

Metrics & regularization

Covered chapters: Prince Ch 8-9

Miles Cranmer

Expectations

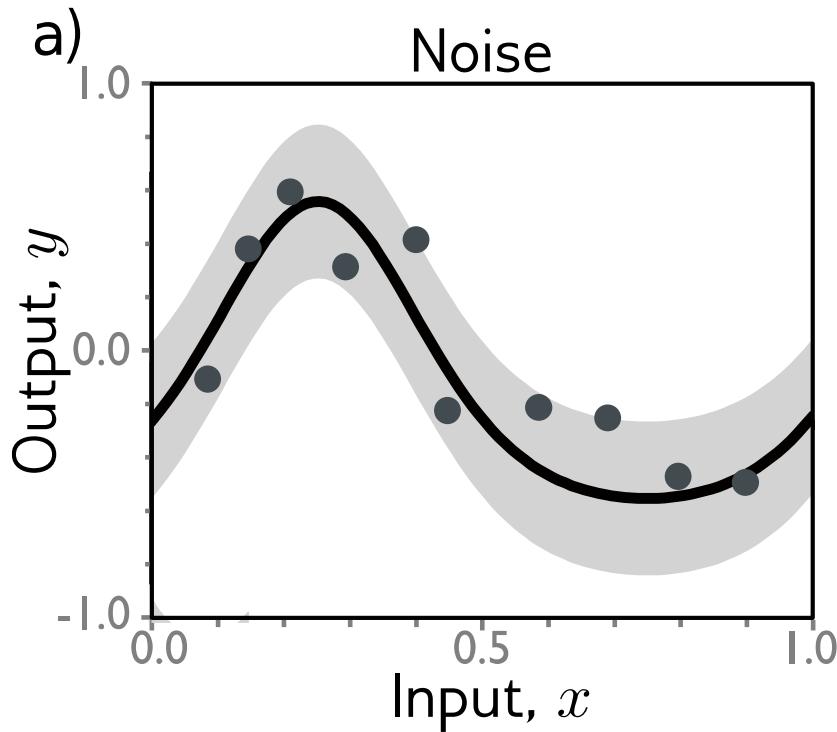
$$\mathbb{E}[g[x]] = \int g[x] Pr(x) dx,$$

Interpretation: what is the average value of $g[x]$ when taking into account the probability of x ?

Could approximate, by sampling many values of x from the distribution, calculating $g[x]$, and taking average:

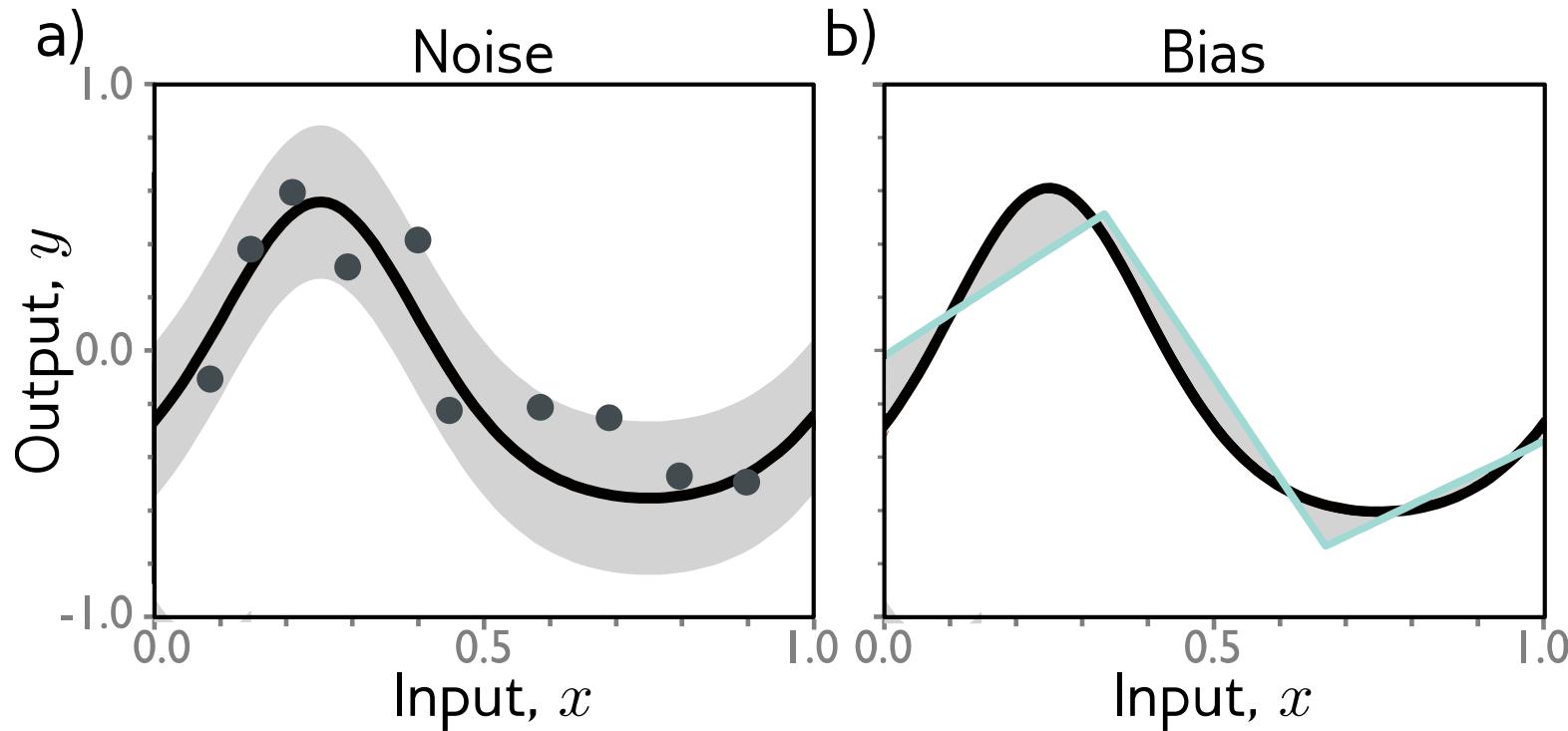
$$\mathbb{E}[g[x]] \approx \frac{1}{N} \sum_{n=1}^N g[x_n^*] \quad \text{where} \quad x_n^* \sim Pr(x)$$

Noise, bias, and variance

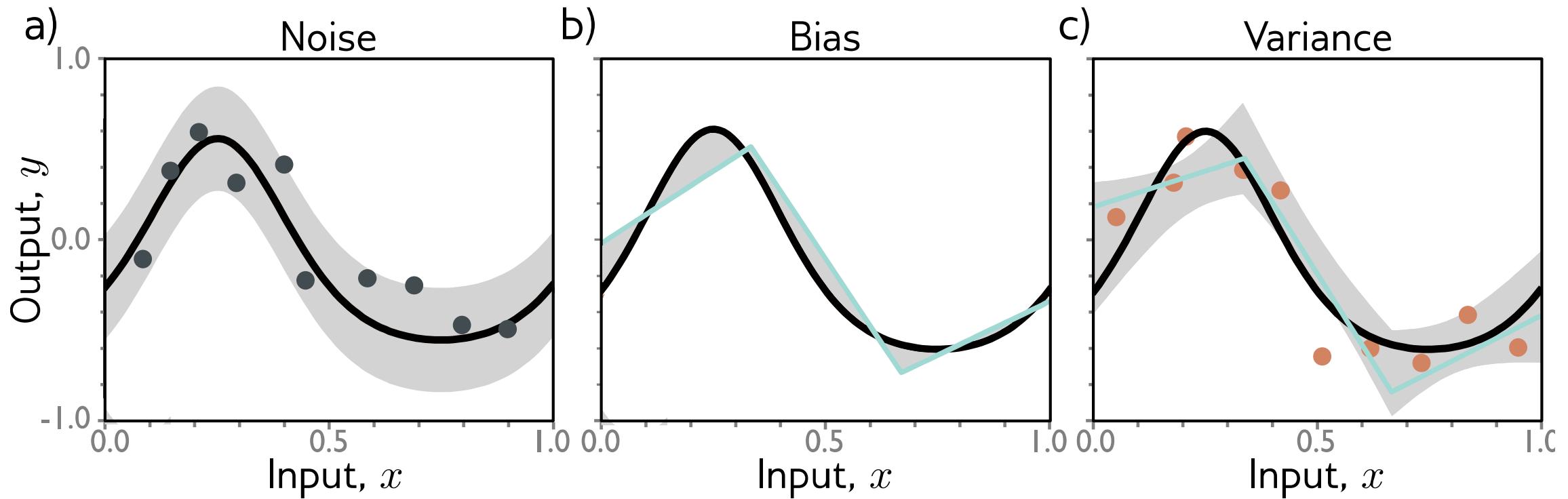


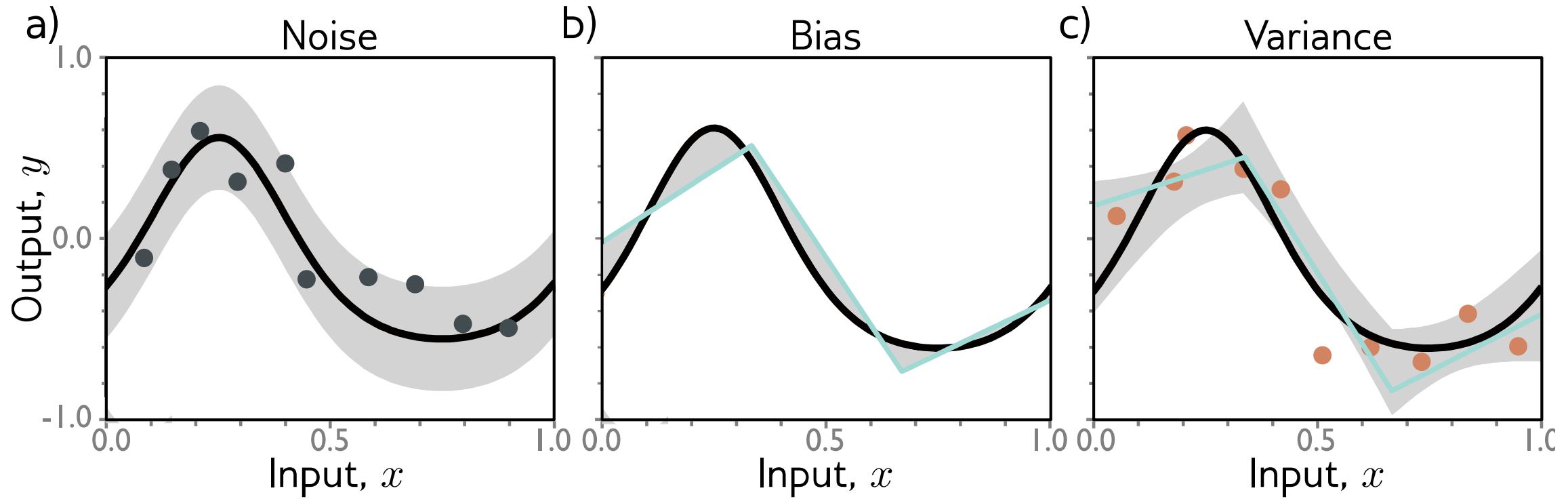
- Noise in measurements
- Some variables not observed
- Data mislabeled

Noise, **bias**, and variance



Noise, bias, and variance





$$\begin{aligned}
 \text{(Total error)} = & \underbrace{\mathbb{E}_{\mathcal{D}} \left[(f[x, \phi[\mathcal{D}]] - f_{\mu}[x])^2 \right]}_{\text{variance}} + \underbrace{\left(f_{\mu}[x] - \mu[x] \right)^2}_{\text{bias}} + \underbrace{\sigma^2}_{\text{noise}}
 \end{aligned}$$

Noise, bias, and variance

Noise, bias, and variance

- Variance is the uncertainty in fitted model due to choice of **training set**

Noise, bias, and variance

- Variance is the uncertainty in fitted model due to choice of **training set**
- Bias is systematic deviation from the mean of the function we are modeling due to **limitations in our model**

Noise, bias, and variance

- Variance is the uncertainty in fitted model due to choice of **training set**
- Bias is systematic deviation from the mean of the function we are modeling due to **limitations in our model**
- Noise is **inherent uncertainty** in the true mapping from input to output

Least squares regression only

$$L[x] = (f[x, \phi] - y[x])^2$$

- We can show that:

$$\mathbb{E}_{\mathcal{D}} [\mathbb{E}_y [L[x]]] = \underbrace{\mathbb{E}_{\mathcal{D}} [(f[x, \phi[\mathcal{D}]] - f_{\mu}[x])^2]}_{\text{variance}} + \underbrace{(f_{\mu}[x] - \mu[x])^2}_{\text{bias}} + \underbrace{\sigma^2}_{\text{noise}}$$

Expectation over noise in training data

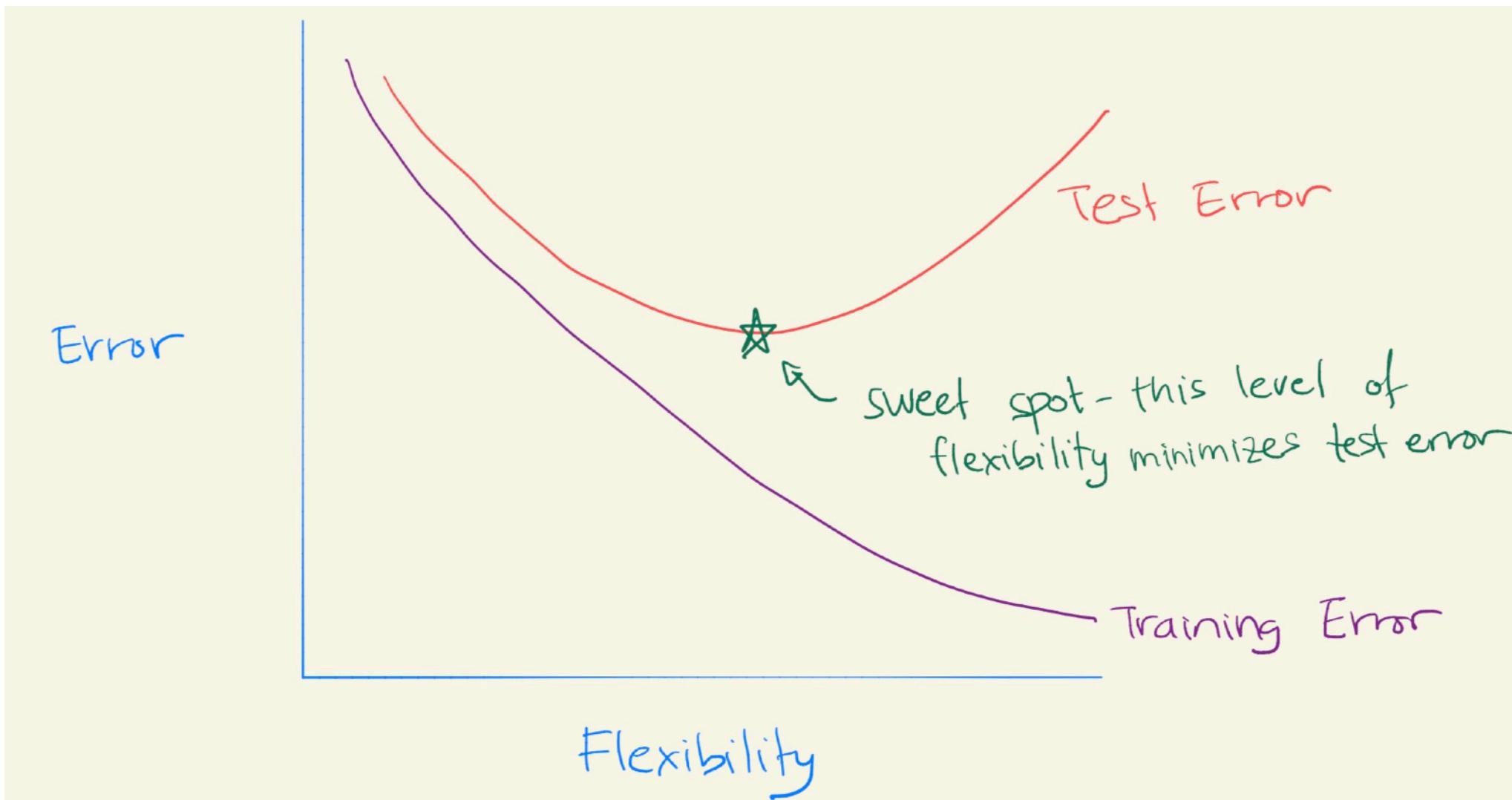
Expectation over noise in test data

Actual model

Best possible model if we had infinite data

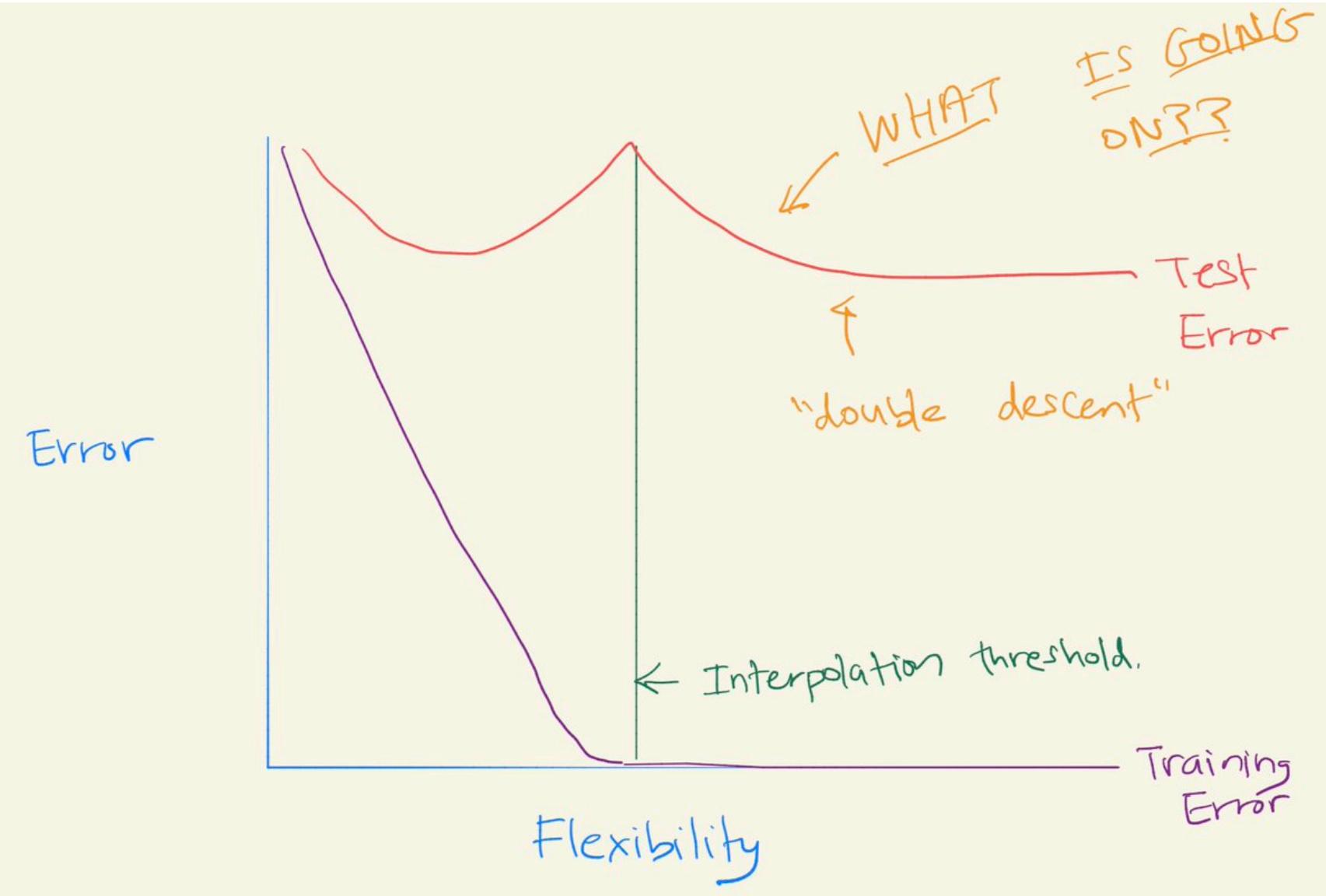
True function

Double descent



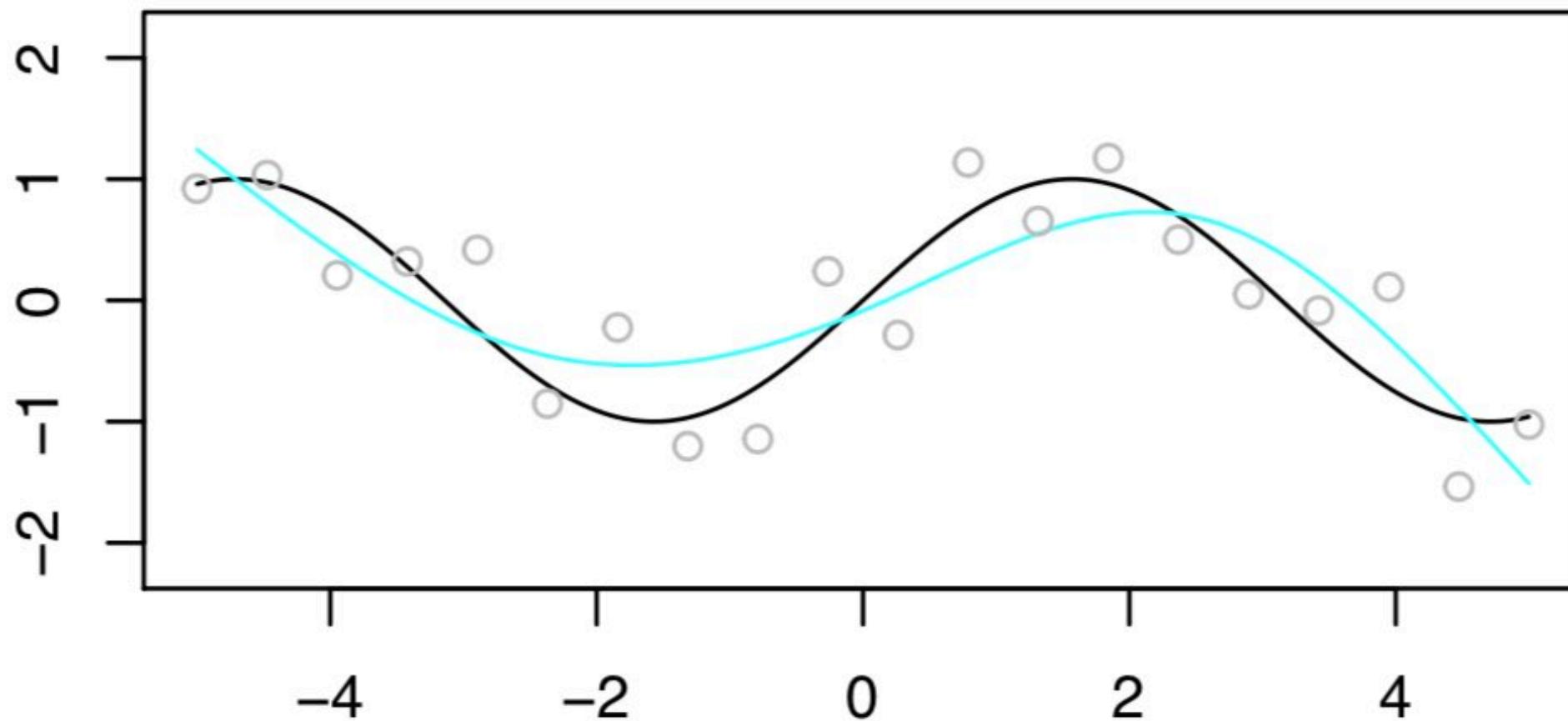
Credit: Daniela Witten
9

More detail – double descent



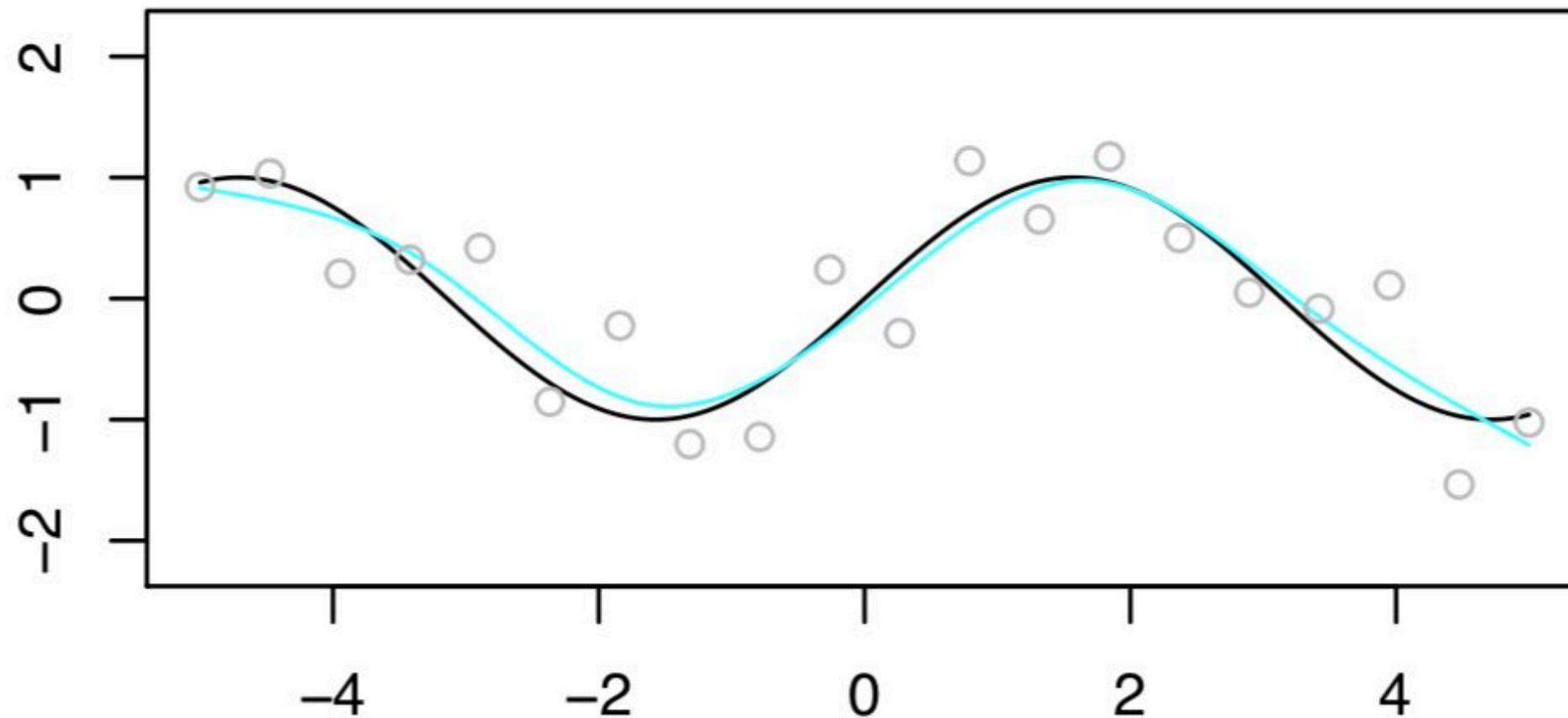
Credit: Daniela Witten

4 Degrees of Freedom



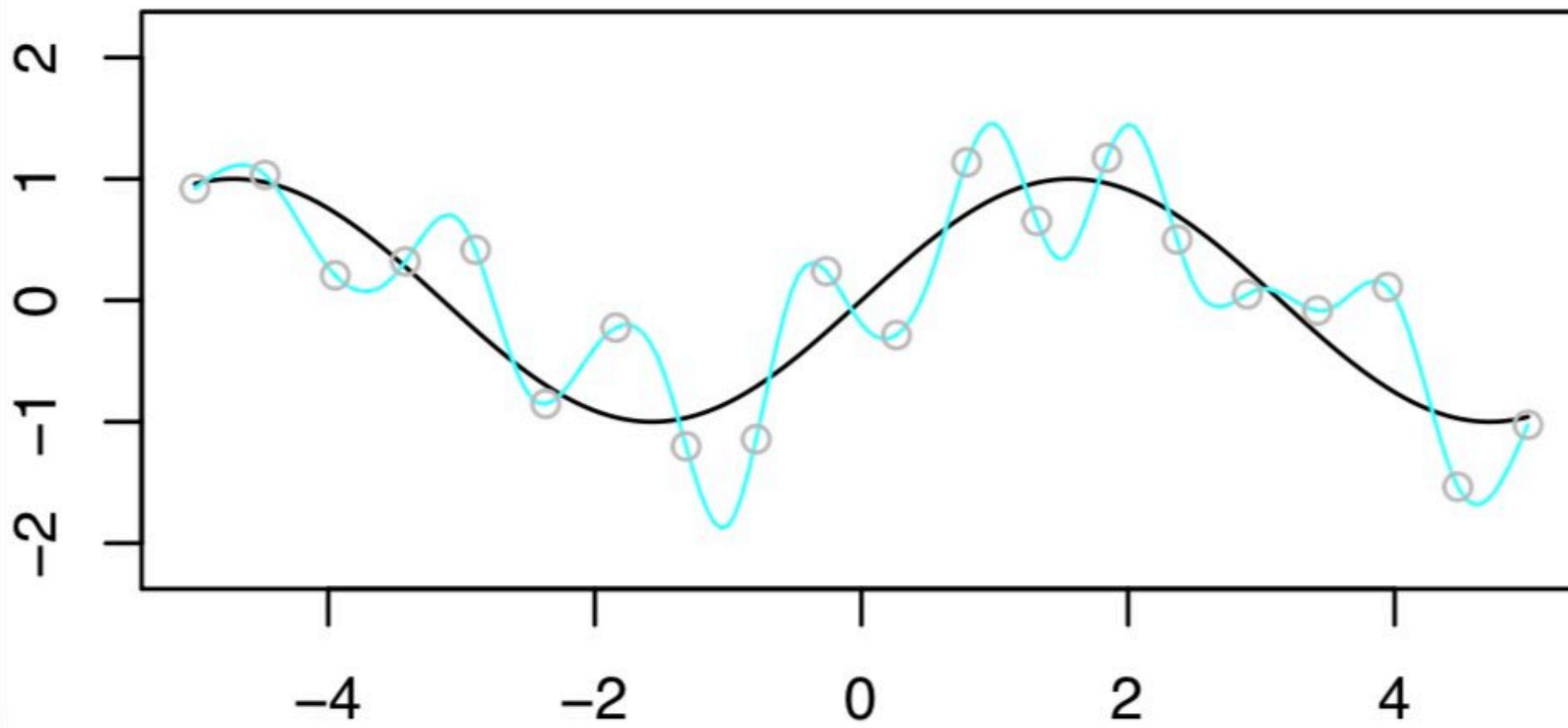
Credit: Daniela Witten
11

6 Degrees of Freedom



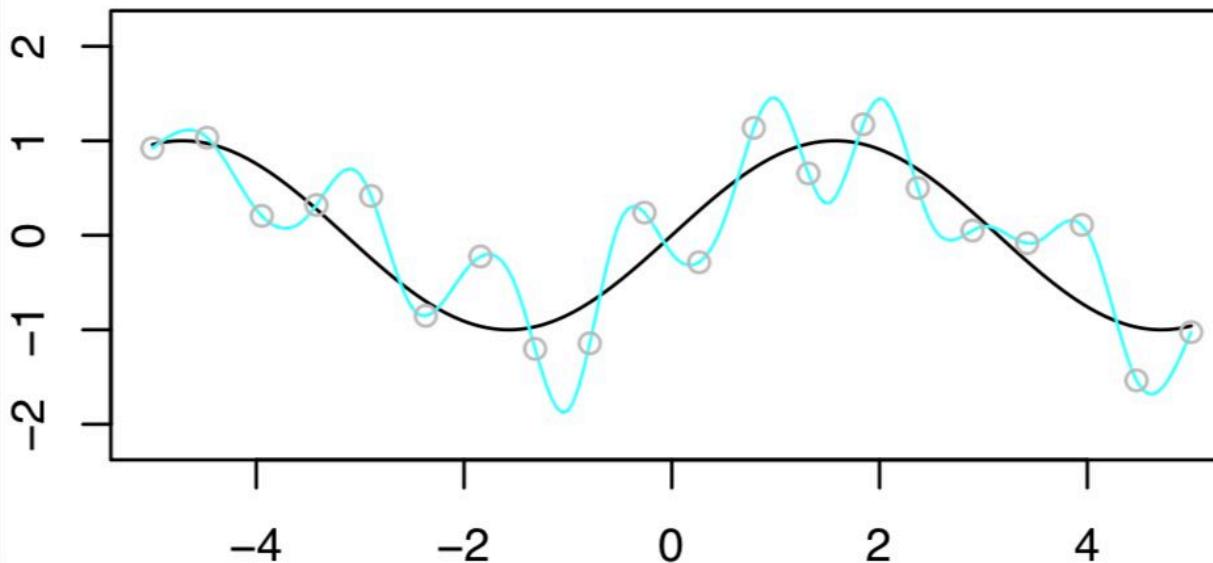
Credit: Daniela Witten
12

20 Degrees of Freedom

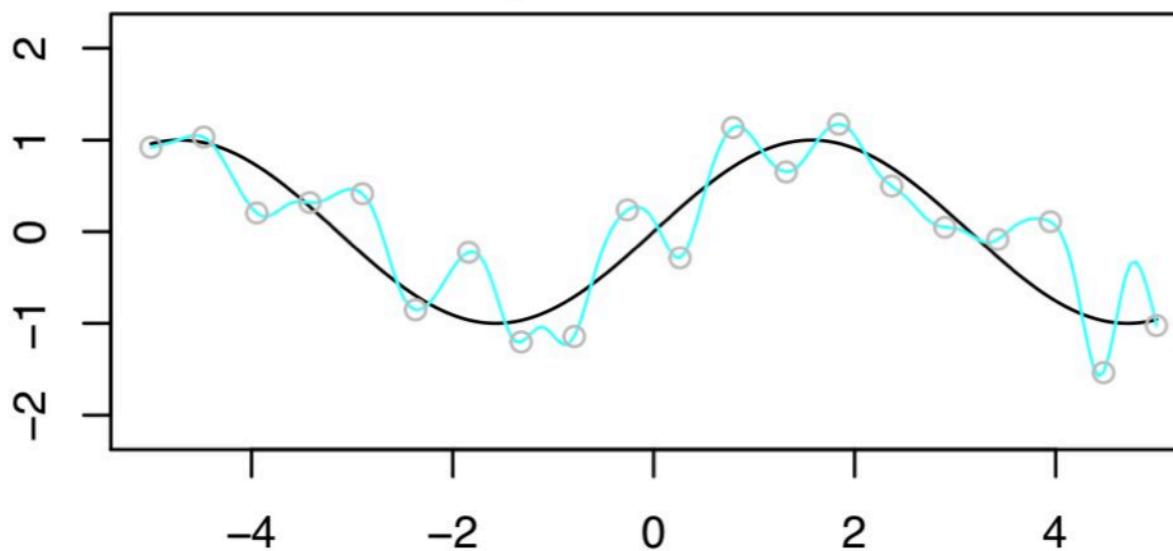


Credit: Daniela Witten
13

20 Degrees of Freedom



36 Degrees of Freedom



Better fit!

Important: requires that we regularize coefficients in polynomial!

(Much more extreme in Neural Networks!
Can actually beat the first descent.)

Choosing hyperparameters

- Don't know bias or variance
- Don't know how much capacity to add
- How do we choose capacity in practice?
 - Or model structure
 - Or training algorithm
 - Or learning rate
- Third data set – validation set
 - Train models with different hyperparameters on training set
 - Choose best hyperparameters with validation set
 - Test once with test set

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:

$$\hat{\phi} = \operatorname{argmin}_{\phi} [L[\phi]]$$

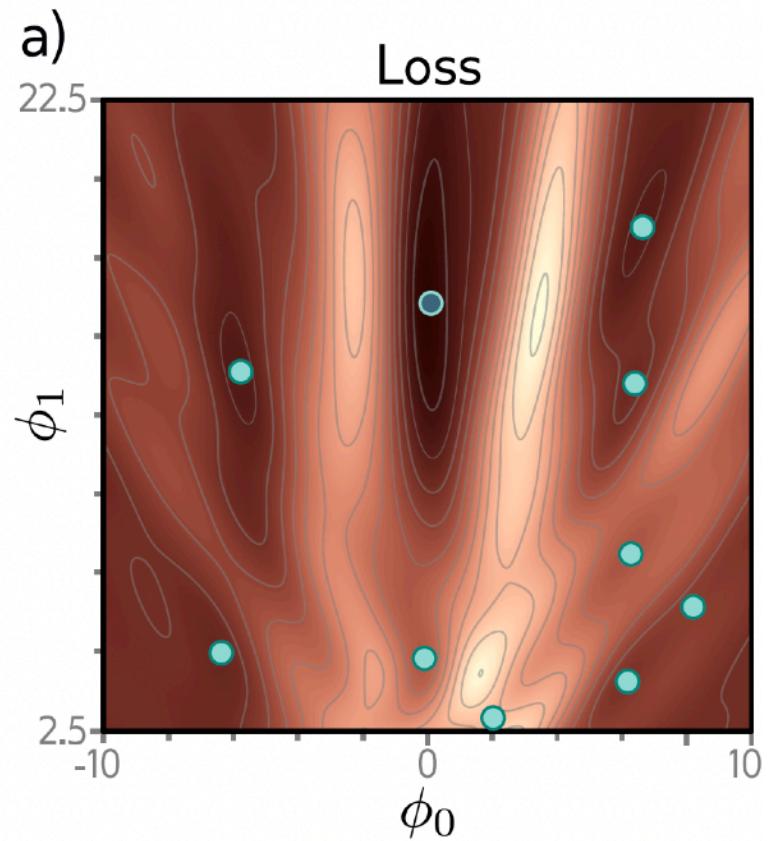
$$= \operatorname{argmin}_{\phi} \left[\sum_{i=1}^I \ell_i[\mathbf{x}_i, \mathbf{y}_i] \right]$$

- Regularization adds an extra term

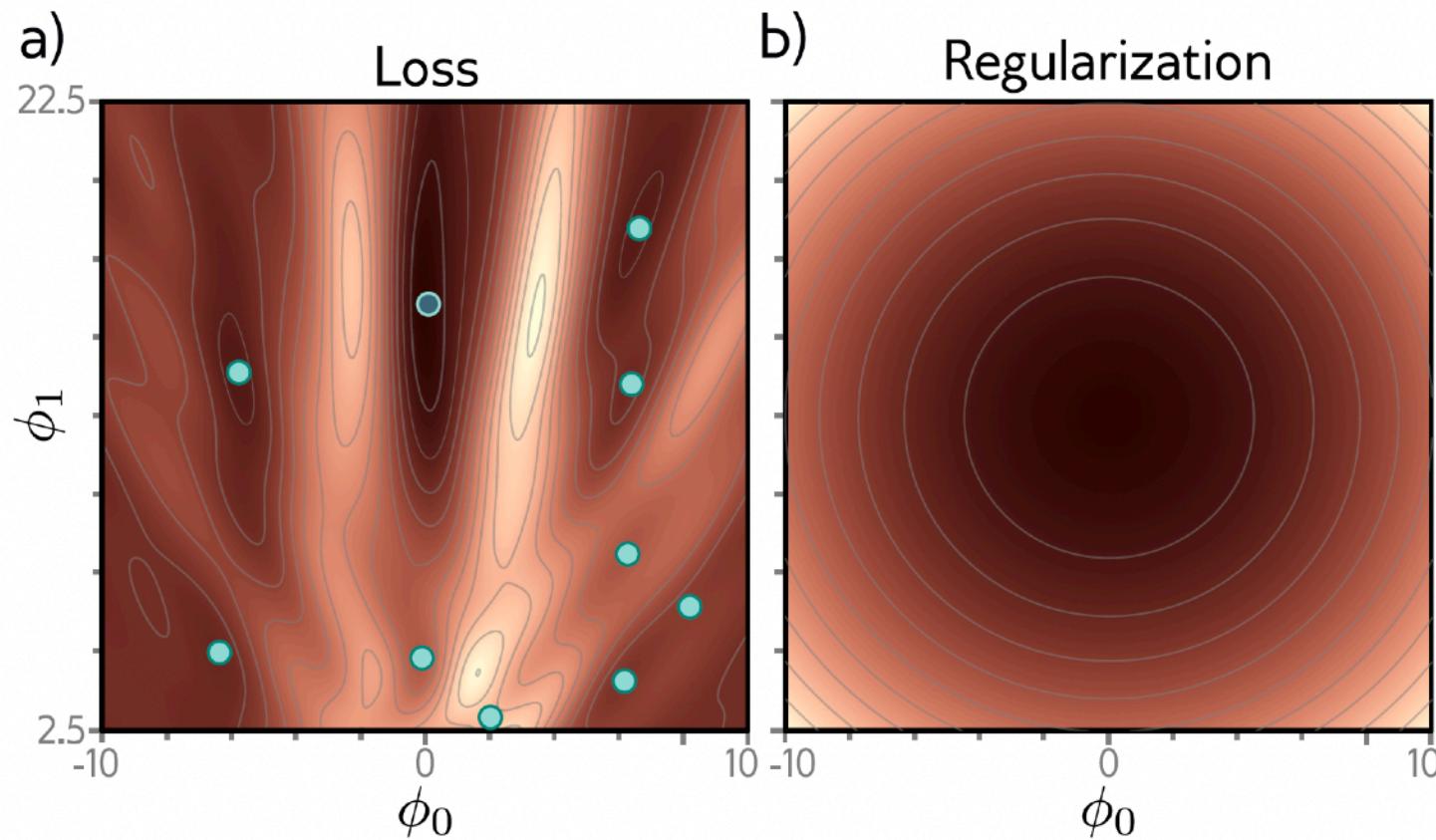
$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[\sum_{i=1}^I \ell_i[\mathbf{x}_i, \mathbf{y}_i] + \lambda \cdot g[\phi] \right]$$

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

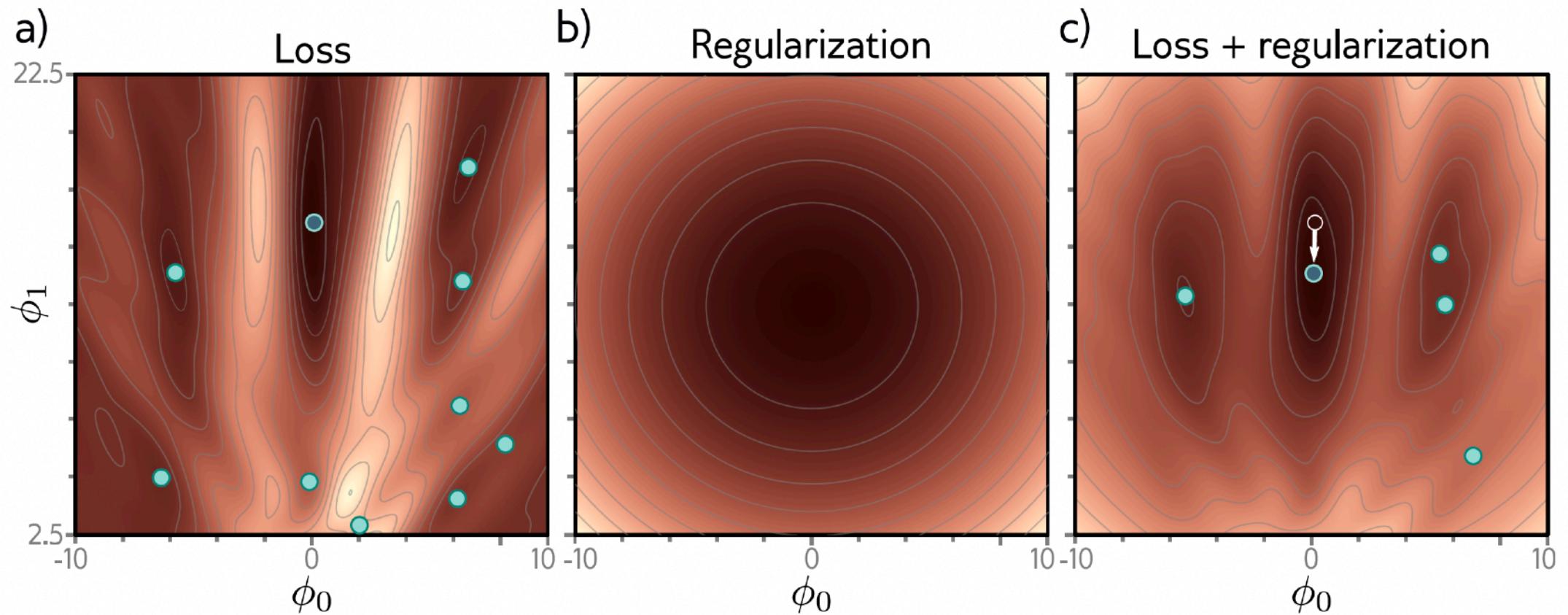
Explicit regularization



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:

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

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

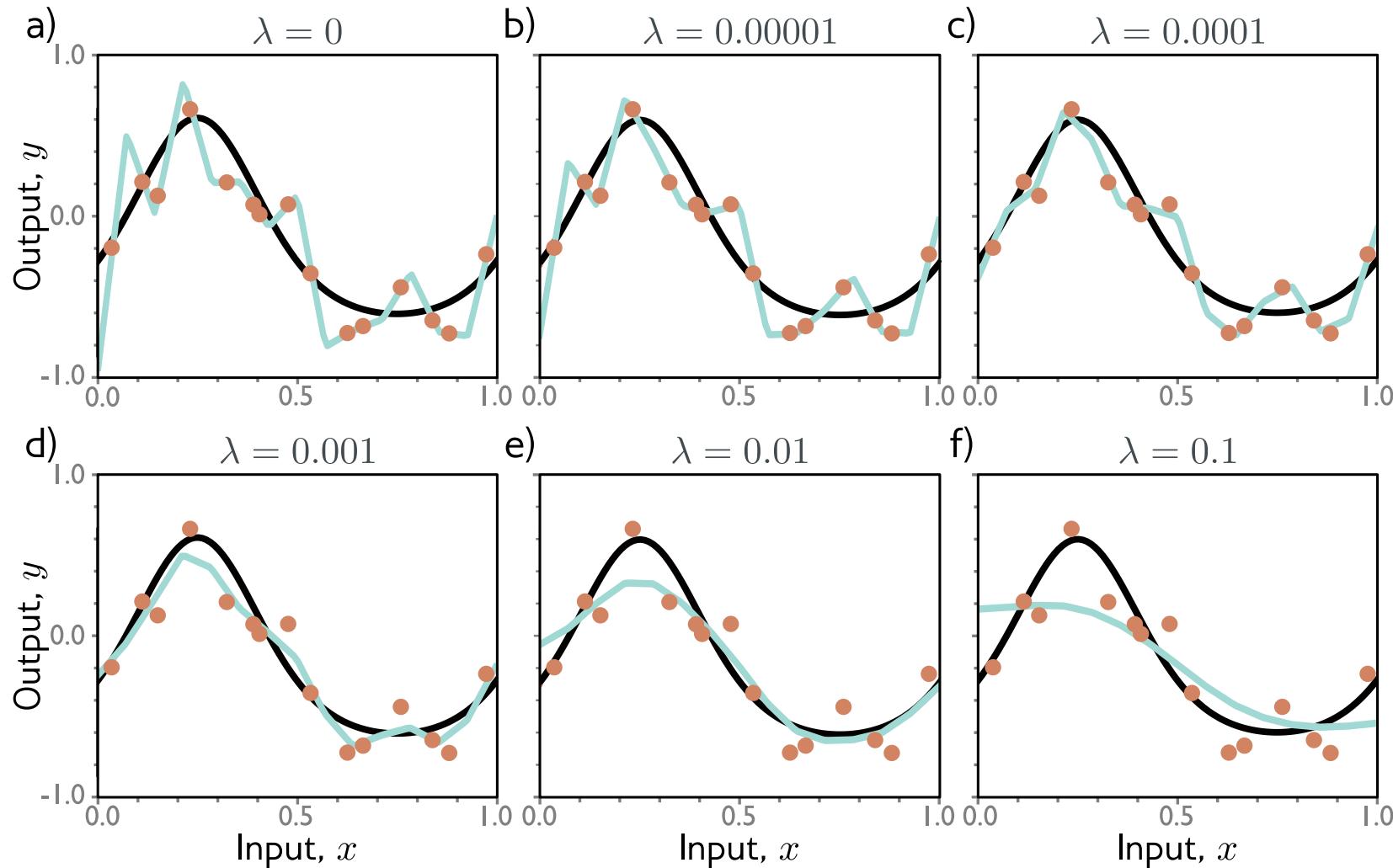
$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[L[\phi, \{\mathbf{x}_i, \mathbf{y}_i\}] + \lambda \sum_j \phi_j^2 \right]$$

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

Why does L2 regularization help?

- Encourages smoothness between datapoints

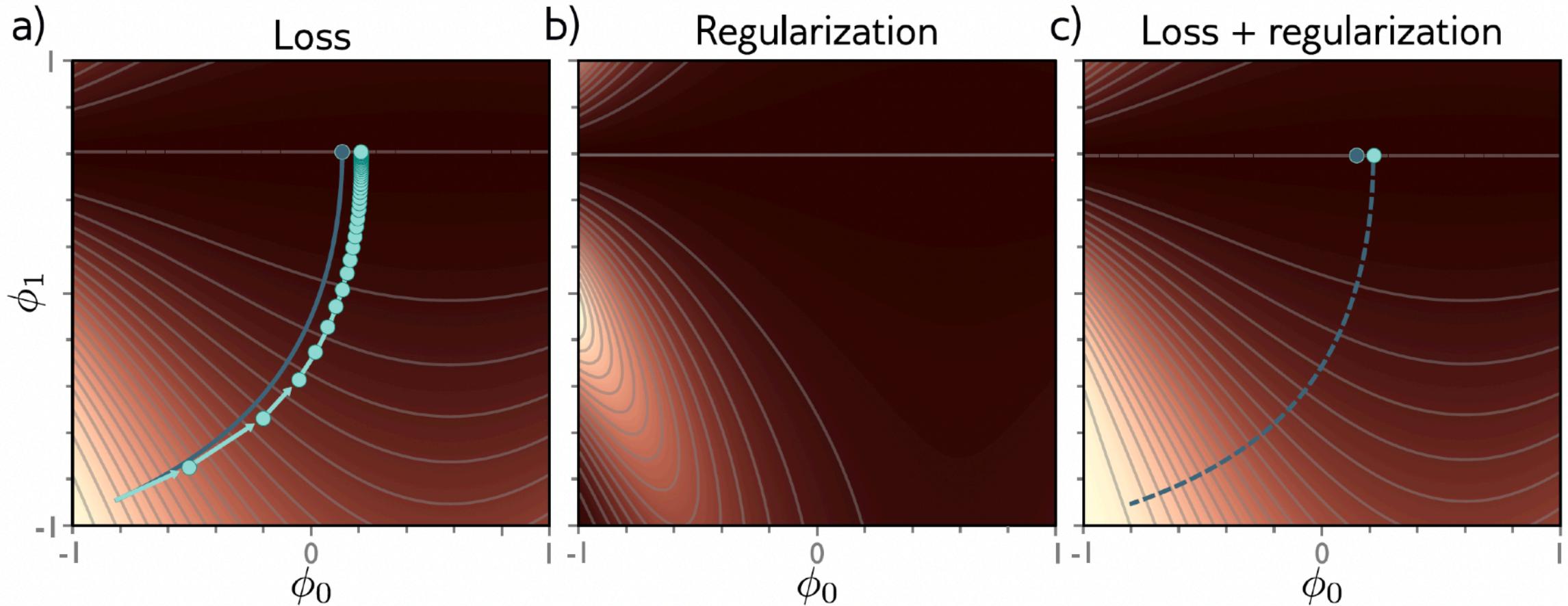
L2 regularization



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.

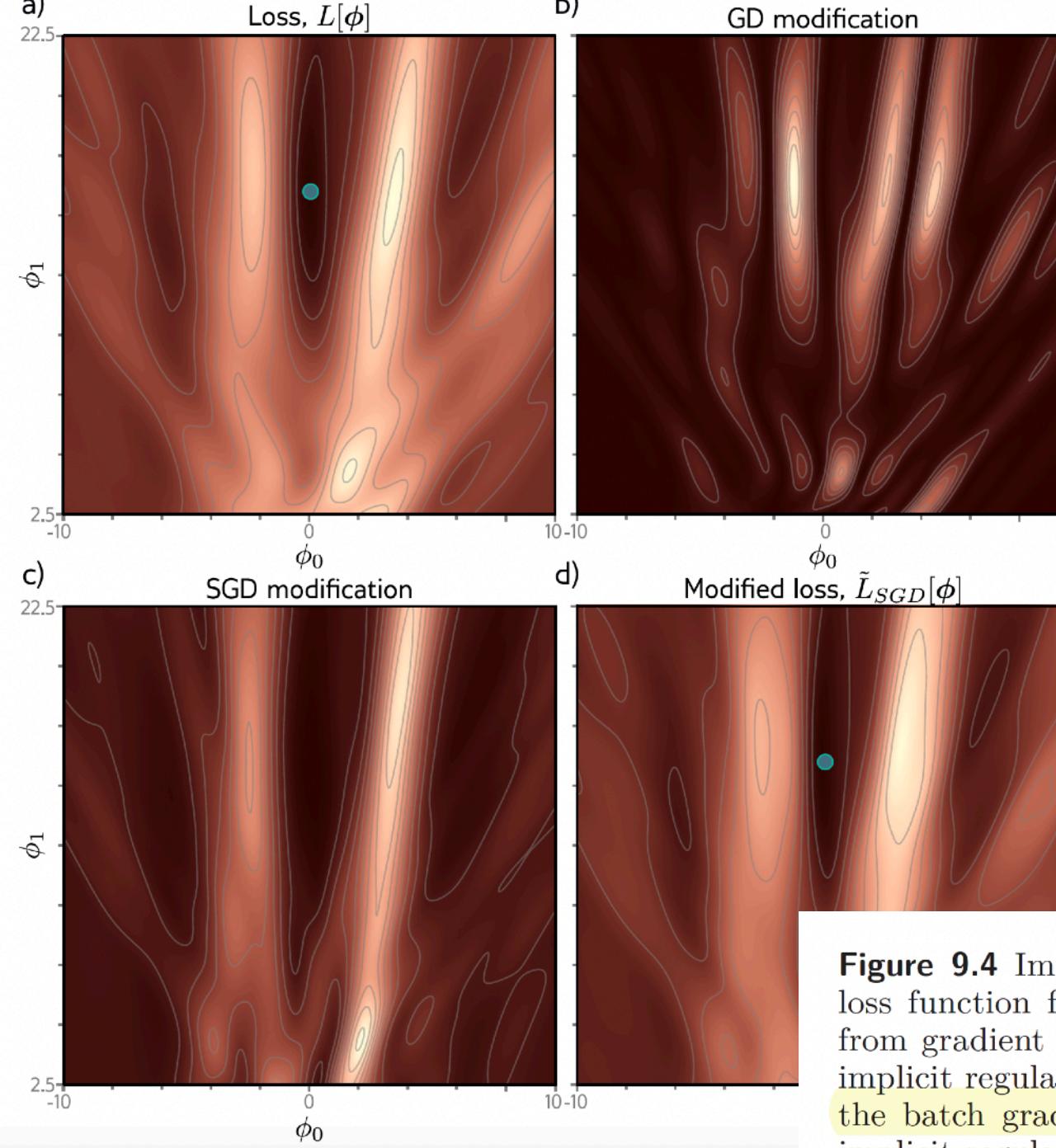
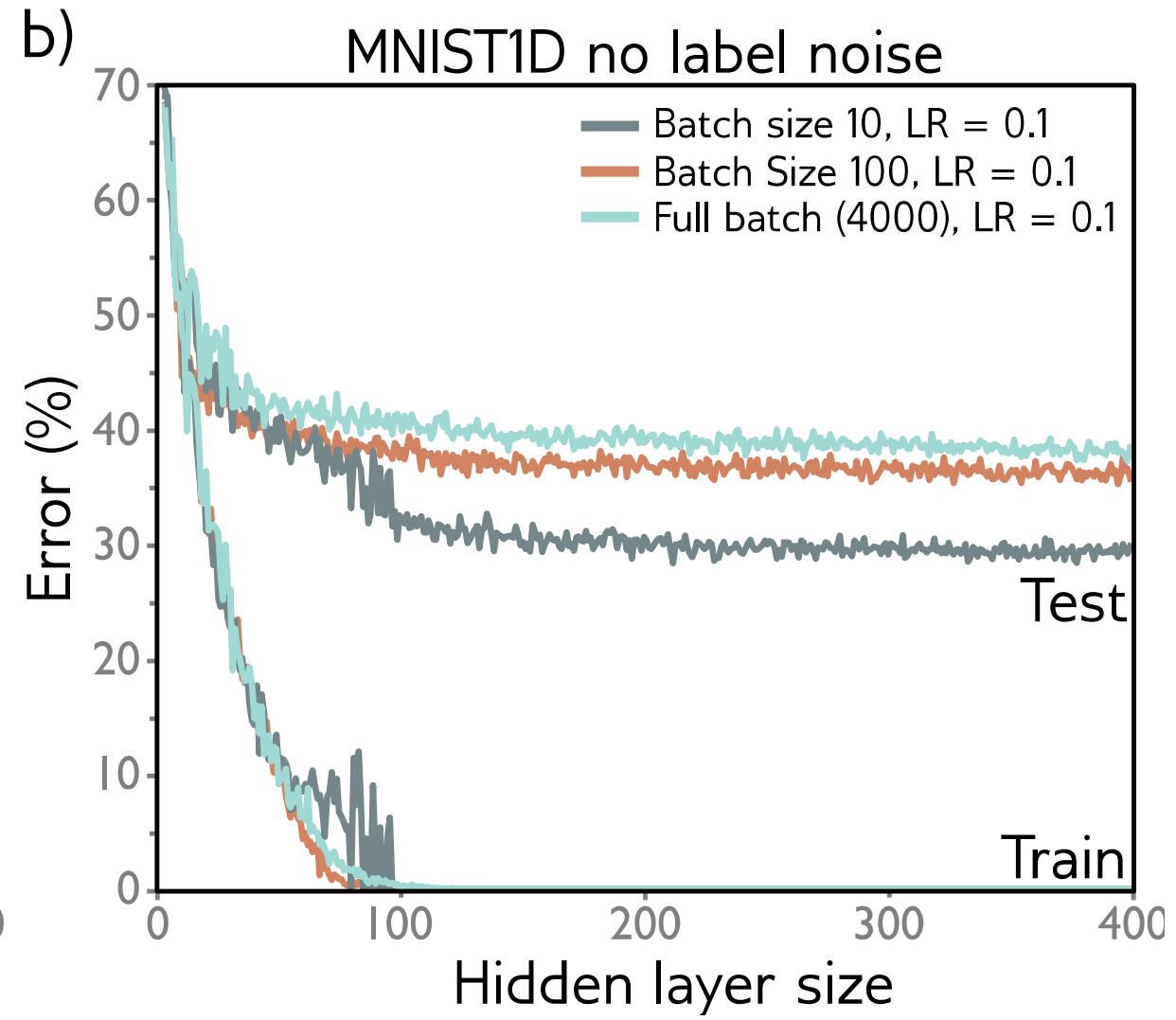
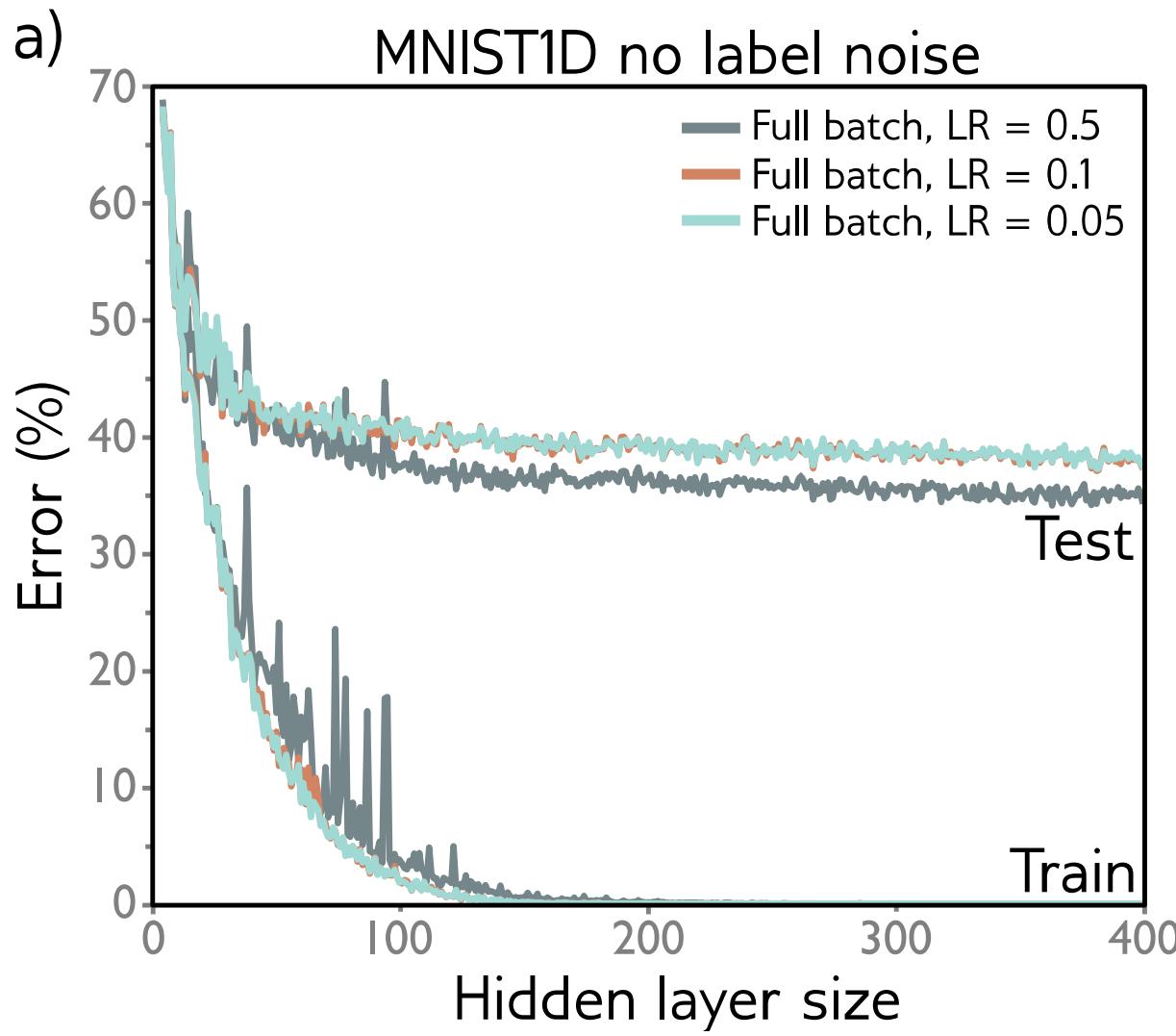


Figure 9.4 Implicit regularization for stochastic gradient descent. a) Original loss function for Gabor model (section 6.1.2). b) Implicit regularization term from gradient descent penalizes the squared gradient magnitude. c) Additional implicit regularization from stochastic gradient descent penalizes the variance of the batch gradients. d) Modified loss function (sum of original loss plus two implicit regularization components).

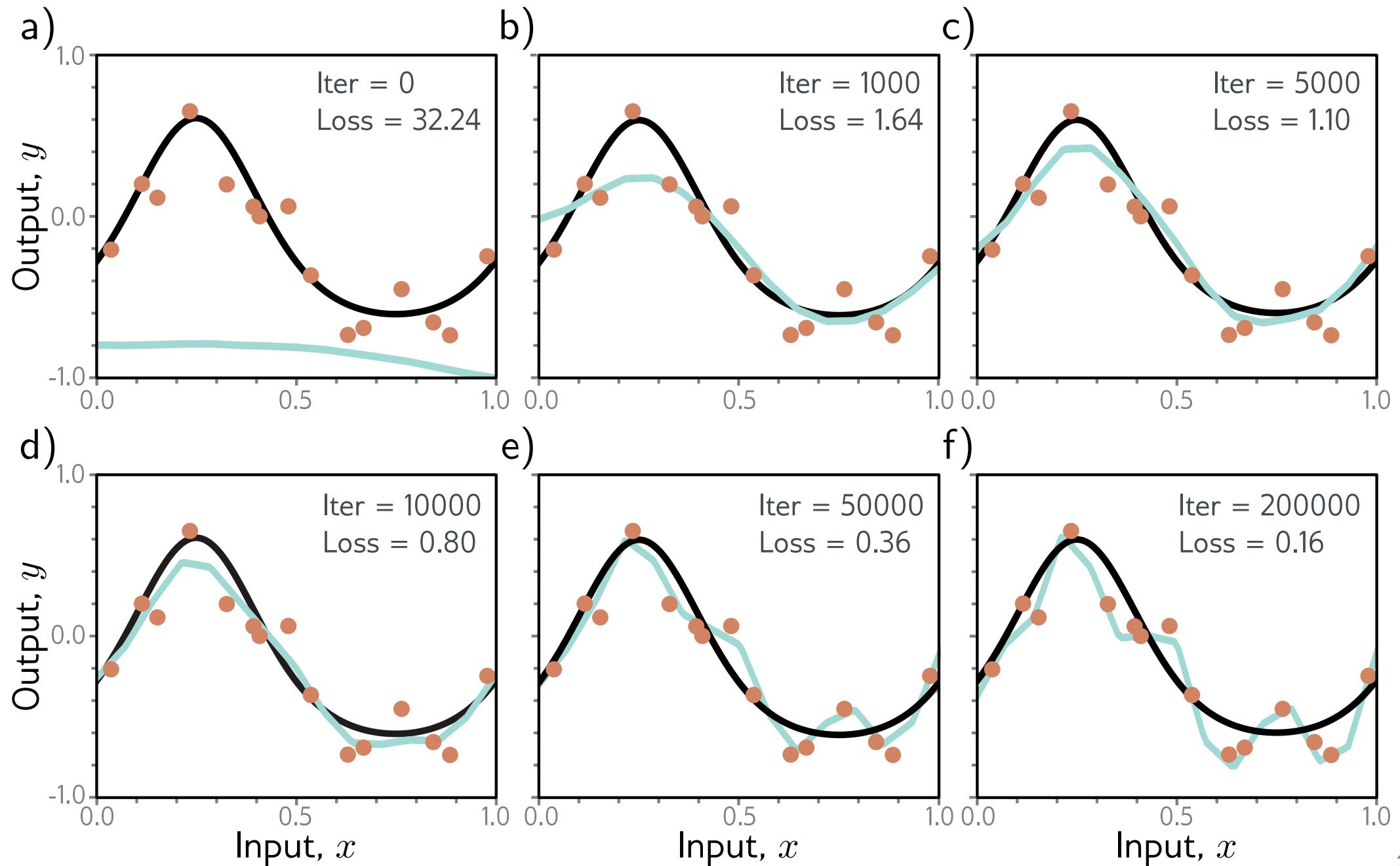


Regularization

- 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

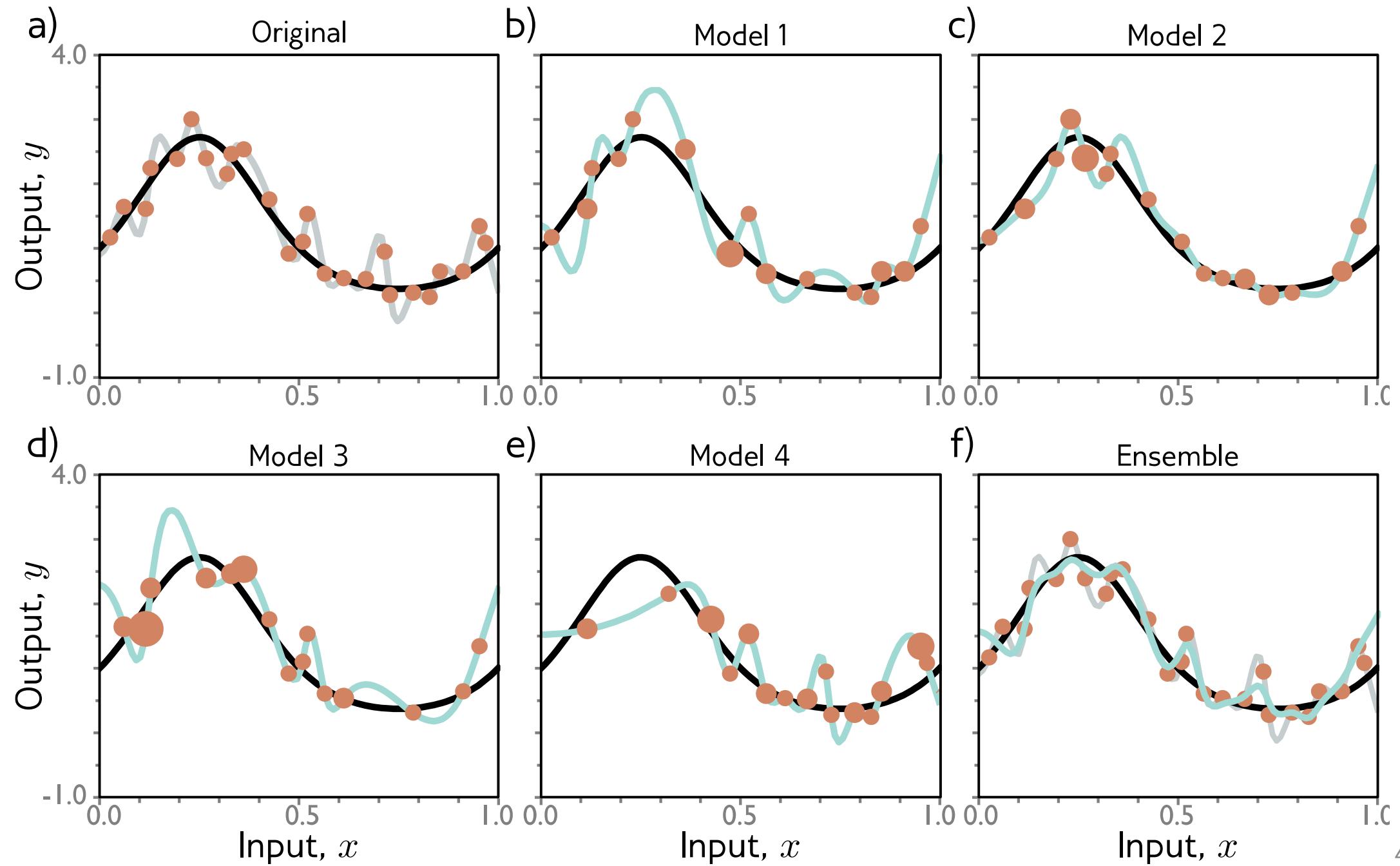


Regularization

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

Ensembling

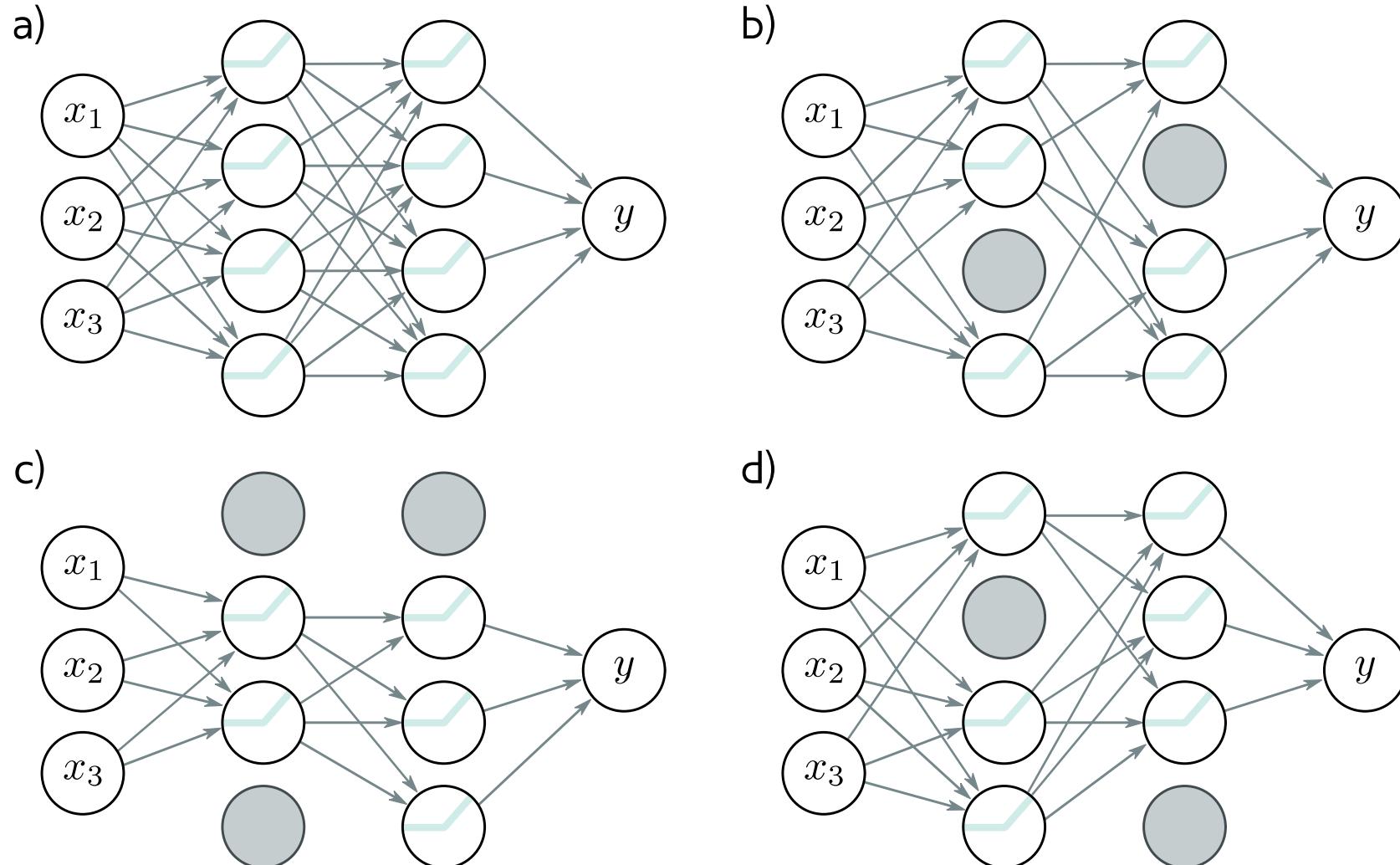
- Average together several models – an **ensemble**
- Can take mean or median
- Different initializations / different models
- Different subsets of the data resampled with replacements -- **bagging**



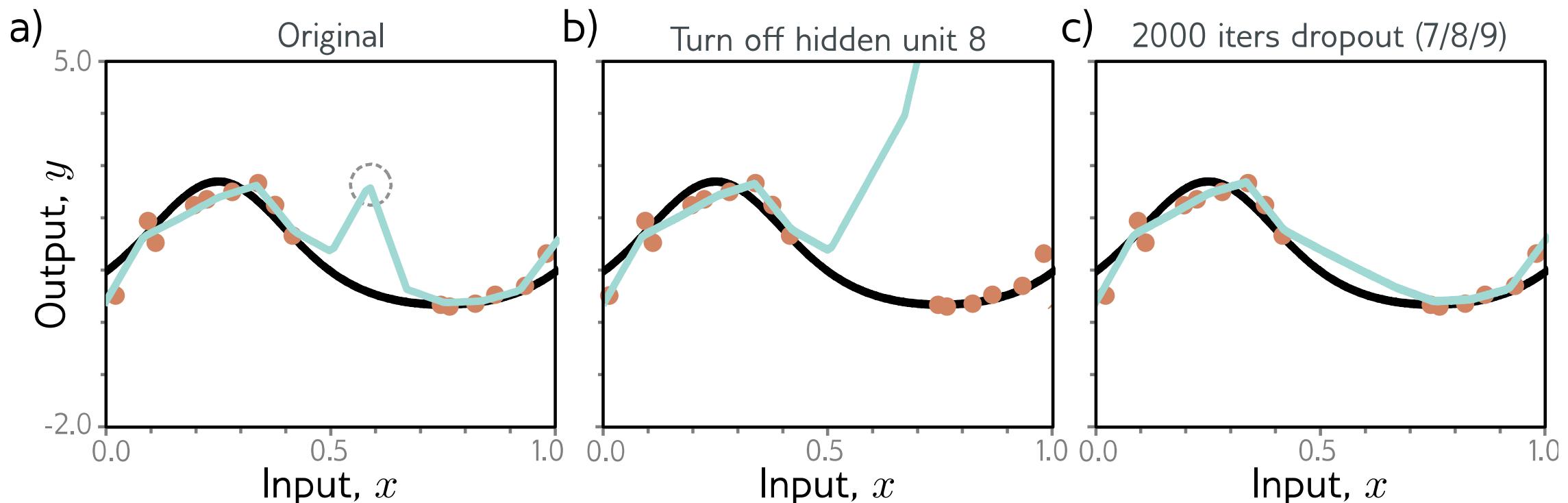
Regularization

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

Dropout



Dropout

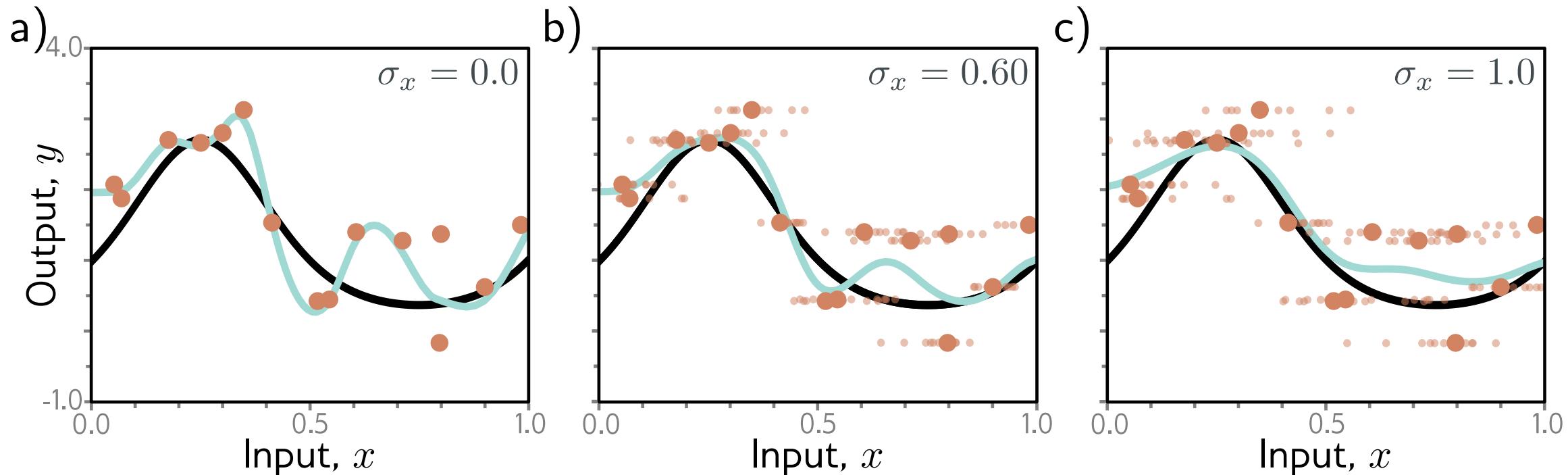


Can eliminate kinks in function that are far from data and don't contribute to training loss

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

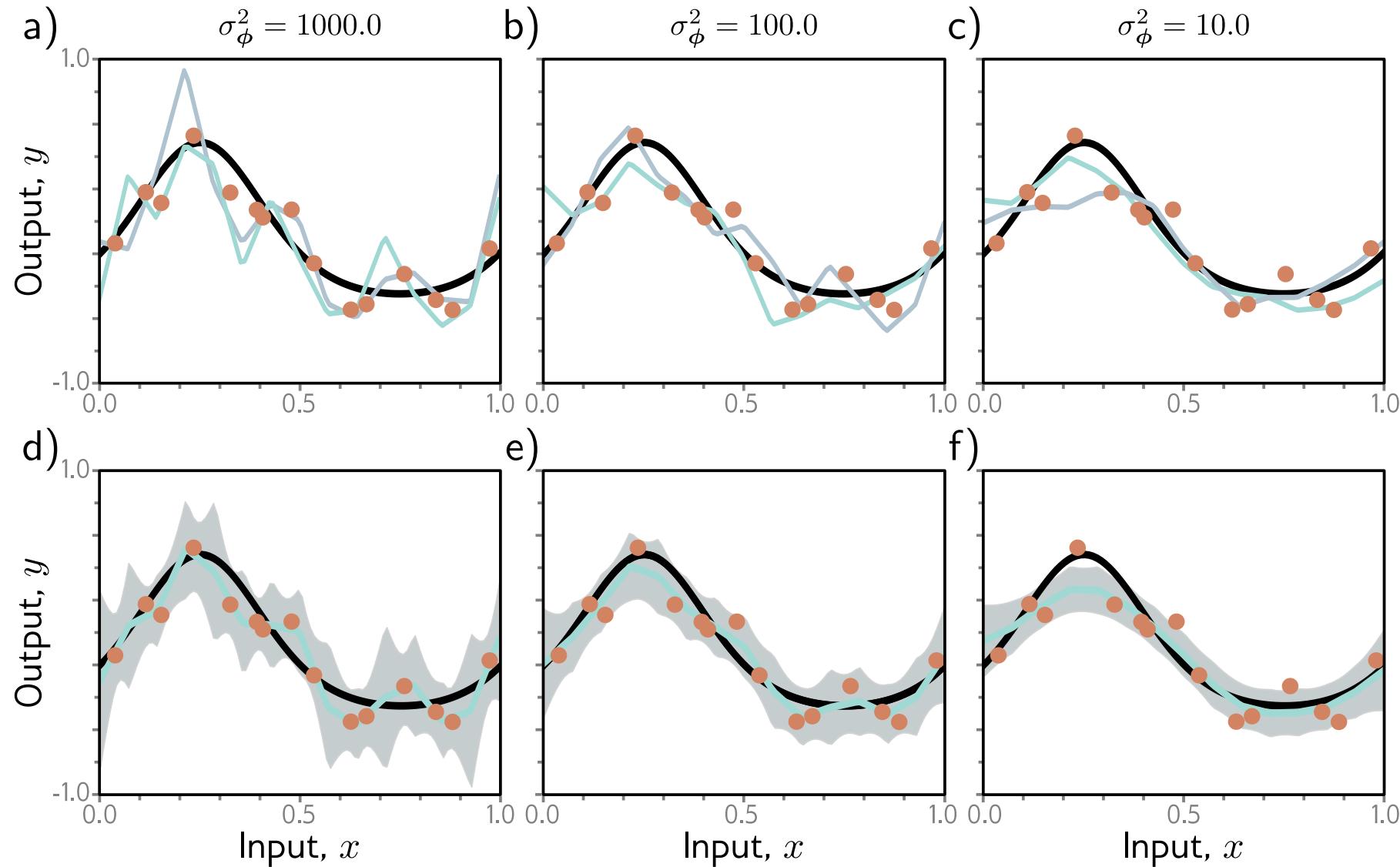
Prior info about parameters

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

- Take all possible parameters into account when make prediction

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

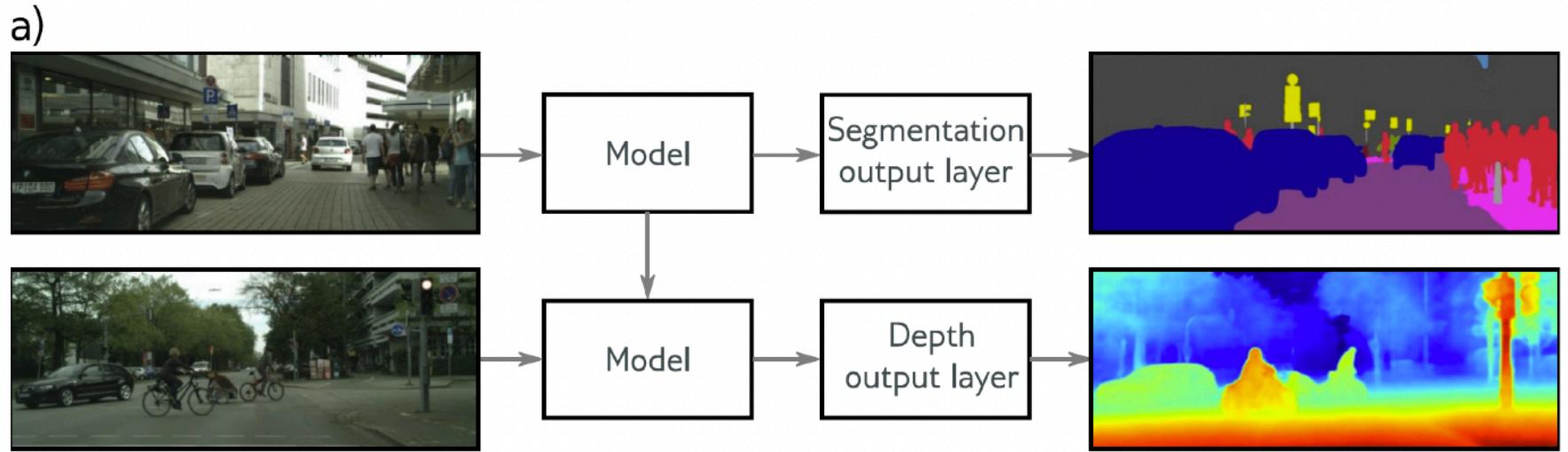
Bayesian approaches



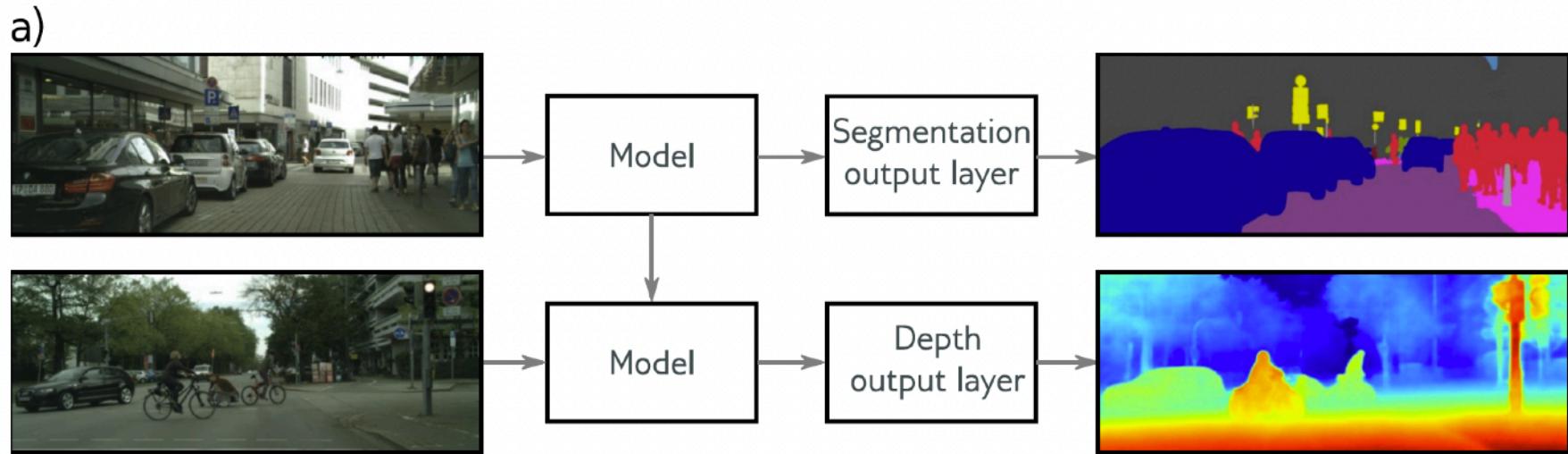
Regularization

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

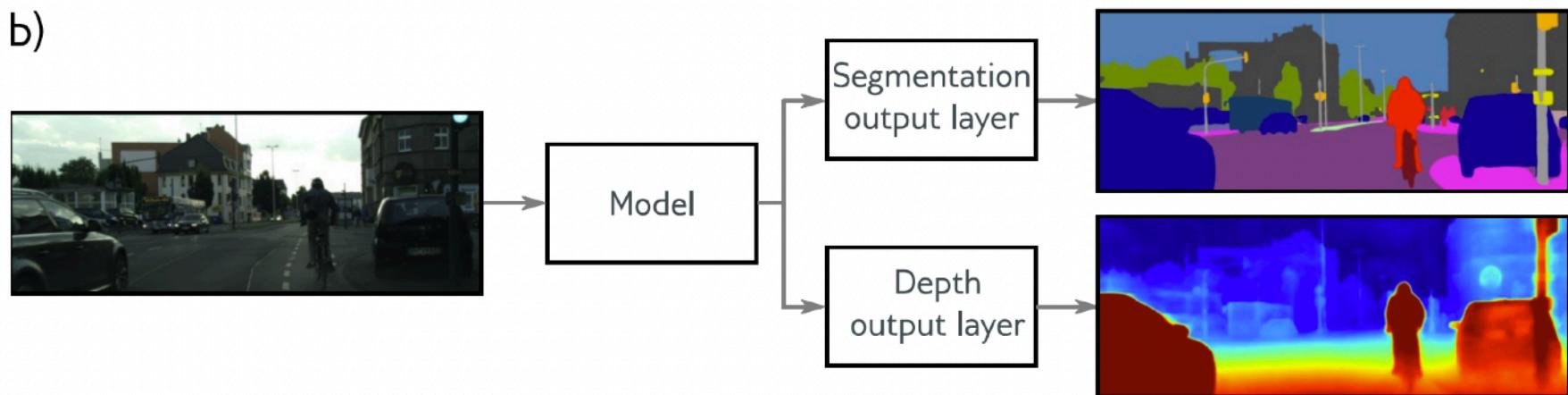
- Transfer learning



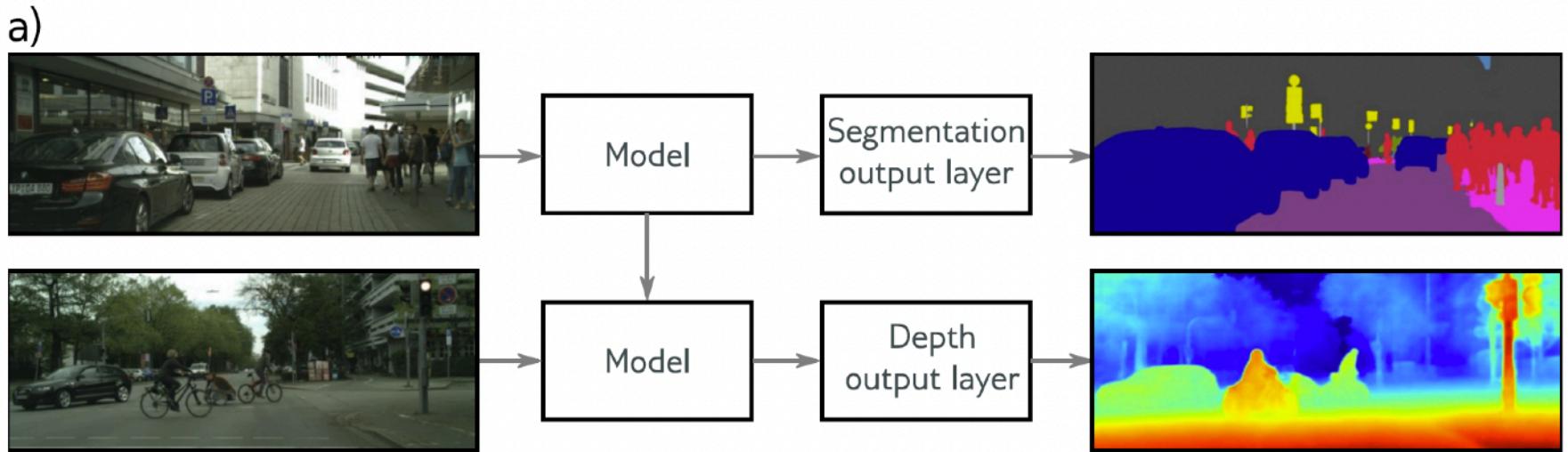
- Transfer learning



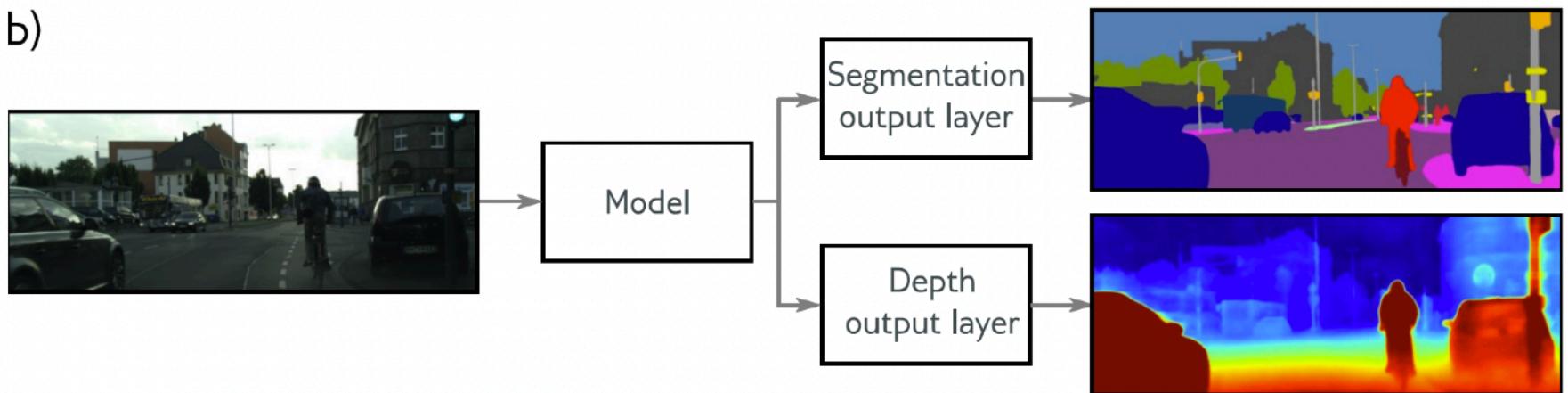
- Multi-task learning



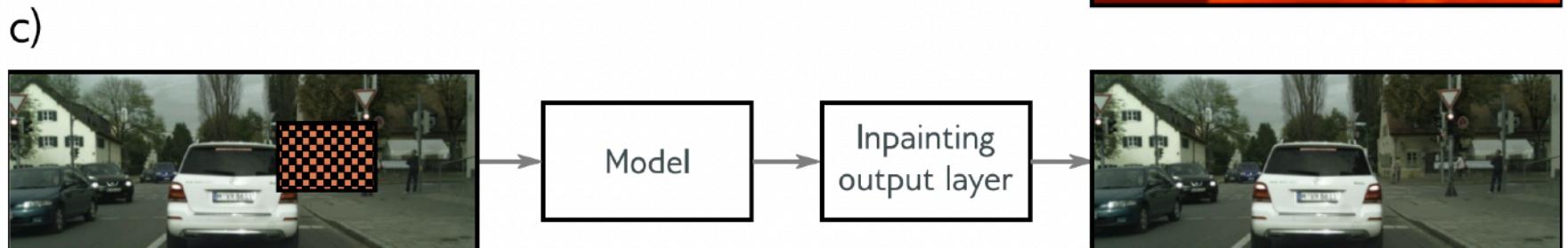
- Transfer learning



- 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



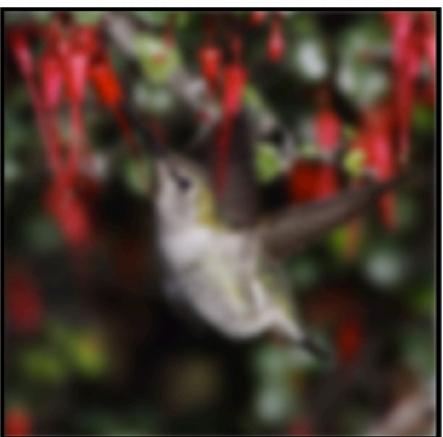
d) Vertical stretch



e) Color balance



f) Blur



g) Vignette



h) Pincushion



Regularization overview

