

The problem of overfitting

(larger)?

MLP Two layers, ReLU
3 hidden units

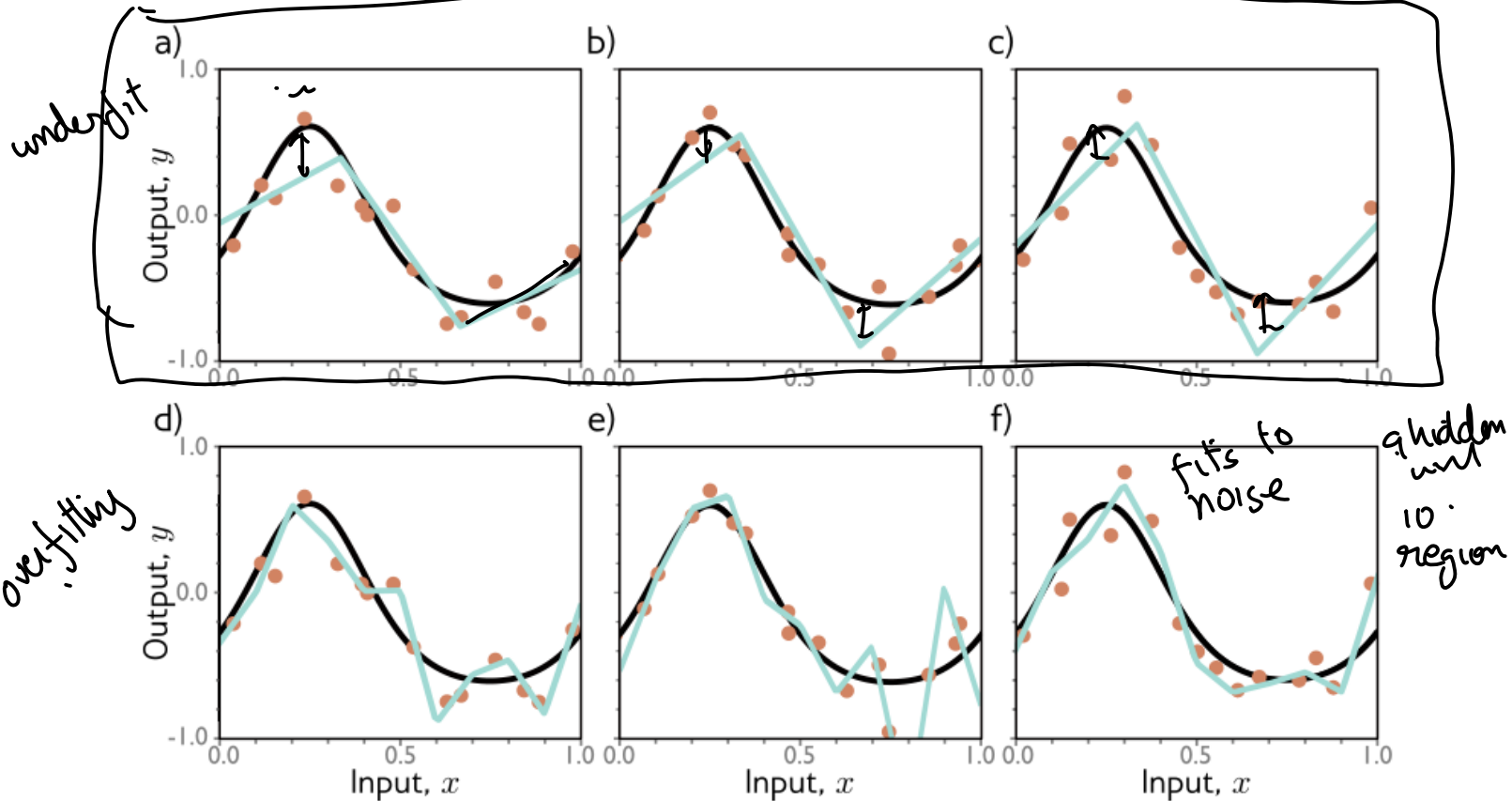


Figure 8.8 Overfitting. a-c) A model with three regions is fit to three different datasets of fifteen points each. The result is very similar in all three cases (i.e., the variance is low). d-f) A model with ten regions is fit to the same datasets. The additional flexibility does not necessarily produce better predictions. While these three models each describe the training data better, they are not necessarily closer to the true underlying function (black curve). Instead, they overfit the data and describe the noise, and the variance (difference between fitted curves) is larger.

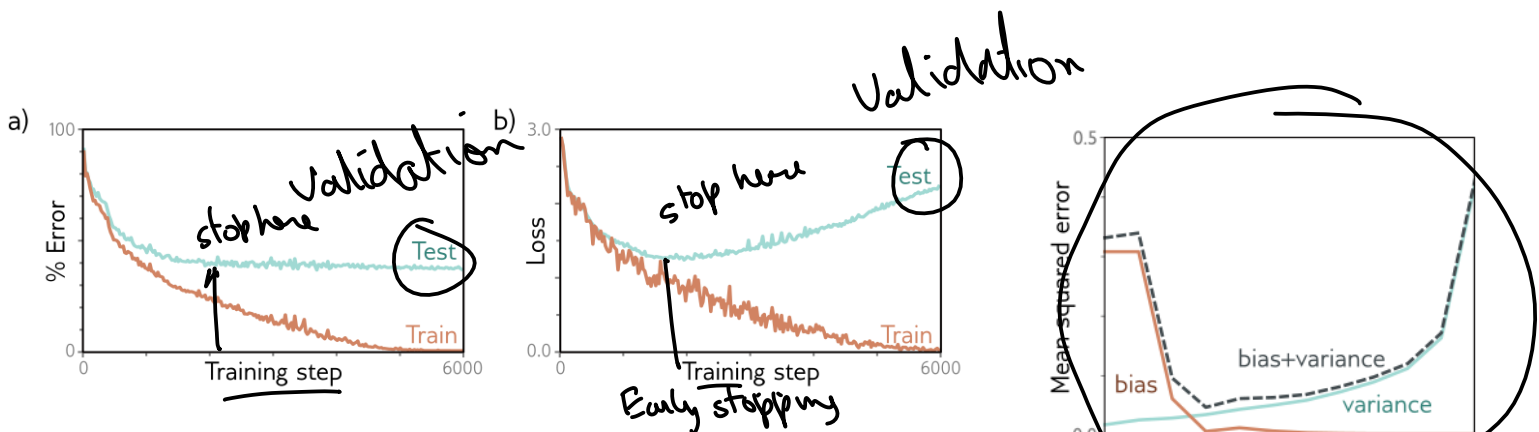
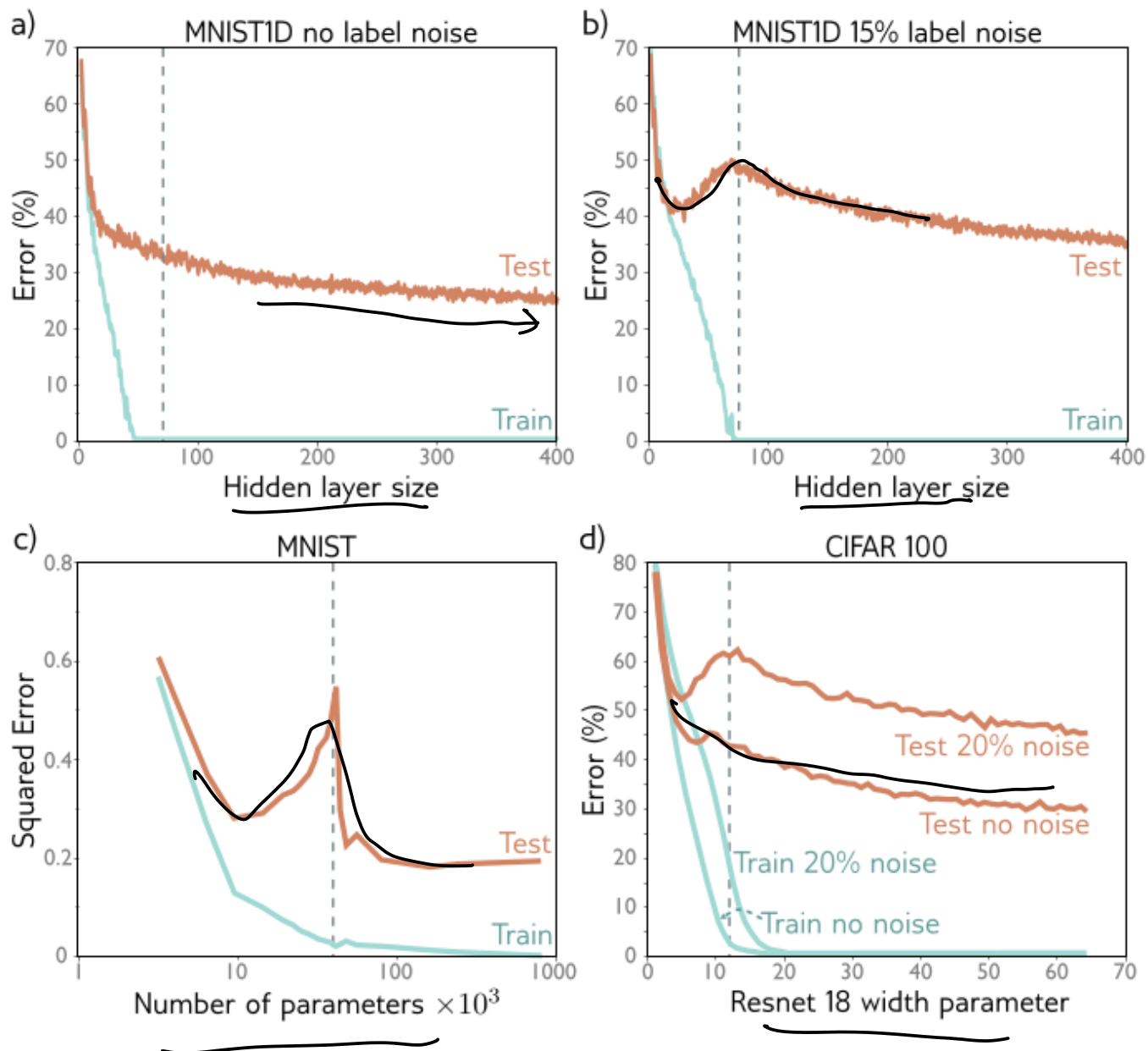
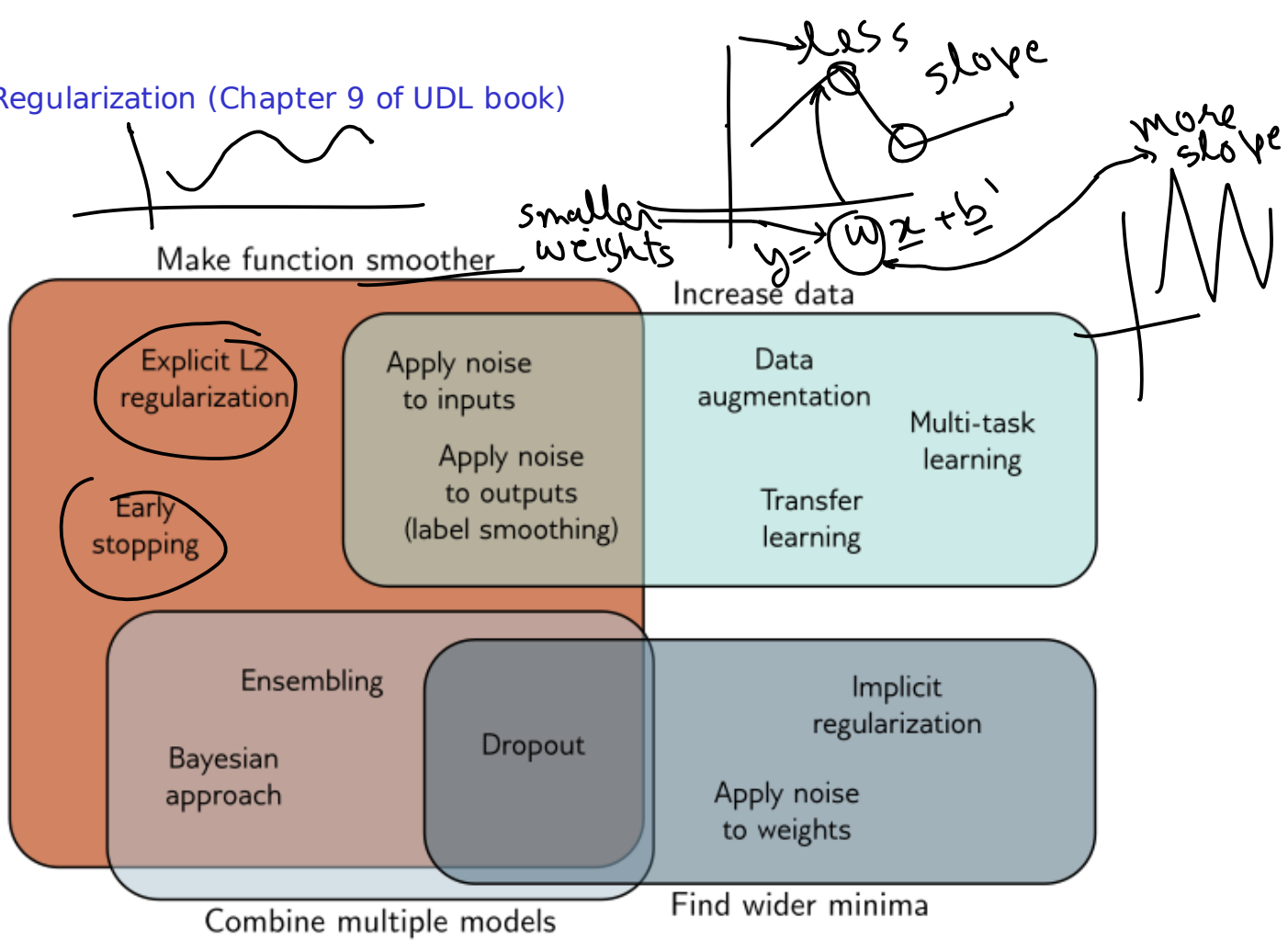


Figure 8.2 MNIST-1D results. a) Percent classification error as a function of the training step. The training set errors decrease to zero, but the test errors do not decrease below $\sim 40\%$. This model doesn't generalize well to new test data. b) Loss as a function of the training step. The training loss decreases steadily towards zero. The test loss decreases at first but then increases as the model becomes increasingly confident about its (wrong) predictions.

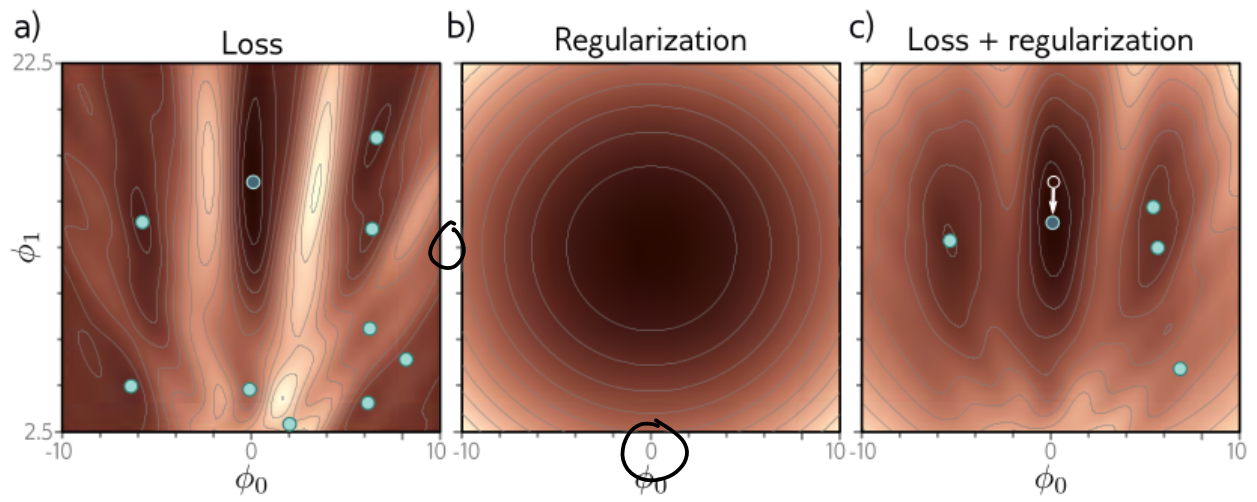
Aside: Model capacity in neural networks is weird though



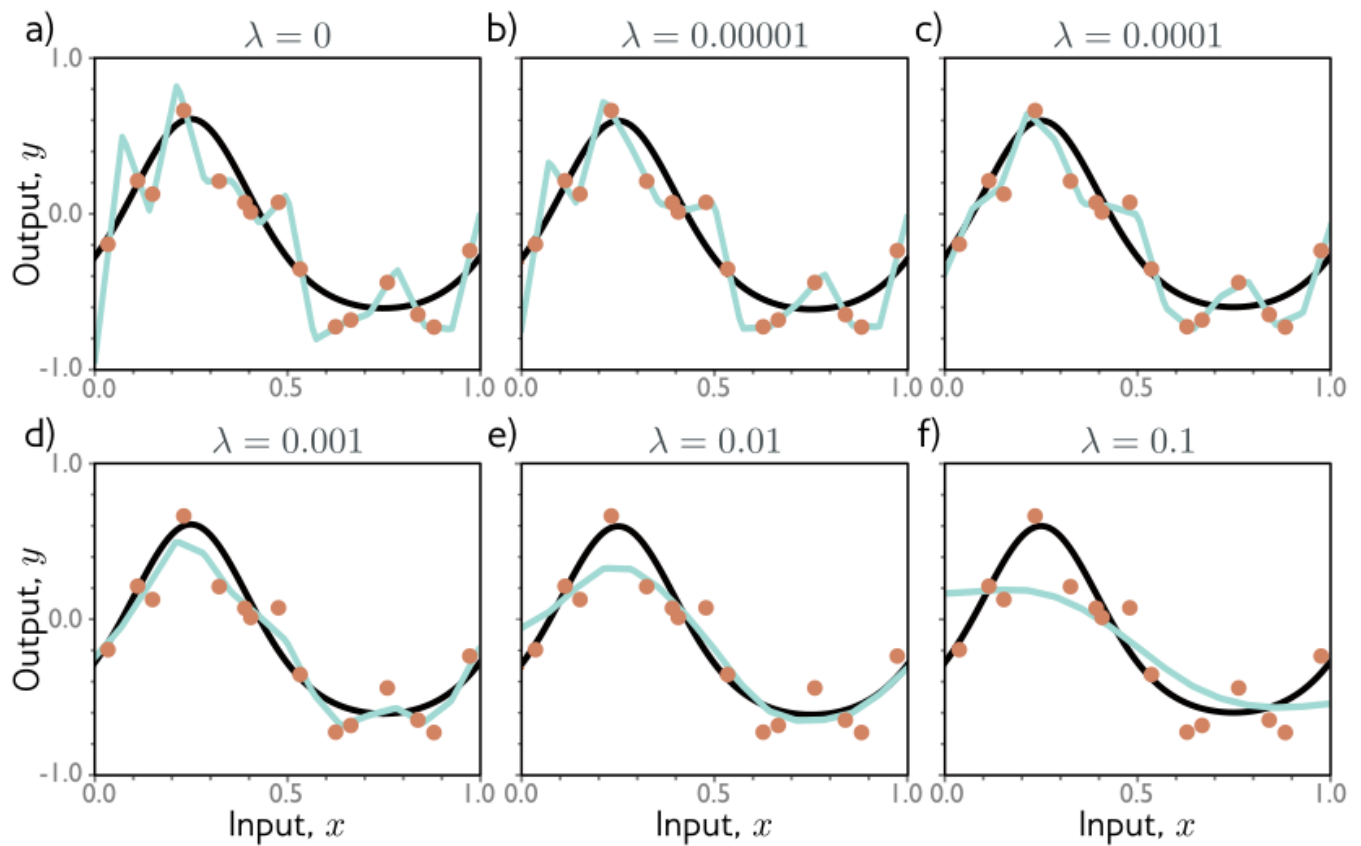
Regularization (Chapter 9 of UDL book)



Explicit regularization



L2 Regularization



L2-regularization

Mean squared loss for regression

$$L(D; W) = \sum_{i=1}^D \|y_i - f(x_i; W)\|_2^2$$

$$f(x_i; W) = \underline{w}_2^T a(\underline{w}_1, x_i + \underline{b}_1)$$

$$W = \{\underline{w}_2, \underline{w}_1, \underline{b}_1\}$$

For L2-regularization

$$L(D; W) = \sum_{i=1}^D \|y_i - f(x_i; W)\|_2^2 + \lambda \|\underline{w}_2\|_2^2$$

fit to data

positive scalar
L2-norm
regularization term
keep the weights small

$$W^* = \arg \min_W L(D; W)$$

L2-norm? = vector magnitude

$$\|\underline{w}_2\|_2 = \sqrt{w_{21}^2 + w_{22}^2 + \dots + w_{2n}^2}$$
$$= (w_{21}^2 + w_{22}^2 + \dots + w_{2n}^2)^{1/2}$$

L1-norm

$$\|\underline{w}_2\|_1 = (|w_{21}| + |w_{22}| + \dots + |w_{2n}|)$$

$$\text{Lp-norm } \|\underline{w}_2\|_p = ((w_{21})^p + w_{22}^p + \dots + w_{2n}^p)^{1/p}$$

$$\underline{w}_2 = \begin{bmatrix} w_{21} \\ w_{22} \\ \vdots \\ w_{2n} \end{bmatrix}$$

For Ex 1

$$w = \{ \underline{w}_2, w_1, b_1 \}$$

$\frac{\partial}{\partial \underline{w}} \lambda \|\underline{w}\|_2^2 = \frac{\partial}{\partial \underline{w}} \lambda \underline{w}^T \underline{w} = 2\lambda \underline{w}$

$$L(D; w) = \underbrace{\| \quad \|}_{\text{Data term}} + \lambda \underbrace{\|\underline{w}\|_2^2}$$

In neural networks, regularization is restricted to weights not biases and it is called weight decay

Regularization term

$$\lambda \|\underline{w}\|_2^2 = \lambda \|\underline{w}_2\|_2^2 + \lambda \|\text{flatten}(w_1)\|_2^2$$

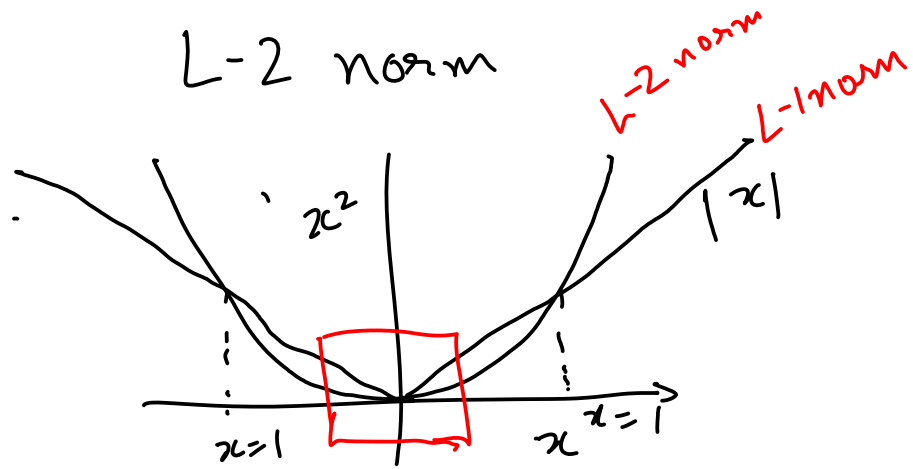
$$= \lambda \left\| \begin{array}{c} \underline{w}_2 \\ \text{flatten}(w_1) \end{array} \right\|_2^2$$

L-1 vs L-2
Regularization

$$x = 10^{-3}$$

$$|x| = 10^{-3}$$

$$x^2 = 10^{-6} \text{ much smaller}$$



Regularizing with L-1 norm, then smaller values are made even smaller

L-1 regularization leads to

sparse weights

large weights are okay but smaller weights become zero

without regularization

$$\underline{w} = \begin{bmatrix} 1.2 \\ 10.4 \\ 10^{-6} \\ \vdots \end{bmatrix}$$

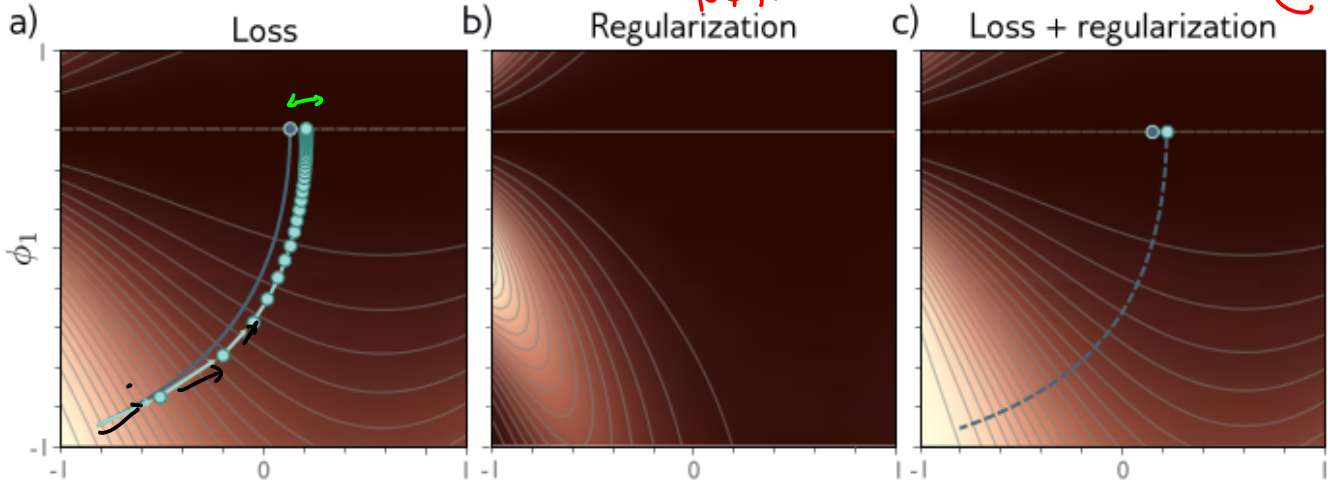
$$\xrightarrow{L2} \begin{bmatrix} 0.8 \\ 5.2 \\ 10^{-4} \\ \vdots \end{bmatrix}$$

$$\xrightarrow{L-1} \begin{bmatrix} 0.9 \\ 8.4 \\ 0 \\ \vdots \end{bmatrix}$$

$$\underline{w}_{t+1} = \underline{w}_t - \alpha \nabla_{\underline{w}_t} \ell(\mathcal{D}; \underline{w}_t) \quad \left| \begin{array}{l} \text{Theoretically, } \underline{w}(t+h) = \underline{w}(t) - h \nabla_{\underline{w}} \ell(\mathcal{D}; \underline{w}(t)) \\ \alpha \rightarrow \frac{\alpha \|\partial L\|}{4 \|\partial \phi\|} \quad \lim_{h \rightarrow 0} \end{array} \right.$$

Implicit regularization in gradient descent

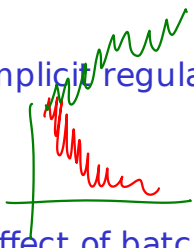
$$\tilde{L}_{GD}[\phi] = L[\phi] + \frac{\alpha}{4} \left\| \frac{\partial L}{\partial \phi} \right\|^2$$



Bigger step size $\alpha \rightarrow$ More regularization \rightarrow Less overfitting
Learning rate

Figure 9.3 Implicit regularization in gradient descent. a) Loss function with family of global minima on horizontal line $\phi_1 = 0.61$. Dashed blue line shows continuous gradient descent path starting in bottom left. Cyan trajectory shows discrete gradient descent with step size 0.1 (first few steps shown explicitly as arrows). The finite step size causes the paths to diverge and to reach a different final position. b) This disparity can be approximated by adding a regularization term to the continuous gradient descent loss function that penalizes the squared gradient magnitude. c) After adding this term, the continuous gradient descent path converges to the same place that the discrete one did on the original function.

Implicit regularization in stochastic gradient descent



Smaller batch sizes \rightarrow More regularization \rightarrow Less overfitting

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

$$= L[\phi] + \frac{\alpha}{4} \left\| \frac{\partial L}{\partial \phi} \right\|^2 + \frac{\alpha}{4B} \sum_{b=1}^B \left\| \frac{\partial L_b}{\partial \phi} - \frac{\partial L}{\partial \phi} \right\|^2$$

Batch loss, Full data loss, variance of batch loss, total batches in the Dataset

Effect of batch size and learning rate

Thumb rule: pick the biggest learning rate for which your loss does not diverge

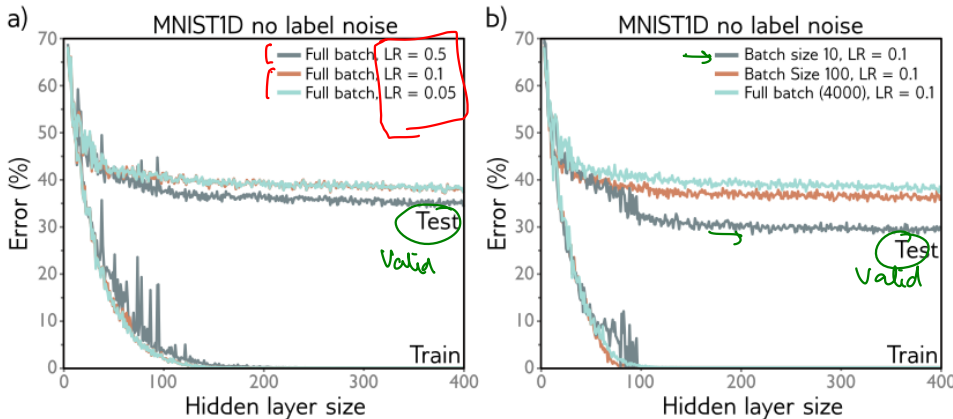


Figure 9.5 Effect of learning rate and batch size for 4000 training and 4000 test examples from MNIST-1D (see figure 8.1) for a neural network with two hidden layers. a) Performance is better for large learning rates than for intermediate or small ones. In each case, the number of iterations is $6000 \times$ the learning rate so each solution has the opportunity to move the same distance. b) Performance is superior for smaller batch sizes. In each case the number of iterations was chosen so that the training data were memorized at roughly the same model capacity.

Ensemble methods

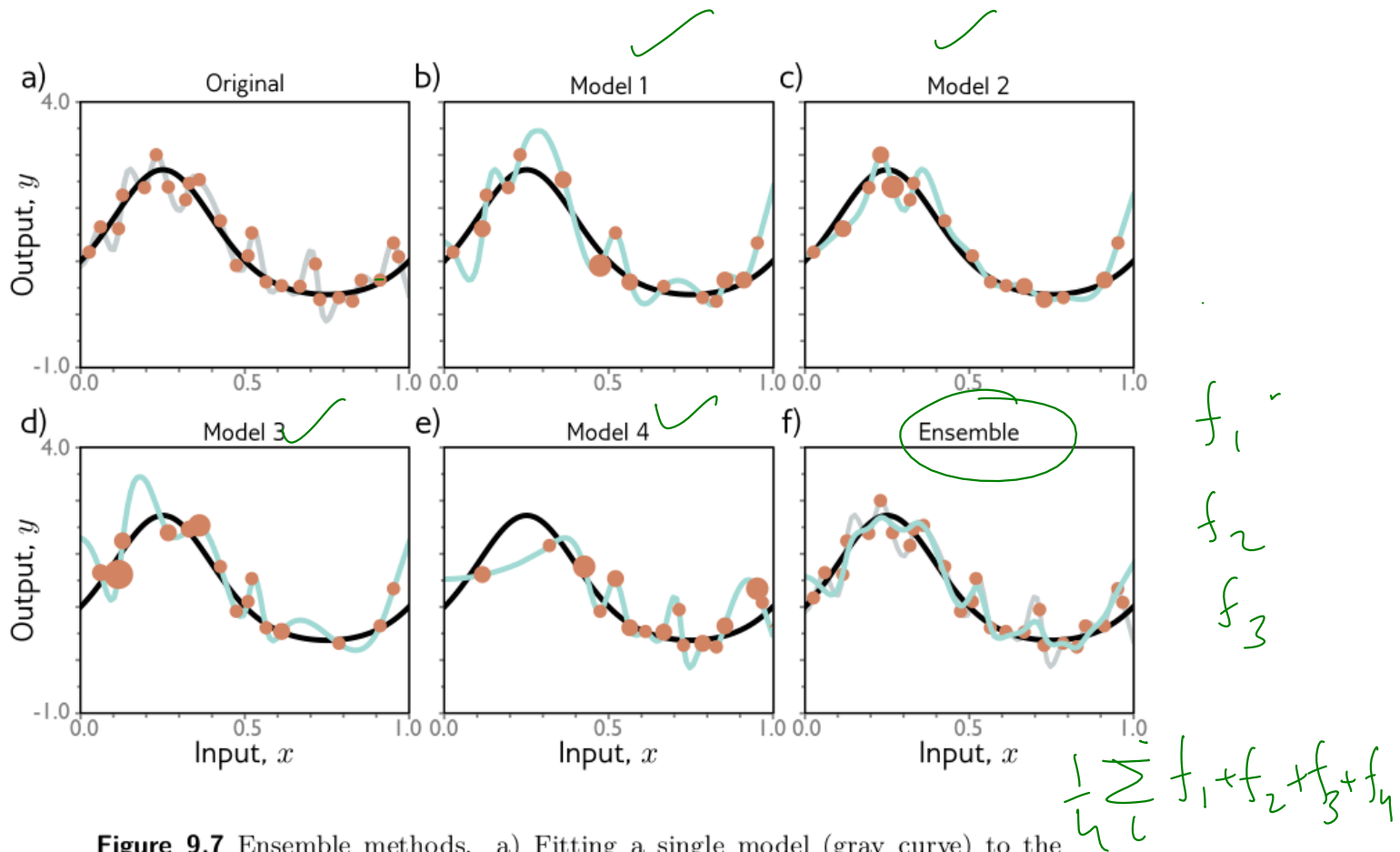


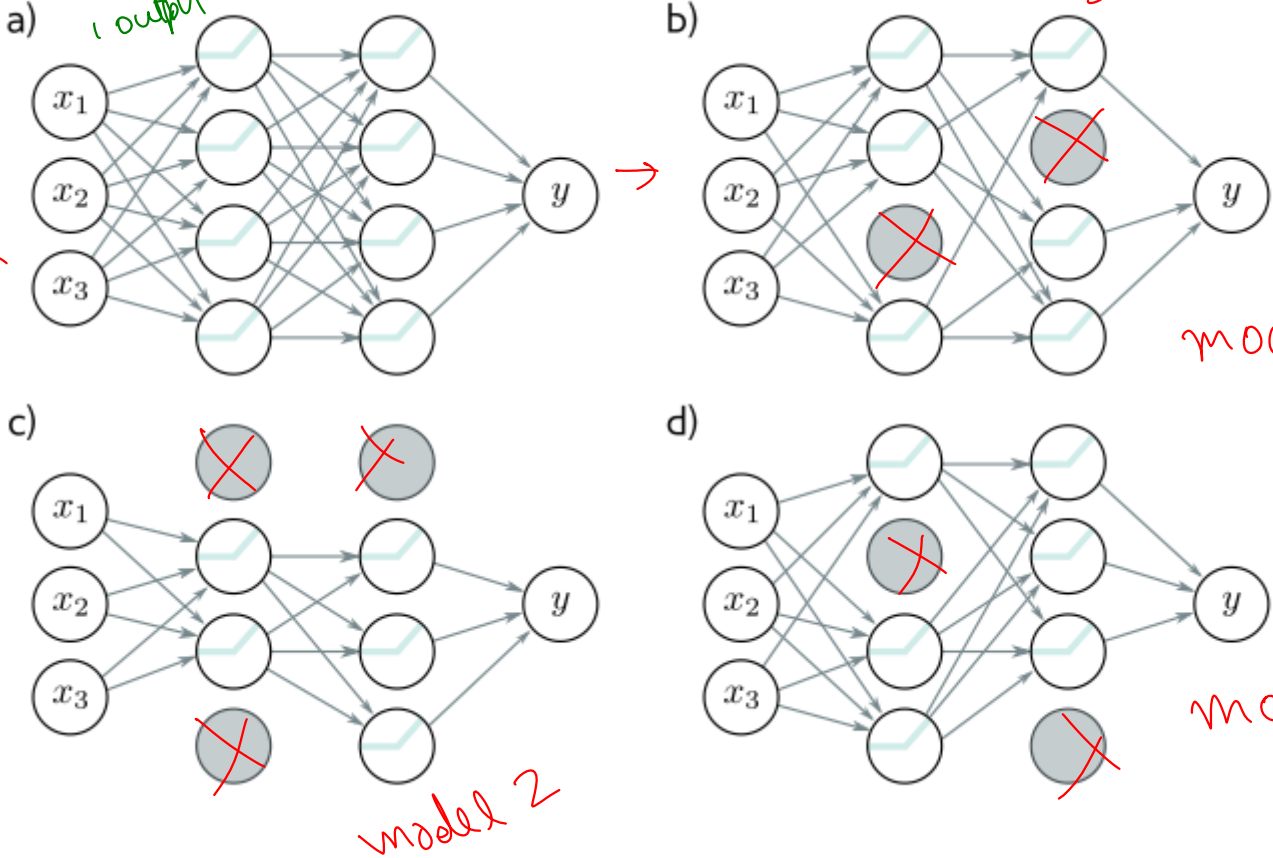
Figure 9.7 Ensemble methods. a) Fitting a single model (gray curve) to the entire dataset (orange points). b-e) Four models created by re-sampling the data with replacement (bagging) four times (size of orange point indicates number of times the data point was re-sampled). f) When we average the predictions of this ensemble, the result (cyan curve) is smoother than the result from panel (a) for the full dataset (gray curve) and will probably generalize better.

Dropout

3 inputs

3 Layer MLP

2 ReLU, (4,4) hidden units
1 output



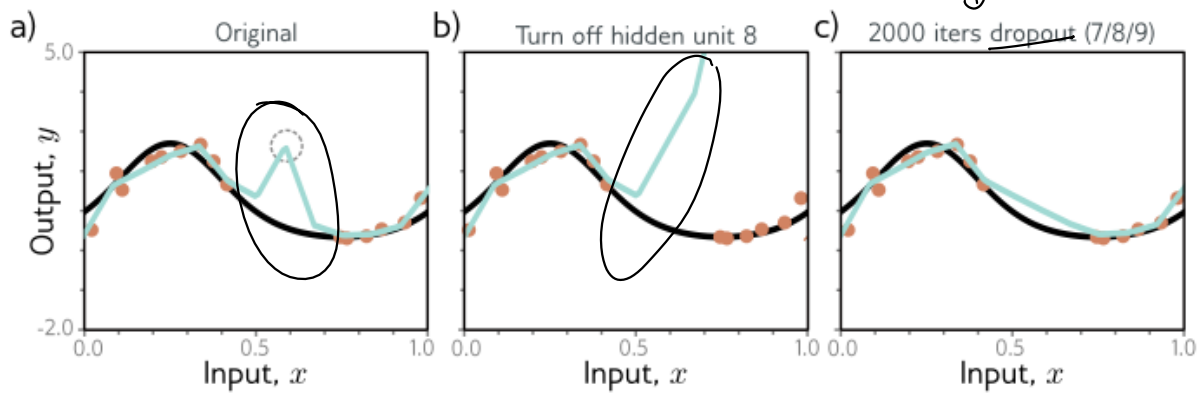
Training a random smaller network than the original

→ Being the average of model 1, 2, 3

Dropout layers → Create an implicit ensemble of subnetworks → Regularize
↓
reduce overfitting

Effect of dropout

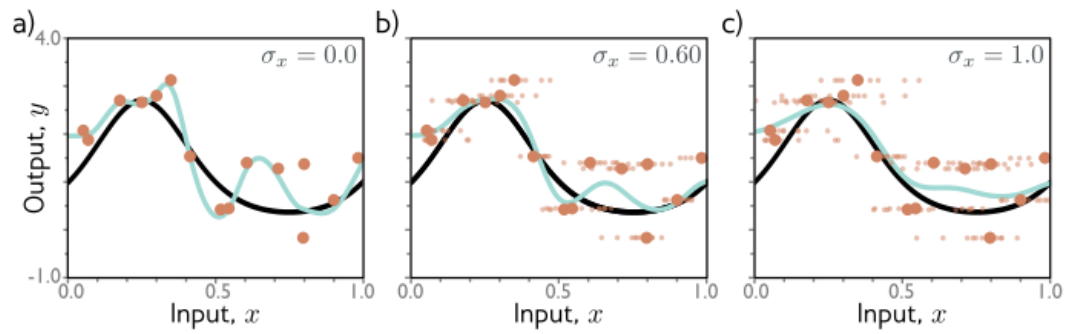
Dropout reduces dependence of a network on specific hidden units. Because any hidden unit can be turned off anytime.



Adding noise to each batch

You can add noise to the inputs, outputs or the weights during training time.

The training process will make your learned function to be robust to these noises making it "smoother"



Data augmentation

