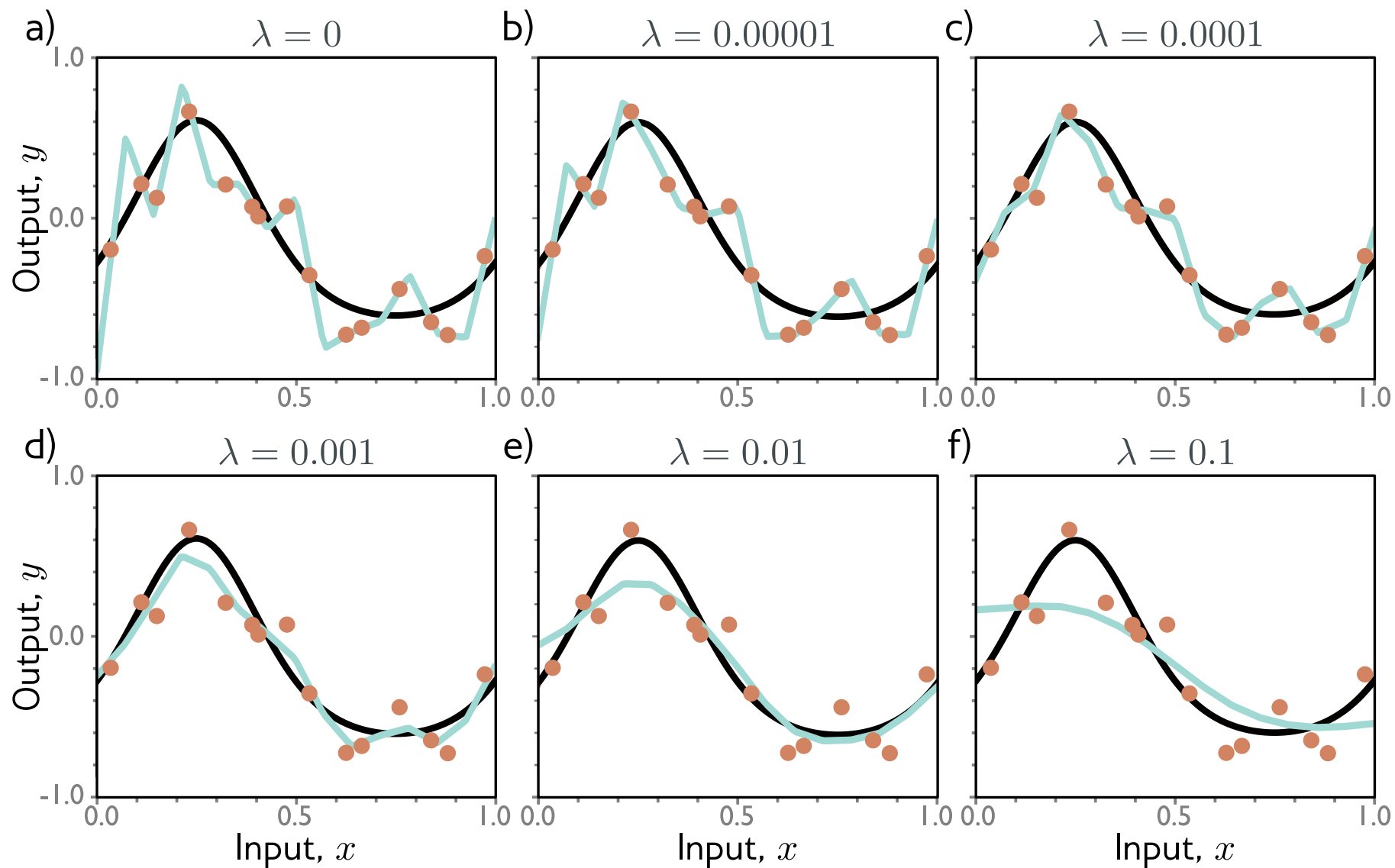
$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[ \mathbf{L}[\phi, \{\mathbf{x}_i, \mathbf{y}_i\}] + \lambda \sum_j \phi_j^2 \right]$$

- Also called **Tikhonov regularization, ridge regression**
- In neural networks, usually just for weights and called **weight decay**

# Why does L2 regularization help?

- Discourages overcommitment to the data (overfitting)
- Encourages smoothness between datapoints

# L2 regularization

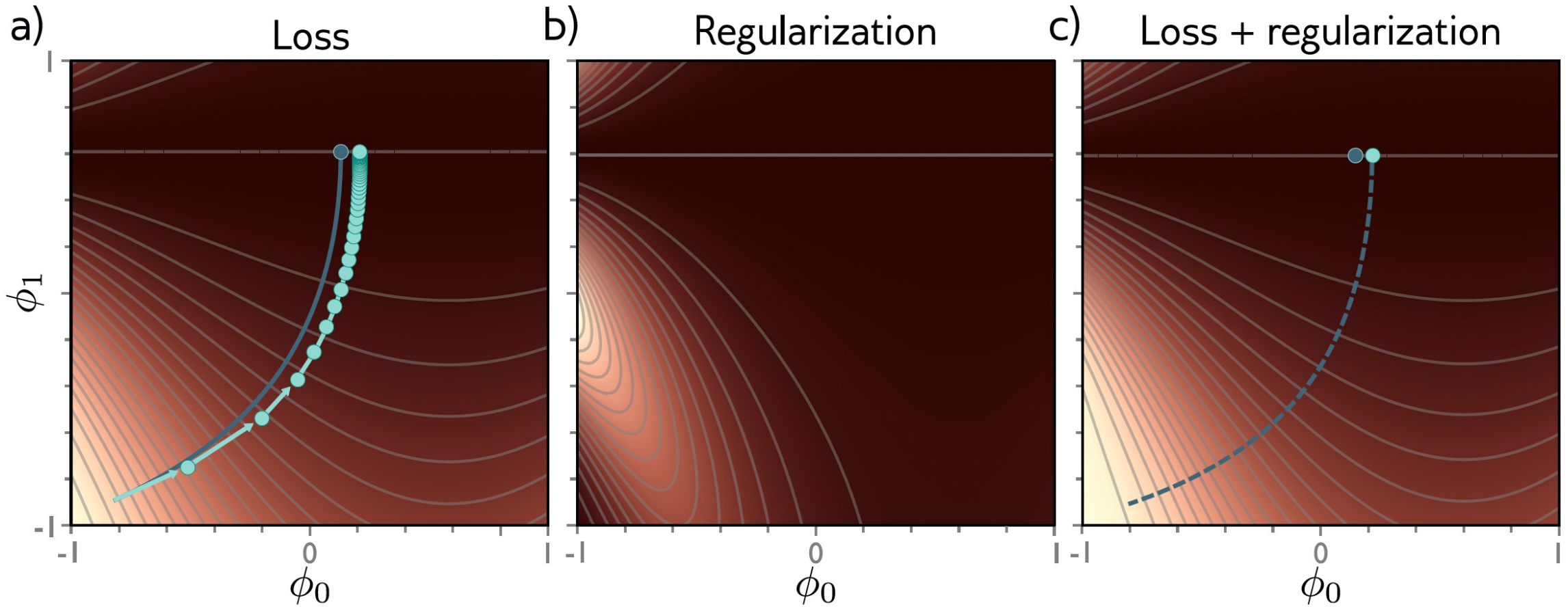


# Regularization

- Explicit regularization
- **Implicit regularization**
- Early stopping
- Ensembling
- Dropout
- Adding noise
- Bayesian approaches
- Transfer learning, multi-task learning, self-supervised learning
- Data augmentation



# Implicit regularization



Gradient descent approximates a differential equation (infinitesimal step size)

Finite step size equivalent to regularization

Add in that regularization and differential equation converges to same place

# Implicit regularization

- Gradient descent disfavors areas where gradients are steep

$$\tilde{L}_{GD}[\phi] = L[\phi] + \frac{\alpha}{4} \left\| \frac{\partial L}{\partial \phi} \right\|^2$$

# Implicit regularization

- Gradient descent disfavors areas where gradients are steep

$$\tilde{L}_{GD}[\phi] = L[\phi] + \frac{\alpha}{4} \left\| \frac{\partial L}{\partial \phi} \right\|^2$$

- SGD likes all batches to have similar gradients

$$\begin{aligned} \tilde{L}_{SGD}[\phi] &= \tilde{L}_{GD}[\phi] + \frac{\alpha}{4B} \sum_{b=1}^B \left\| \frac{\partial L_b}{\partial \phi} - \frac{\partial L}{\partial \phi} \right\|^2 \\ &= L[\phi] + \frac{\alpha}{4} \left\| \frac{\partial L}{\partial \phi} \right\|^2 + \frac{\alpha}{4B} \sum_{b=1}^B \left\| \frac{\partial L_b}{\partial \phi} - \frac{\partial L}{\partial \phi} \right\|^2 \end{aligned}$$

# Implicit regularization

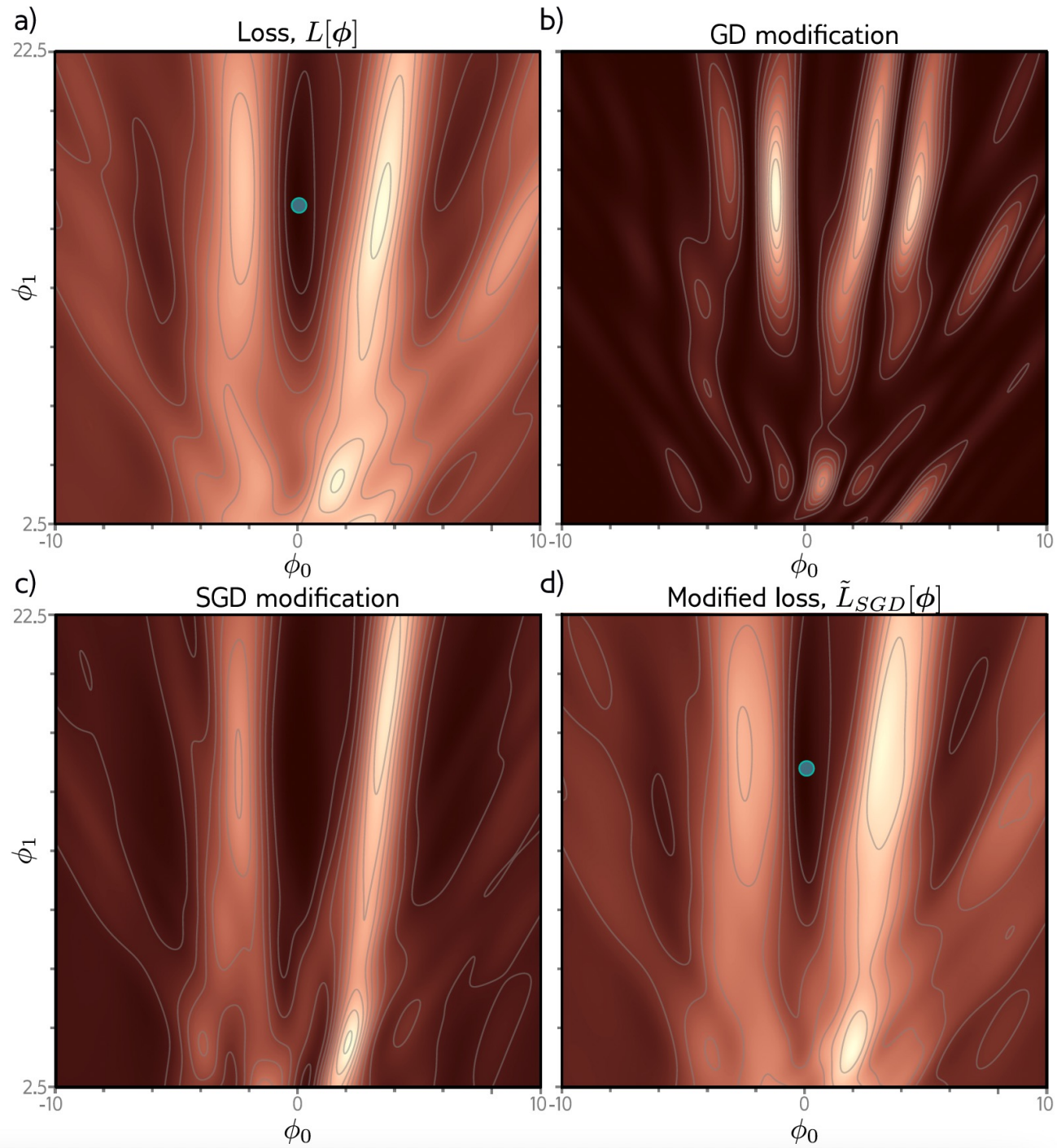
- Gradient descent disfavors areas where gradients are steep

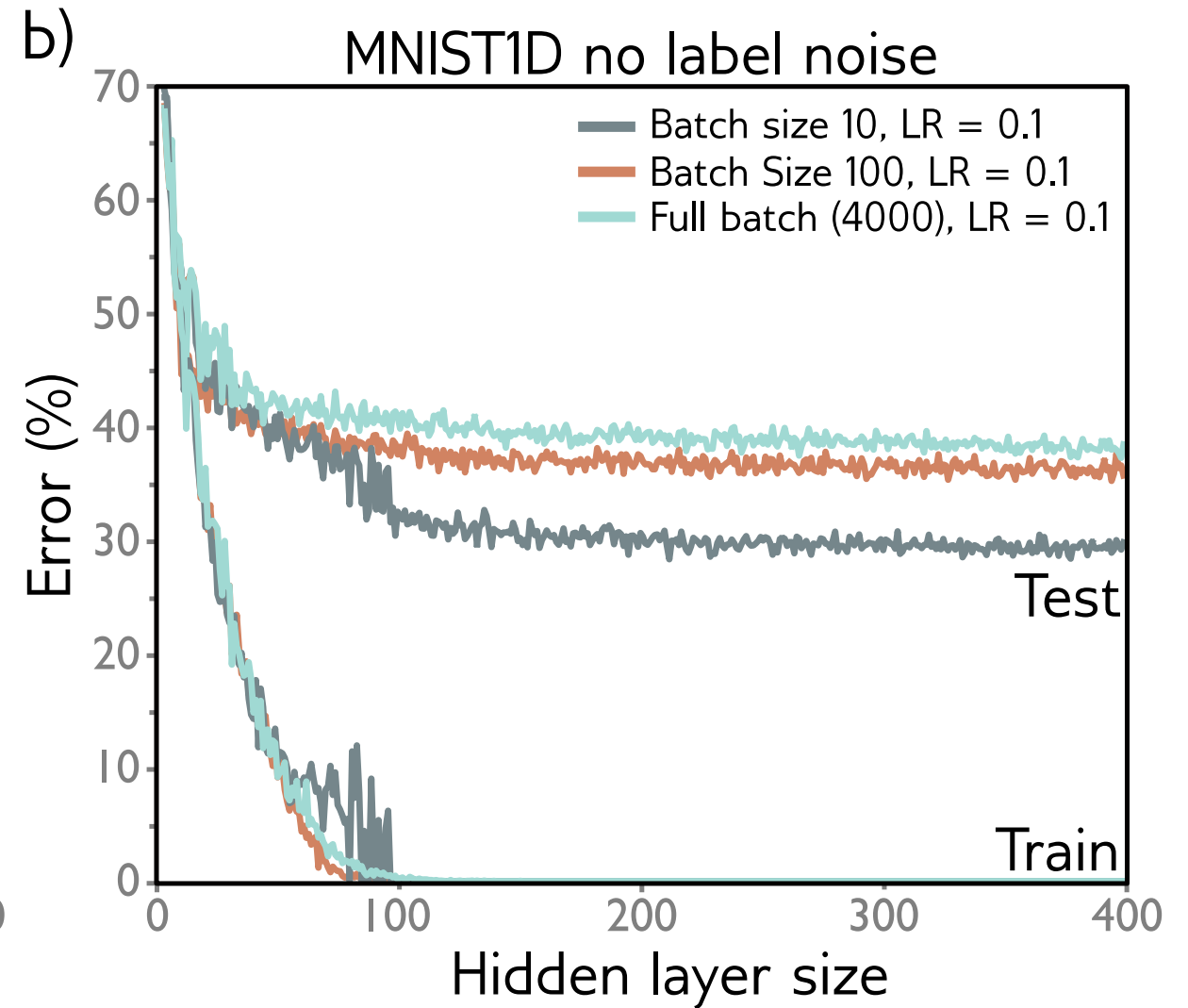
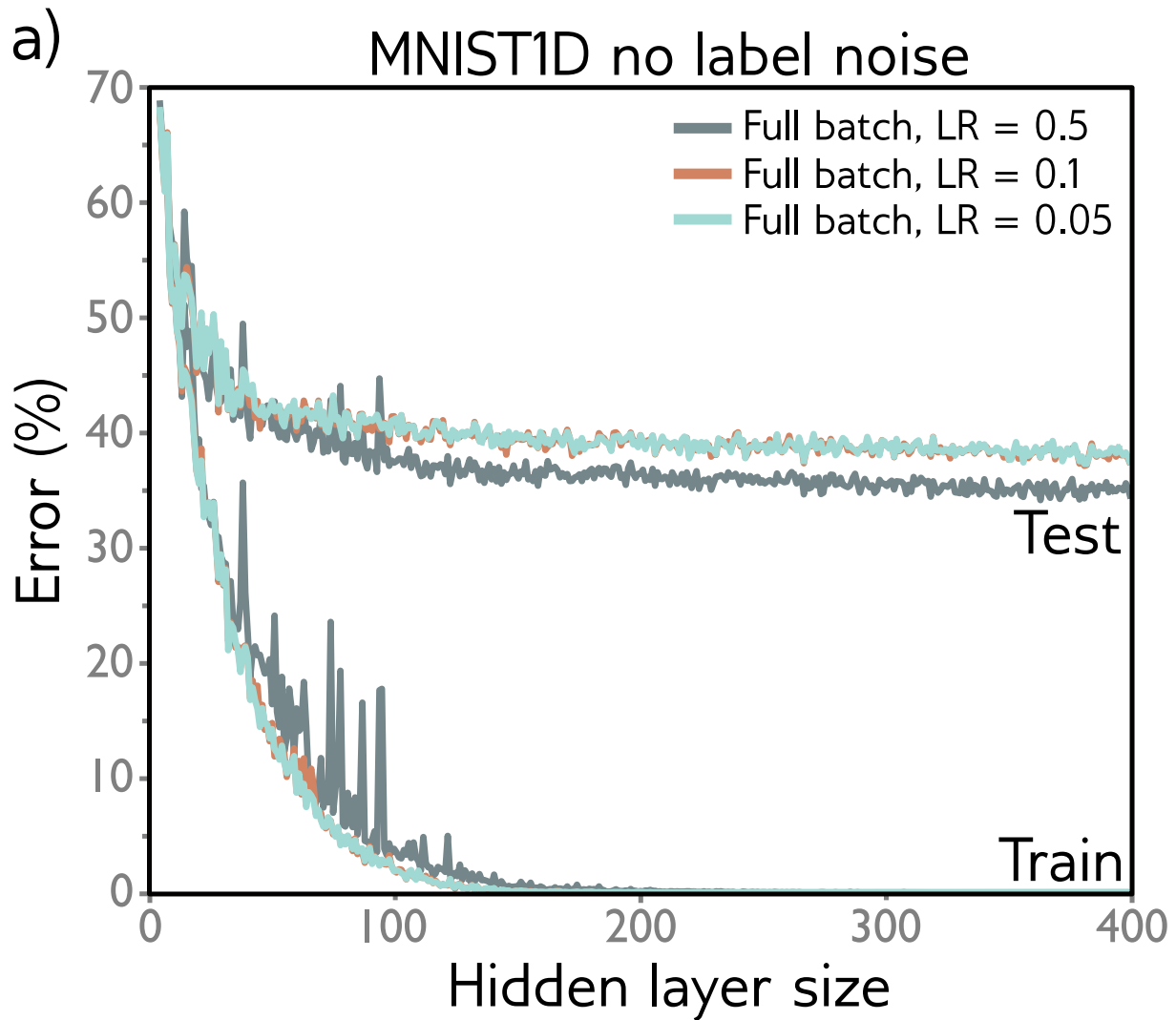
$$\tilde{L}_{GD}[\phi] = L[\phi] + \frac{\alpha}{4} \left\| \frac{\partial L}{\partial \phi} \right\|^2$$

- SGD likes all batches to have similar gradients

$$\begin{aligned} \tilde{L}_{SGD}[\phi] &= \tilde{L}_{GD}[\phi] + \frac{\alpha}{4B} \sum_{b=1}^B \left\| \frac{\partial L_b}{\partial \phi} - \frac{\partial L}{\partial \phi} \right\|^2 \\ &= L[\phi] + \frac{\alpha}{4} \left\| \frac{\partial L}{\partial \phi} \right\|^2 + \frac{\alpha}{4B} \sum_{b=1}^B \left\| \frac{\partial L_b}{\partial \phi} - \frac{\partial L}{\partial \phi} \right\|^2 \end{aligned}$$

- Depends on learning rate – perhaps why larger learning rates generalize better.





Generally performance

- best for larger learning rates
- best with smaller batches

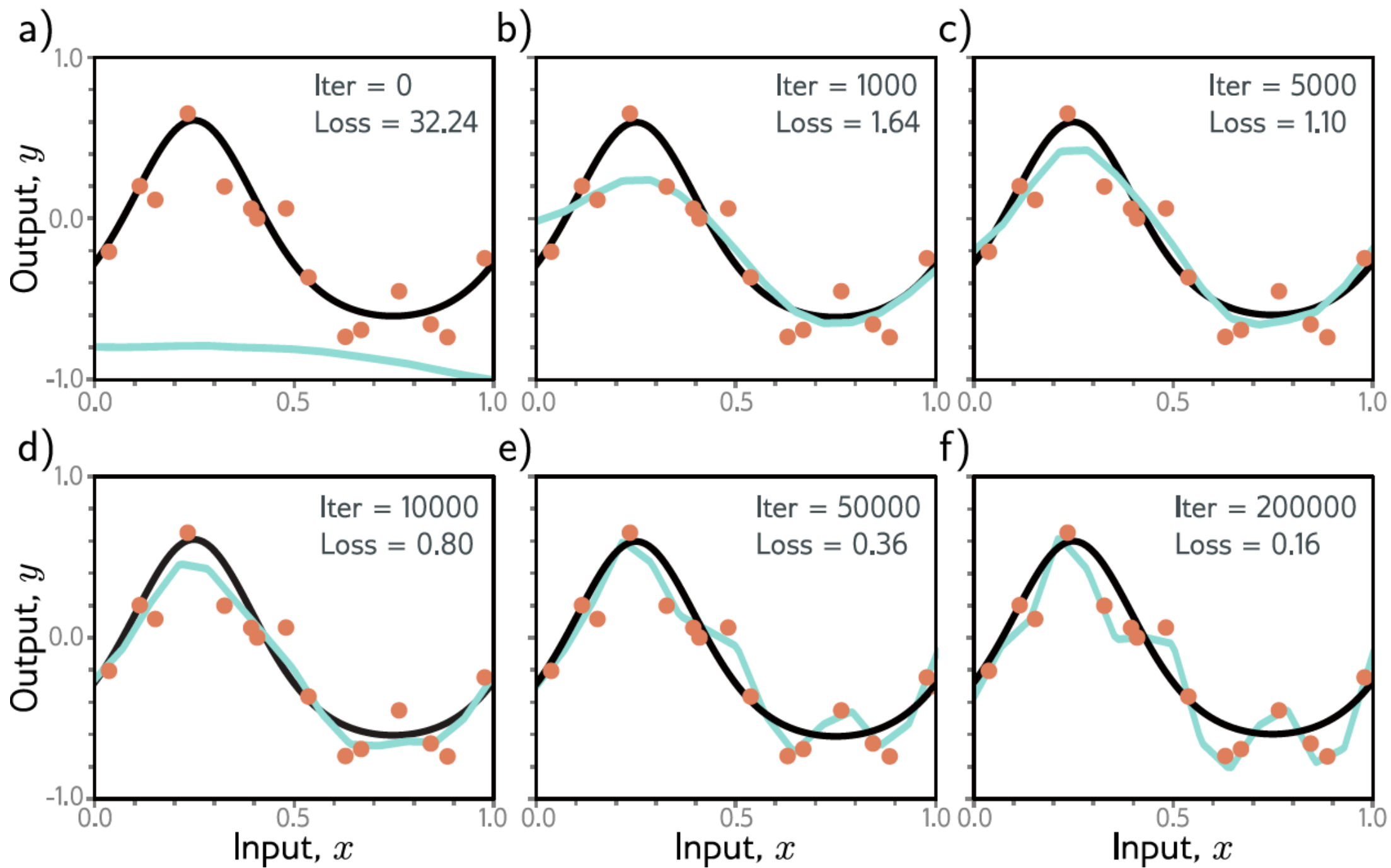
# Regularization

- Explicit regularization
- Implicit regularization
- **Early stopping**
- Ensembling
- Dropout
- Adding noise
- Bayesian approaches
- Transfer learning, multi-task learning, self-supervised learning
- Data augmentation

# Early stopping

- If we stop training early, weights don't have time to overfit to noise
- Weights start small, don't have time to get large
- Reduces effective model complexity
- Known as **early stopping**
- Don't have to re-train



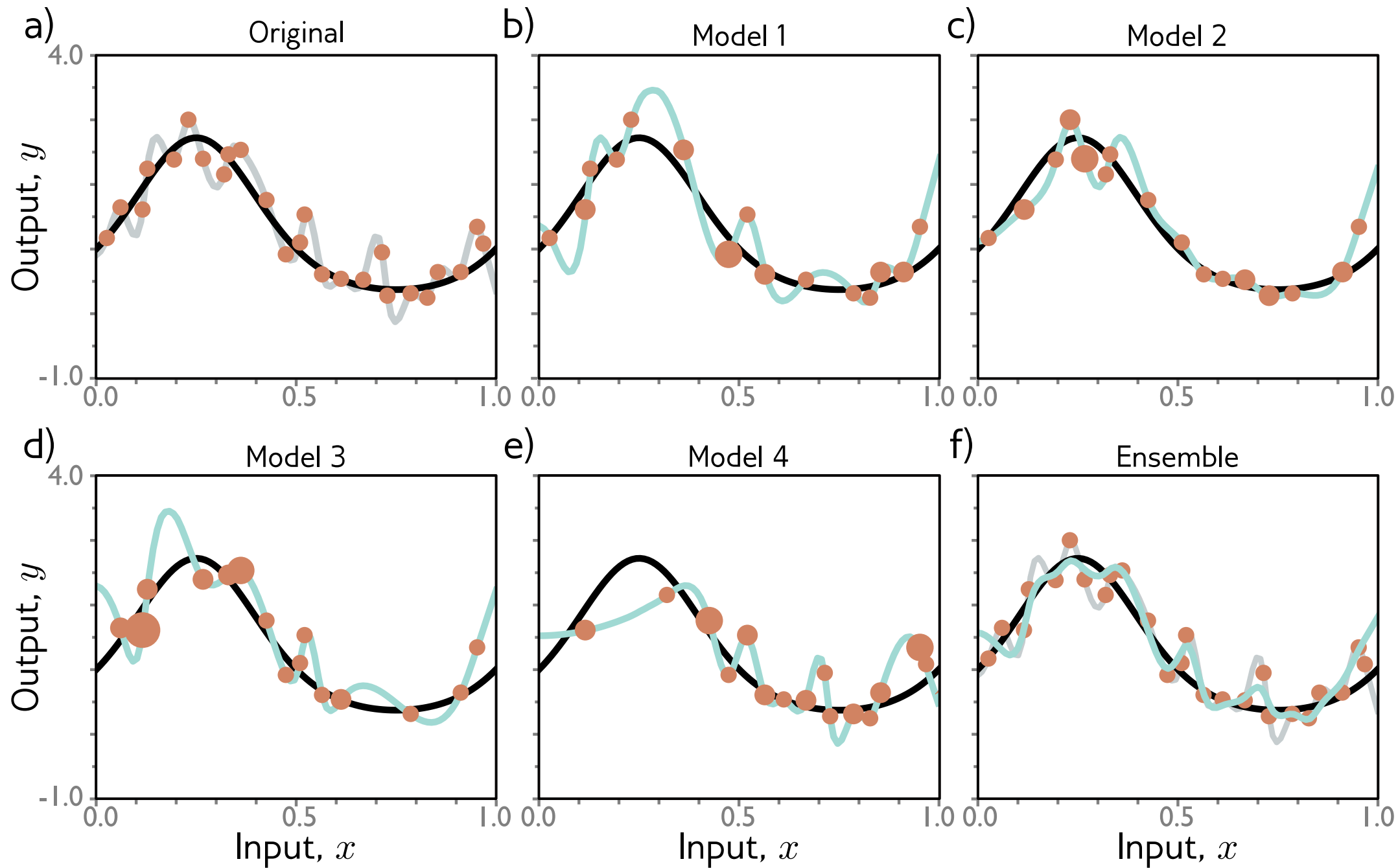


# Regularization

- Explicit regularization
- Implicit regularization
- Early stopping
- **Ensembling**
- Dropout
- Adding noise
- Bayesian approaches
- Transfer learning, multi-task learning, self-supervised learning
- Data augmentation

# Ensembling

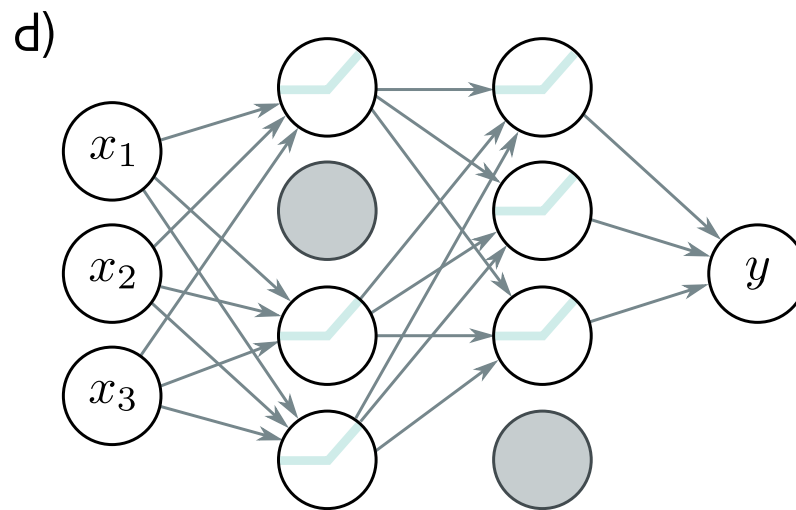
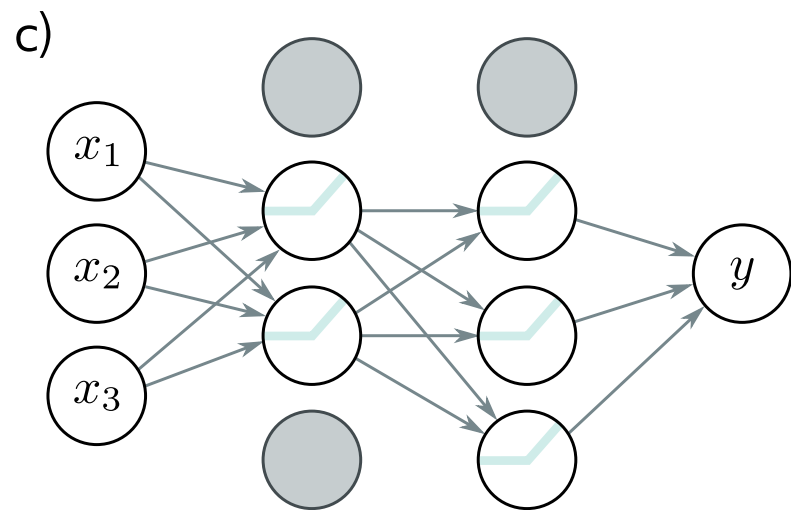
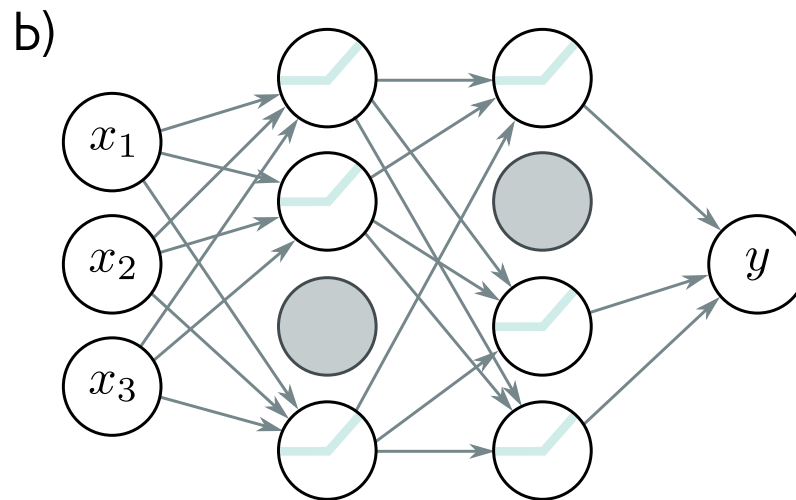
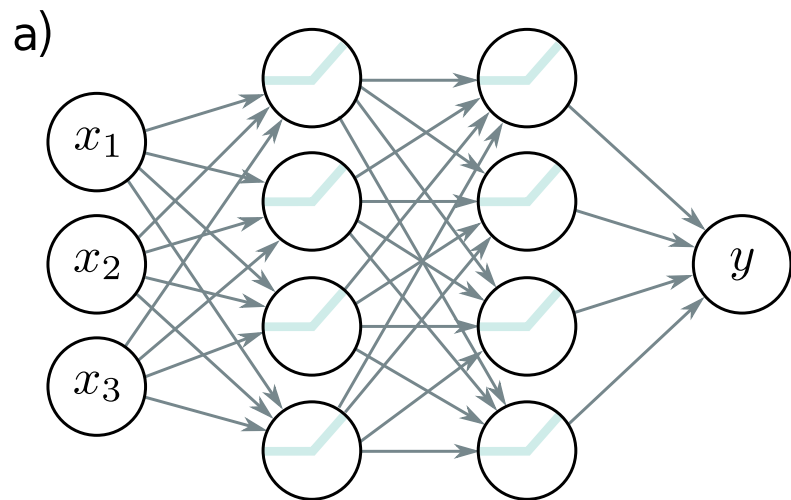
- Average together several models – an **ensemble**
- Can take mean or median
- Different initializations / different models
- Different subsets of the data resampled with replacements -- **bagging**



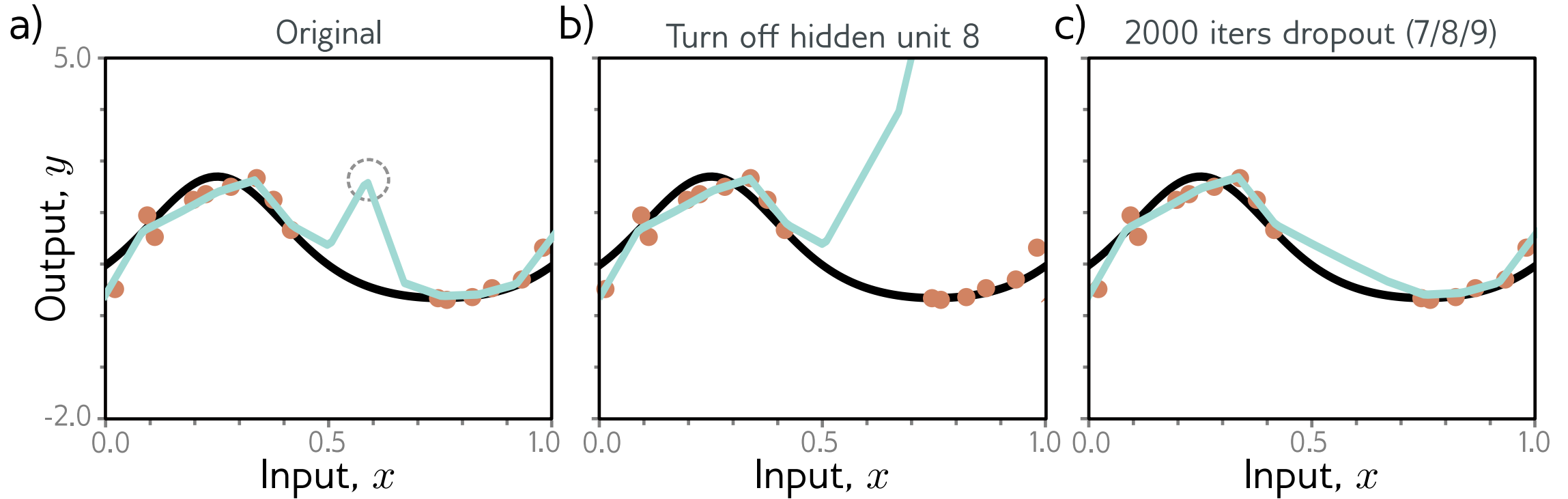
# Regularization

- Explicit regularization
- Implicit regularization
- Early stopping
- Ensembling
- Dropout
- Adding noise
- Bayesian approaches
- Transfer learning, multi-task learning, self-supervised learning
- Data augmentation

# Dropout



# Dropout



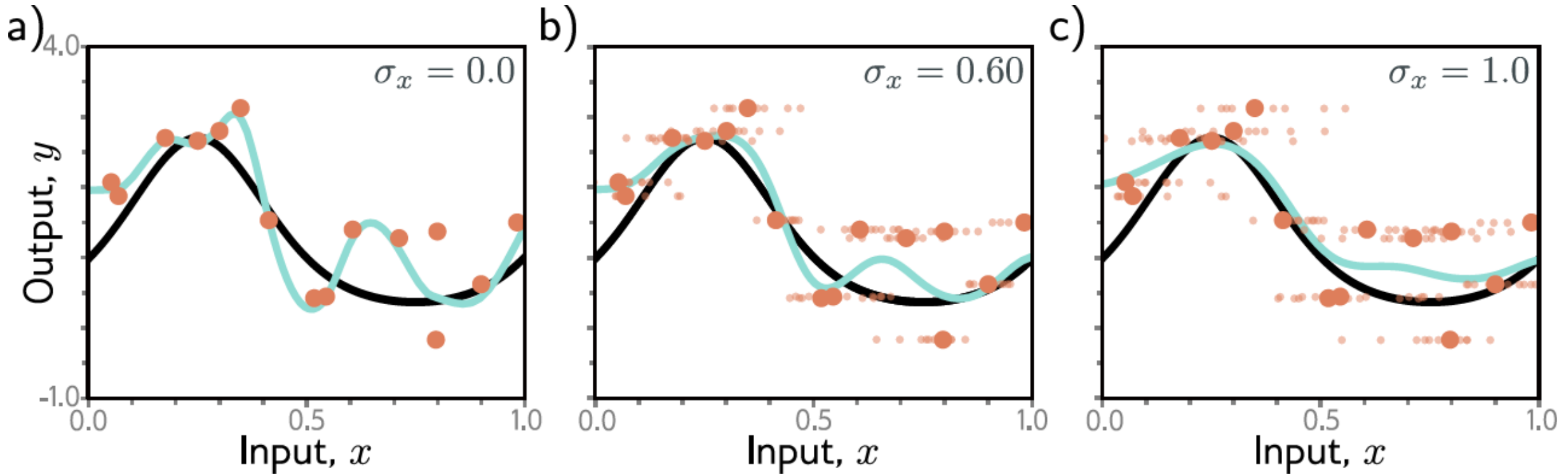
Can eliminate kinks in function that are far from data and don't contribute to training loss

# Regularization

- Explicit regularization
- Implicit regularization
- Early stopping
- Ensembling
- Dropout
- Adding noise
- Bayesian approaches
- Transfer learning, multi-task learning, self-supervised learning
- Data augmentation



# Adding noise



- to inputs
- to weights
- to outputs (labels)


# Regularization

- Explicit regularization
- Implicit regularization
- Early stopping
- Ensembling
- Dropout
- Adding noise
- Bayesian approaches
- Transfer learning, multi-task learning, self-supervised learning
- Data augmentation

# Bayesian approaches

- There are many parameters compatible with the data
- Can find a probability distribution over them

Prior info about parameters

$$Pr(\phi|\{\mathbf{x}_i, \mathbf{y}_i\}) = \frac{\prod_{i=1}^I Pr(\mathbf{y}_i|\mathbf{x}_i, \phi) Pr(\phi)}{\int \prod_{i=1}^I Pr(\mathbf{y}_i|\mathbf{x}_i, \phi) Pr(\phi) d\phi}$$


# Bayesian approaches

- There are many parameters compatible with the data
- Can find a probability distribution over them

Prior info about parameters

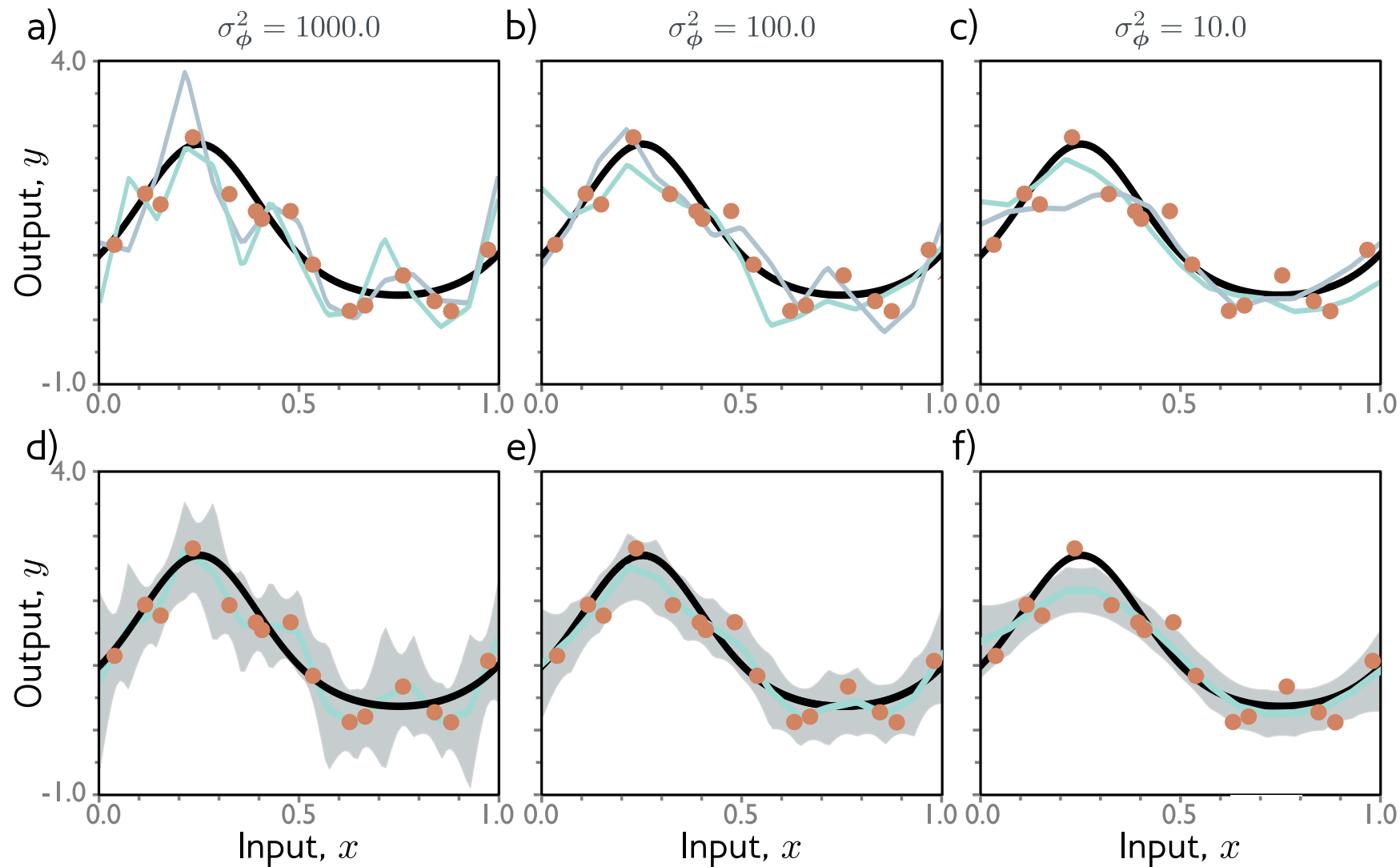


$$Pr(\phi|\{\mathbf{x}_i, \mathbf{y}_i\}) = \frac{\prod_{i=1}^I Pr(\mathbf{y}_i|\mathbf{x}_i, \phi) Pr(\phi)}{\int \prod_{i=1}^I Pr(\mathbf{y}_i|\mathbf{x}_i, \phi) Pr(\phi) d\phi}$$

- Take all possible parameters into account when make prediction

$$Pr(\mathbf{y}|\mathbf{x}, \{\mathbf{x}_i, \mathbf{y}_i\}) = \int Pr(\mathbf{y}|\mathbf{x}, \phi) Pr(\phi|\{\mathbf{x}_i, \mathbf{y}_i\}) d\phi$$

# Bayesian approaches

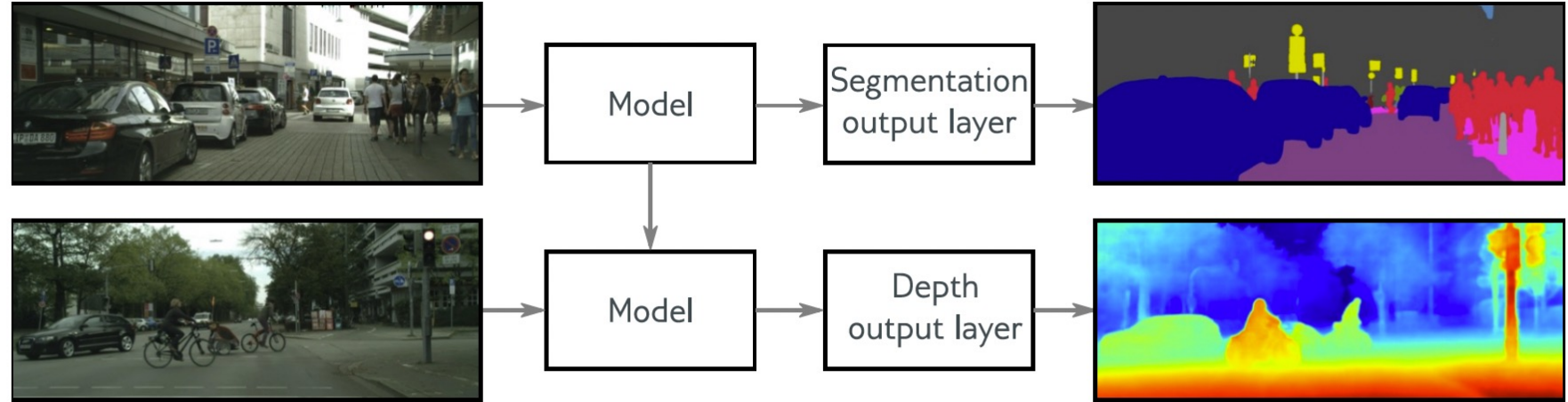


# Regularization

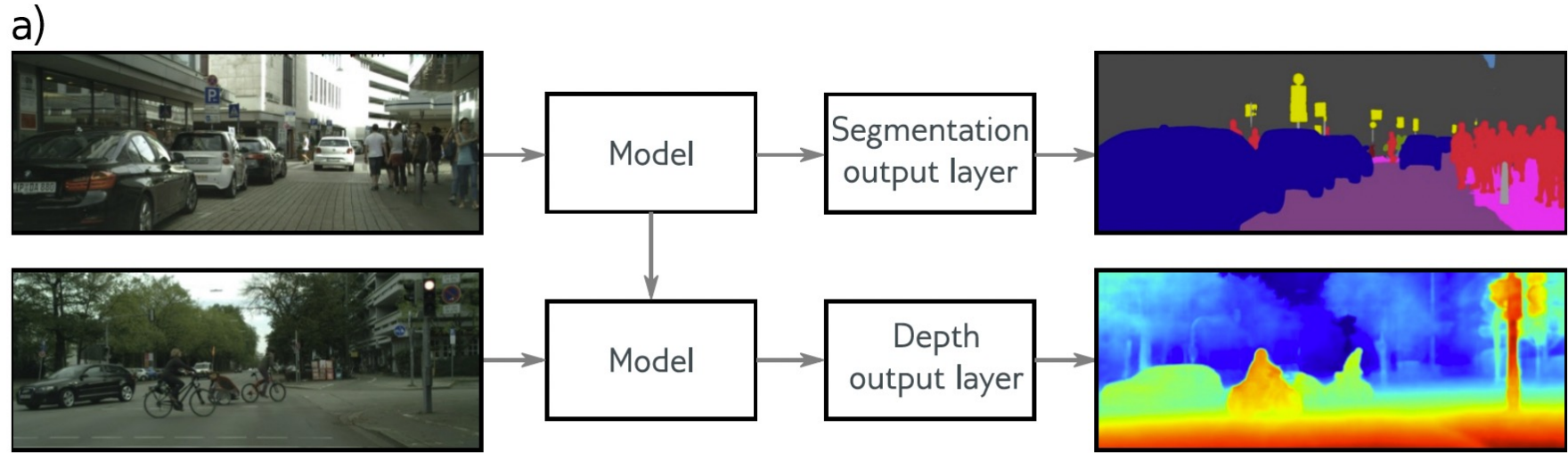
- Explicit regularization
- Implicit regularization
- Early stopping
- Ensembling
- Dropout
- Adding noise
- Bayesian approaches
- Transfer learning, multi-task learning, self-supervised learning
- Data augmentation

- Transfer learning

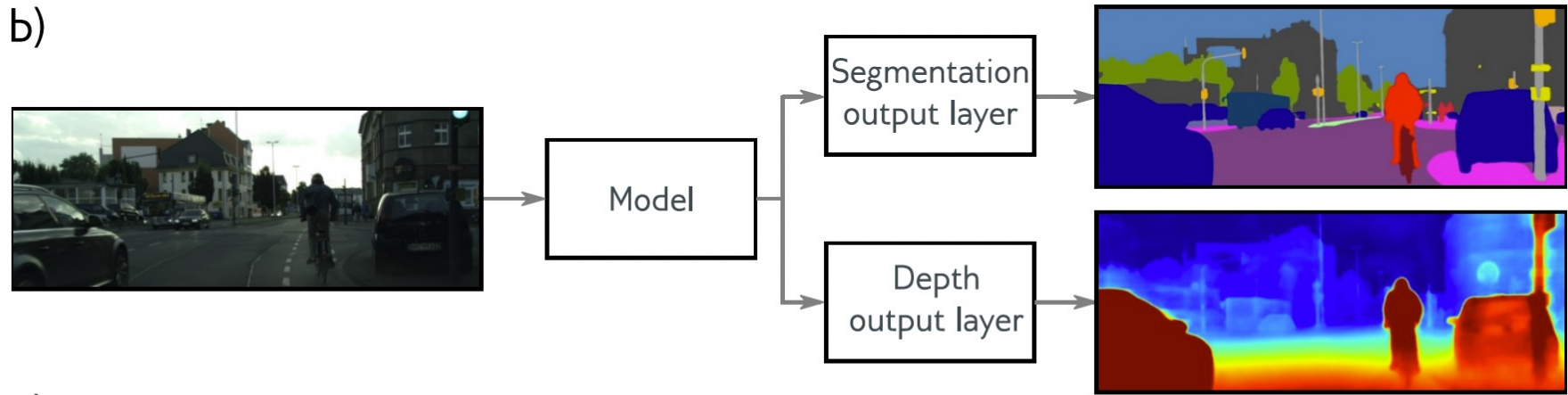
a)



- Transfer learning

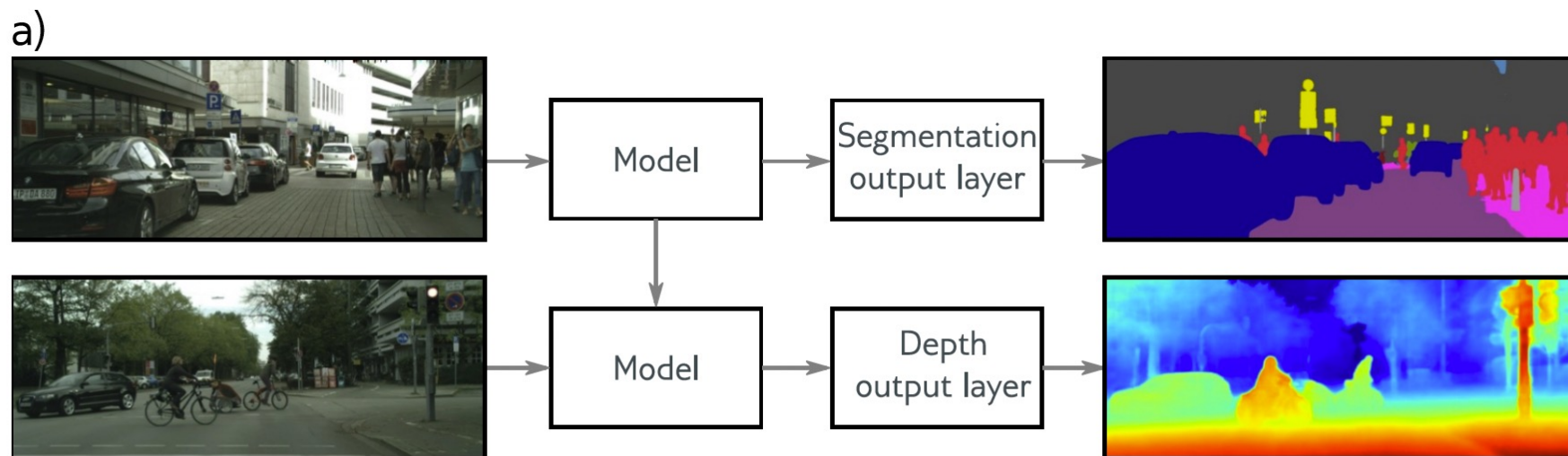


- Multi-task learning

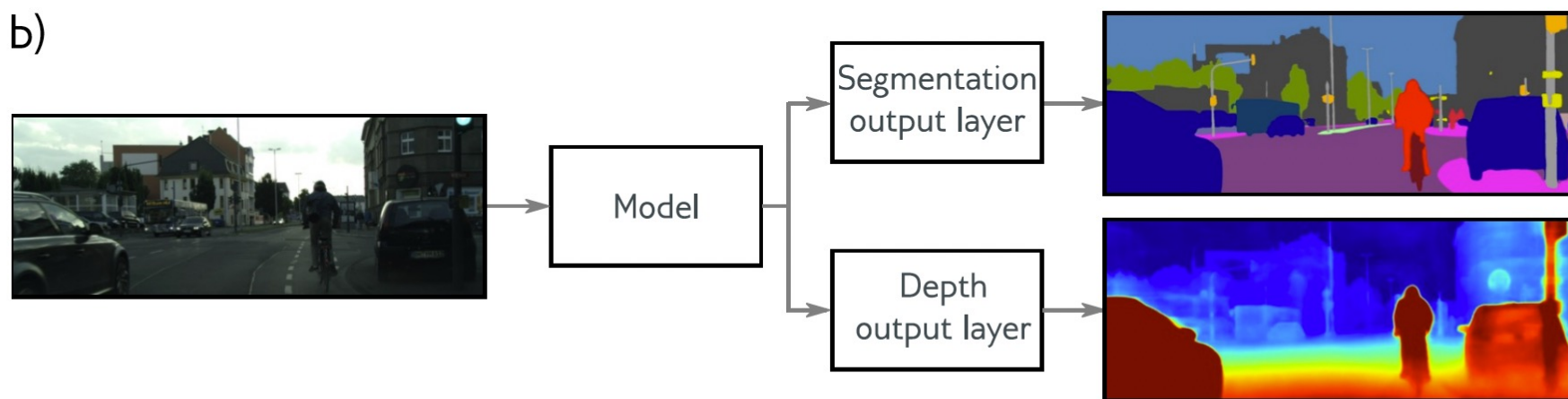




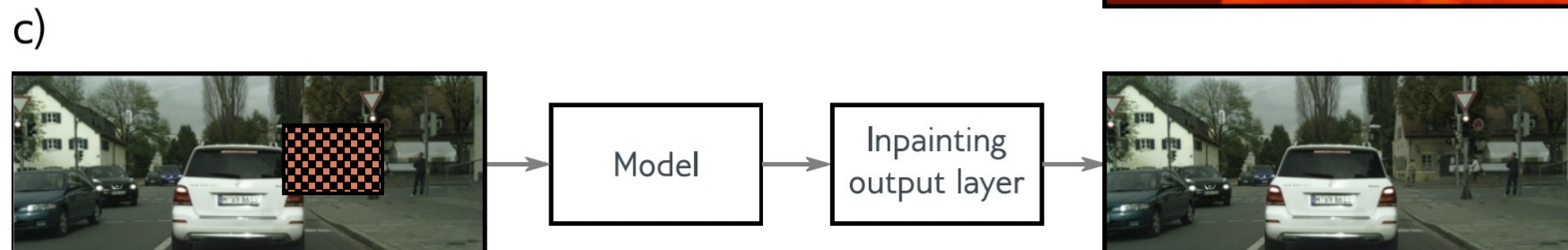
- Transfer learning



- Multi-task learning



- Self-supervised learning



# Regularization

- Explicit regularization
- Implicit regularization
- Early stopping
- Ensembling
- Dropout
- Adding noise
- Bayesian approaches
- Transfer learning, multi-task learning, self-supervised learning
- Data augmentation

# Data augmentation

a) Original



b) Flip



c) Rotate and crop



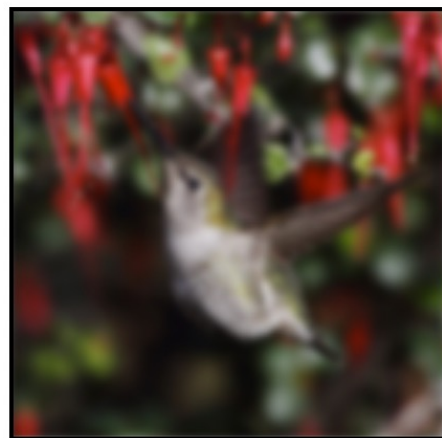
d) Vertical stretch



e) Color balance



f) Blur



g) Vignette



h) Pincushion



# Regularization overview

