

- Which function best fits this data? Why?
- What are some of the choices that we can make when designing a neural network?
- What might we infer about the differences in the neural networks used for each function?
- What can we conclude about the neural network weights in the third figure from looking at the highlighted region?

Data =  $\mathbf{X}, \mathbf{y}$



Training data =  $\mathbf{X}_{train}, \mathbf{y}_{train}$



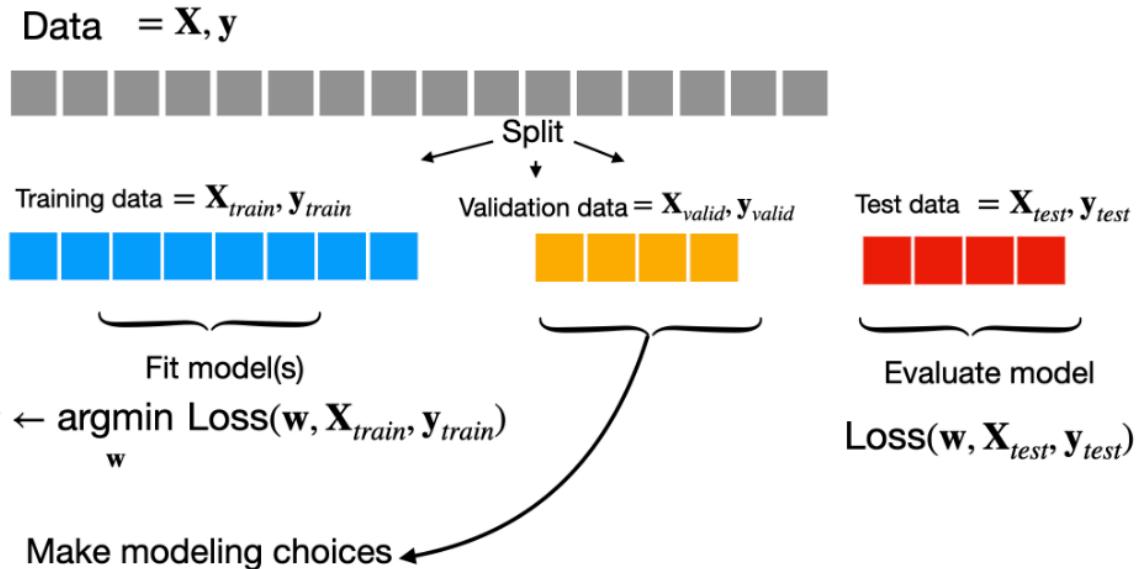
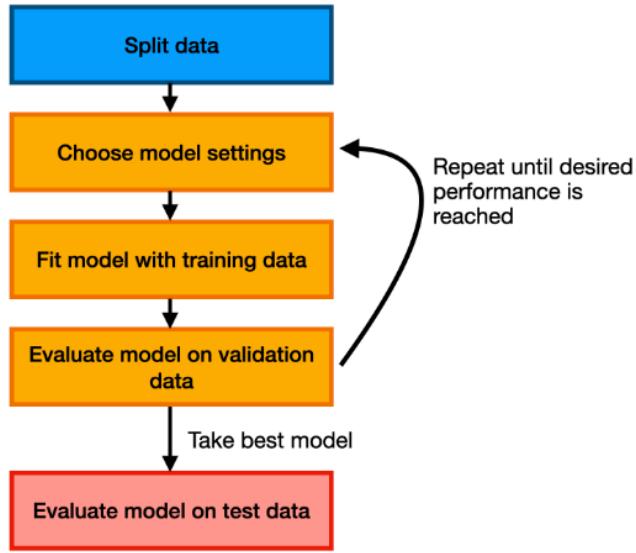
Split

Test data =  $\mathbf{X}_{test}, \mathbf{y}_{test}$

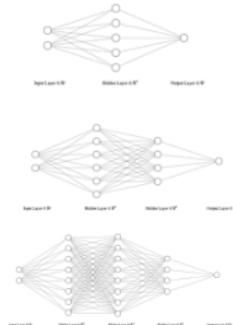


$\mathbf{w} \leftarrow \underset{\mathbf{w}}{\operatorname{argmin}} \text{Loss}(\mathbf{w}, \mathbf{X}_{train}, \mathbf{y}_{train})$

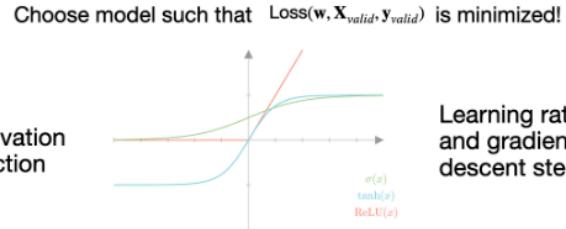
**Loss**( $\mathbf{w}, \mathbf{X}_{test}, \mathbf{y}_{test}$ )



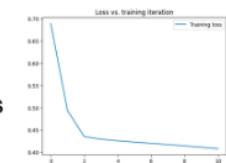
Number of layers and neurons

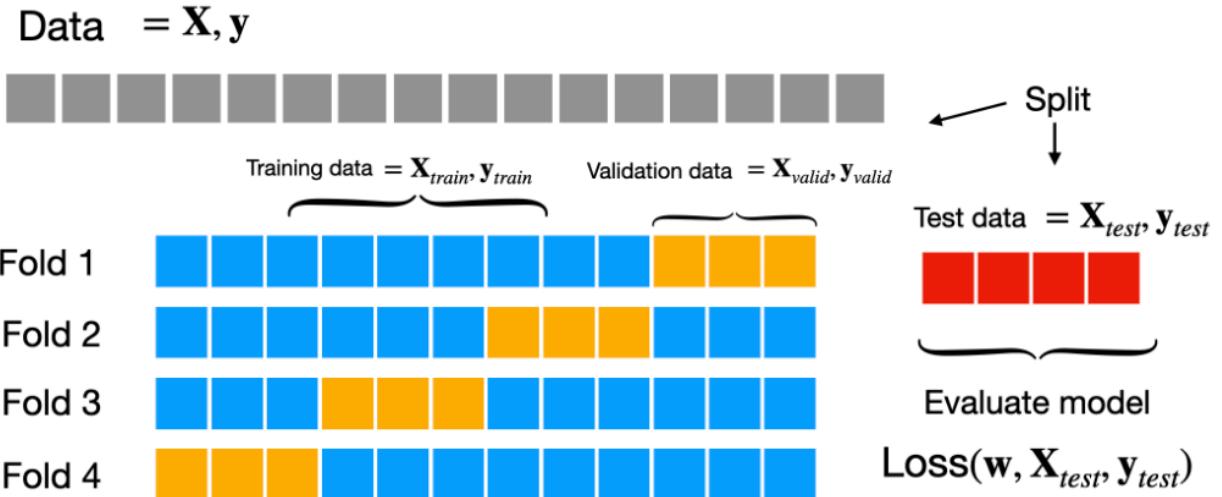
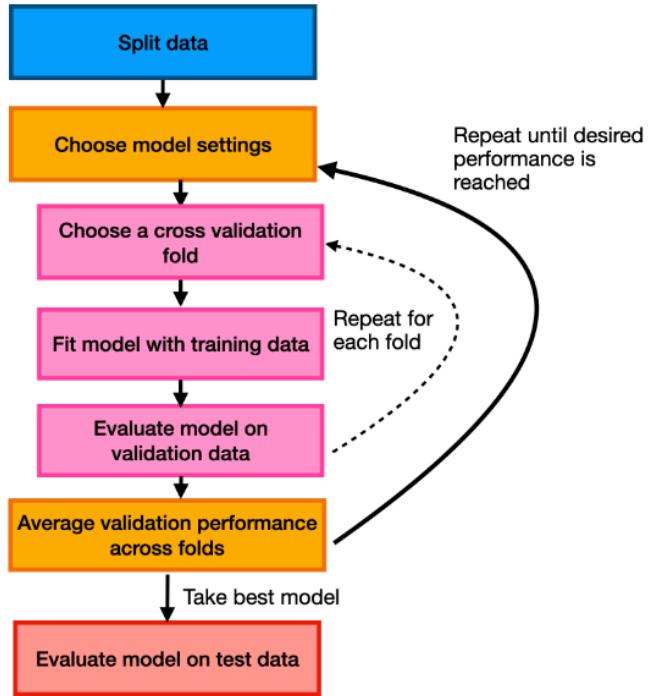


Activation function



Learning rate and gradient descent steps





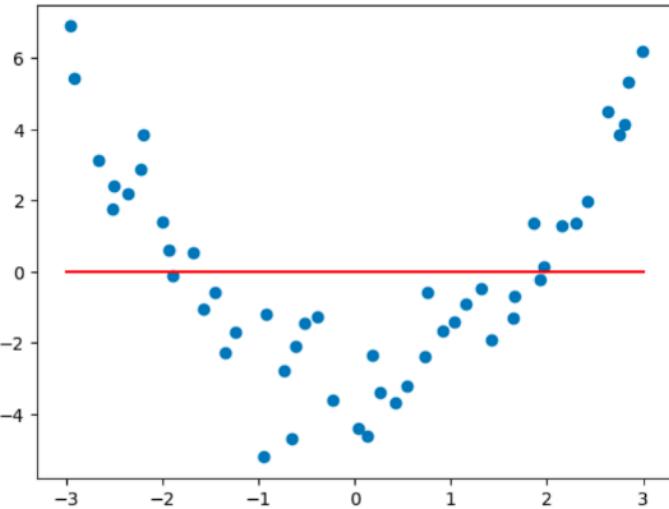
**For each fold:**

Fit model(s)

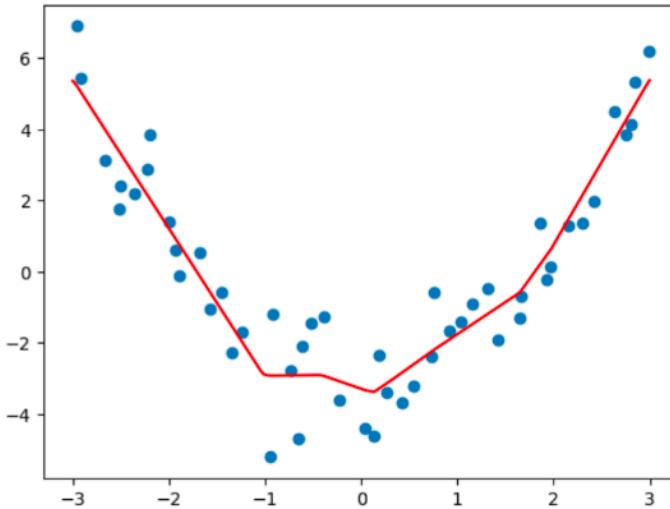
$$\mathbf{w} \leftarrow \underset{\mathbf{w}}{\operatorname{argmin}} \text{Loss}(\mathbf{w}, \mathbf{X}_{train}, \mathbf{y}_{train})$$

Evaluate model

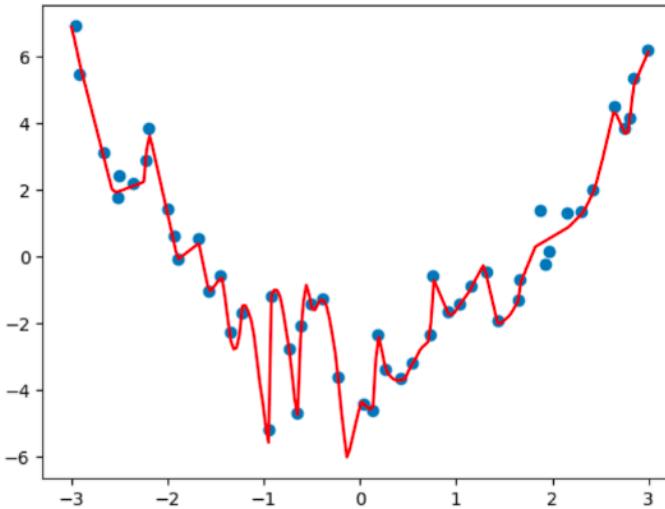
$$\text{Loss}(\mathbf{w}, \mathbf{X}_{valid}, \mathbf{y}_{valid})$$



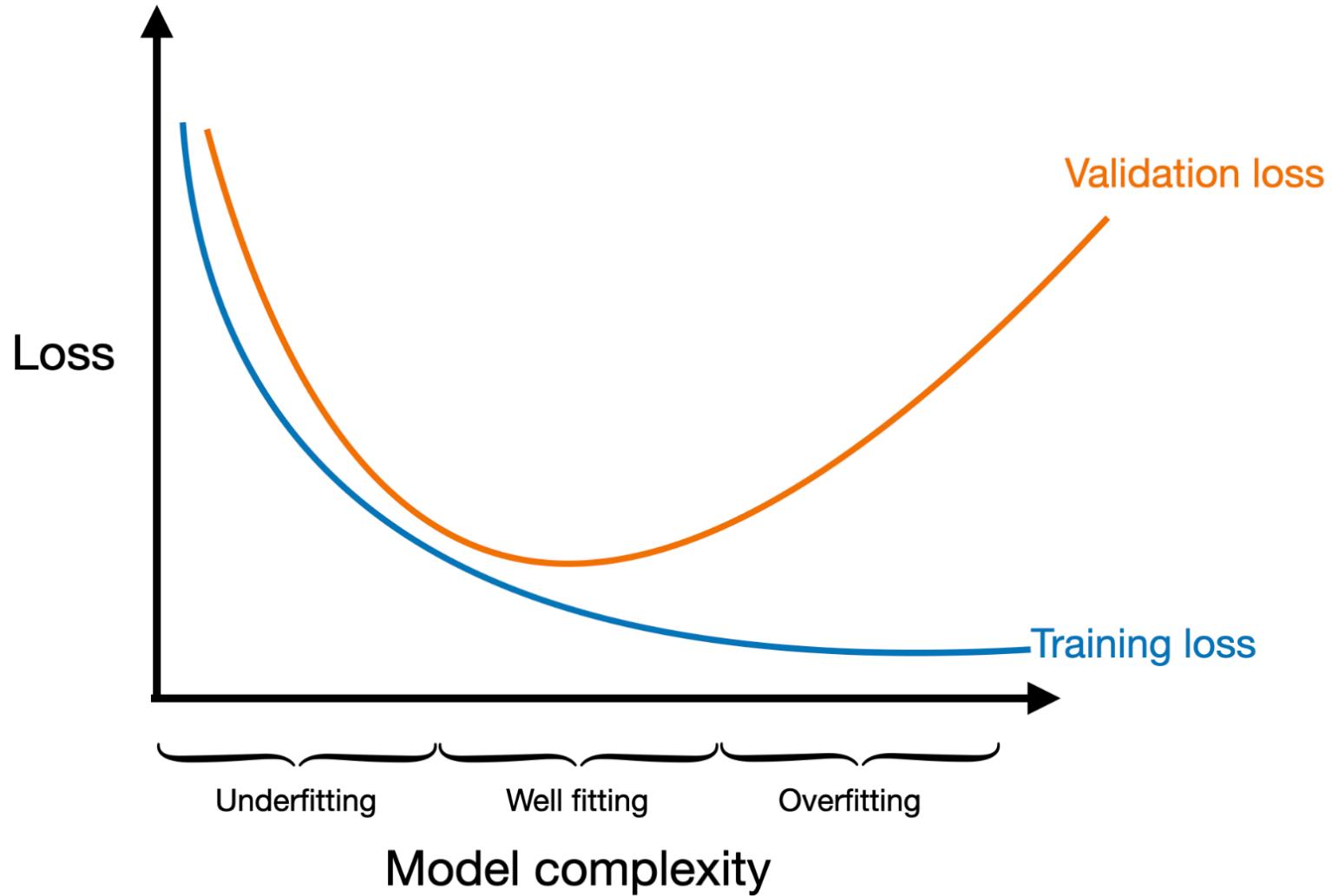
Underfitting

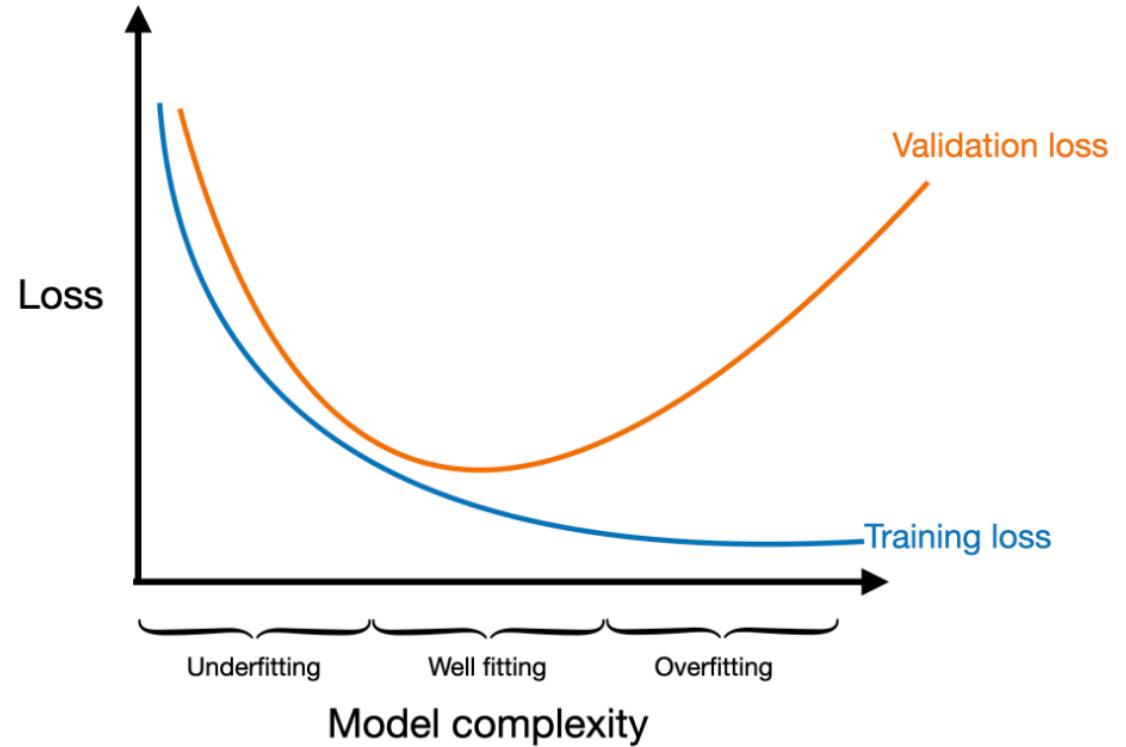
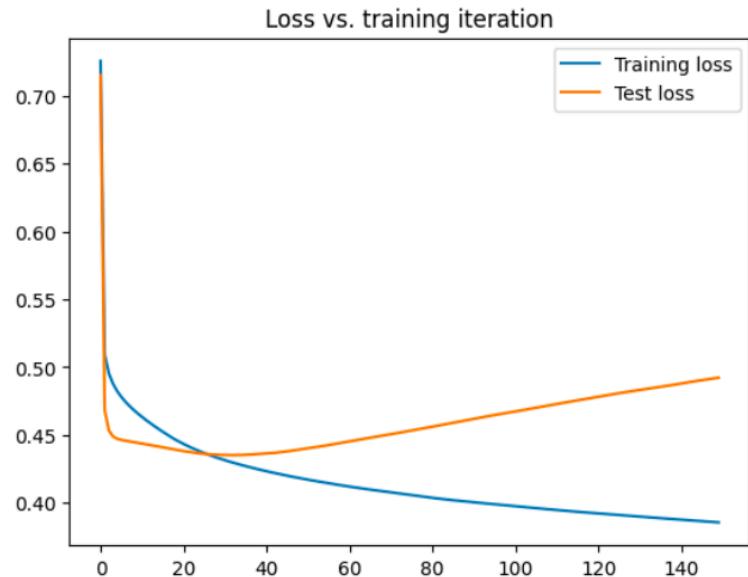


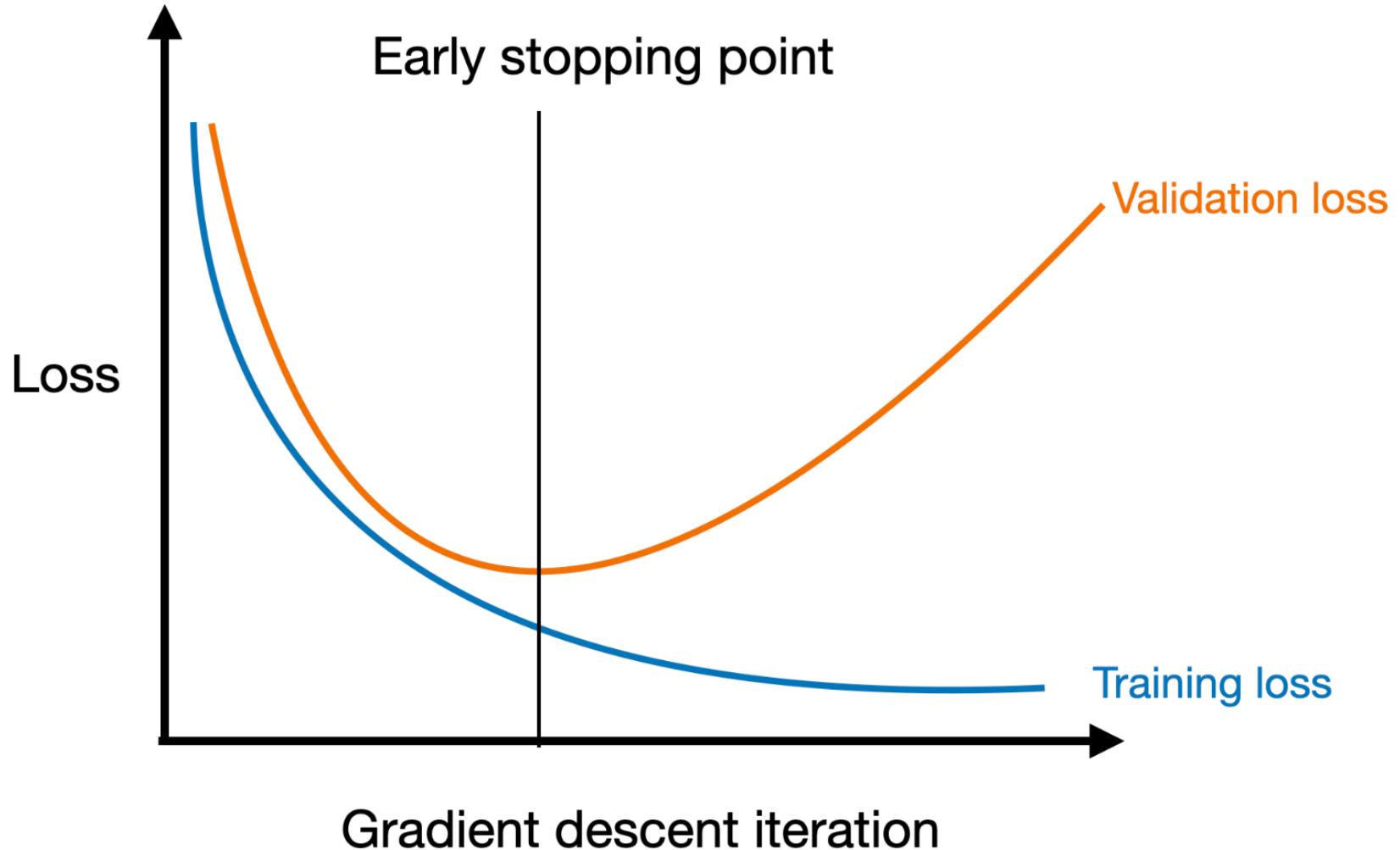
Well-fit

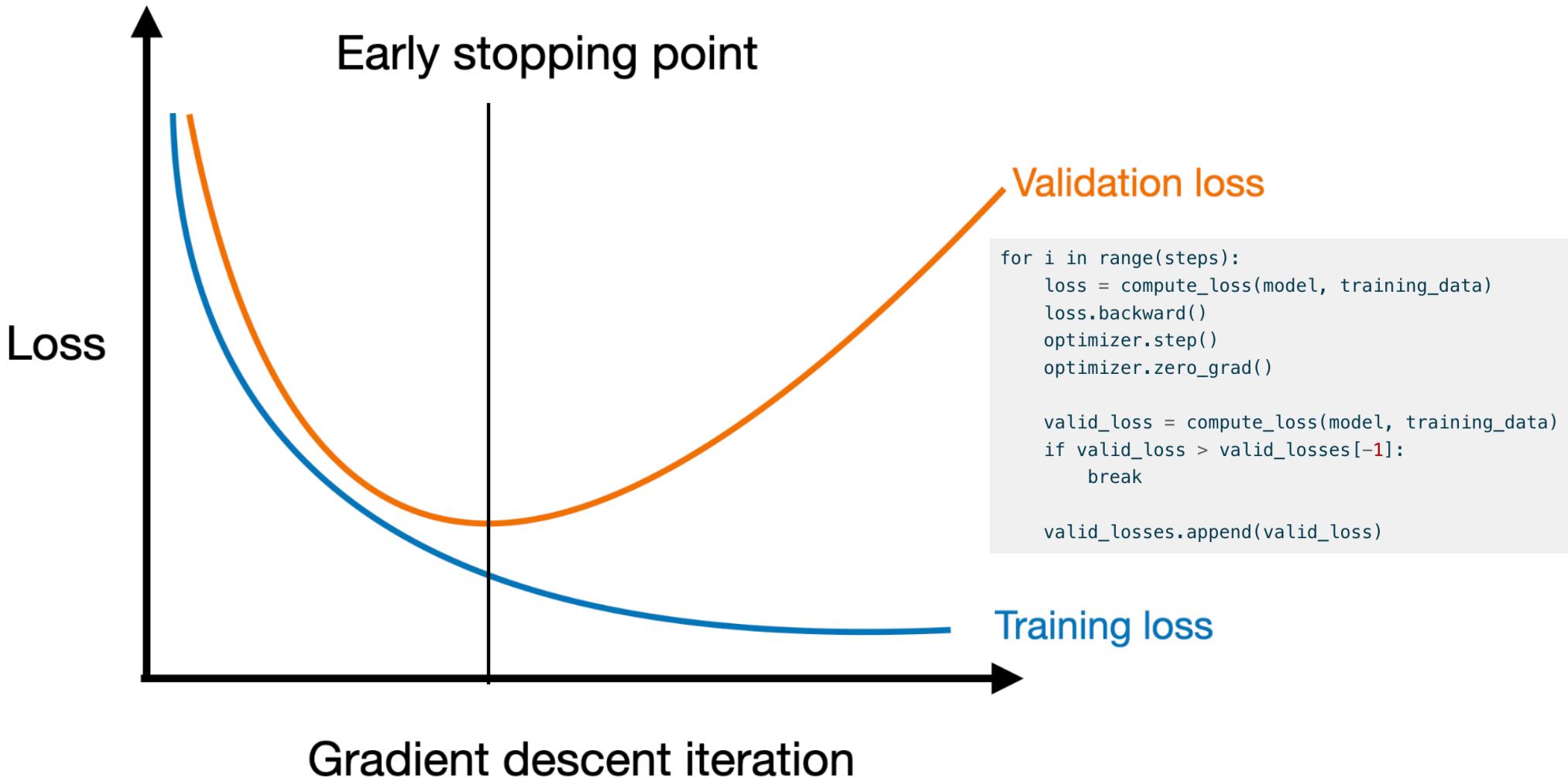


Overfit

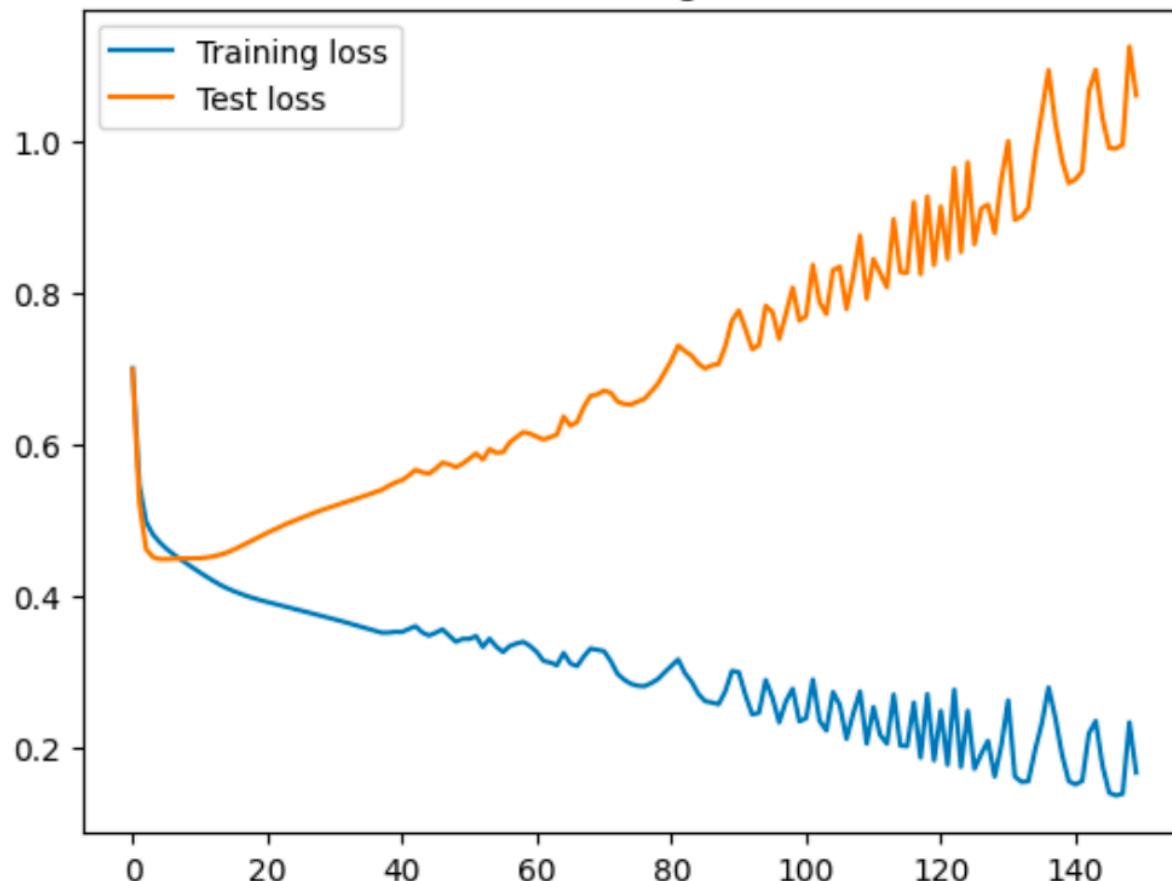




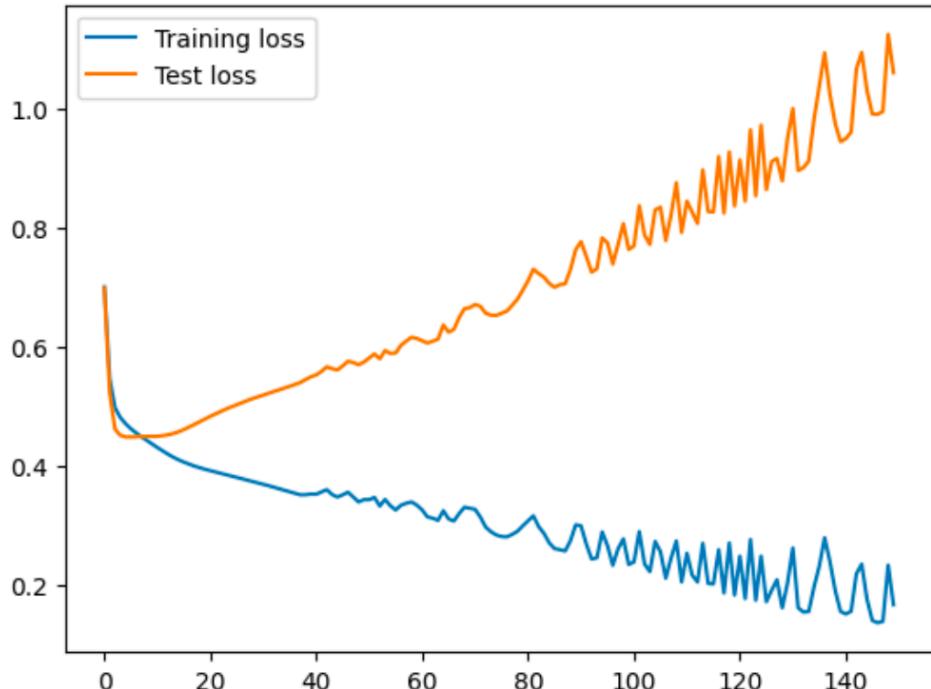




### Loss vs. training iteration



Loss vs. training iteration



```
patience = 5          # Number of steps to wait before stopping
steps_since_improvement = 0 # Steps since validation loss improved
min_loss = 1e8          # Minimum loss seen so far (start large)

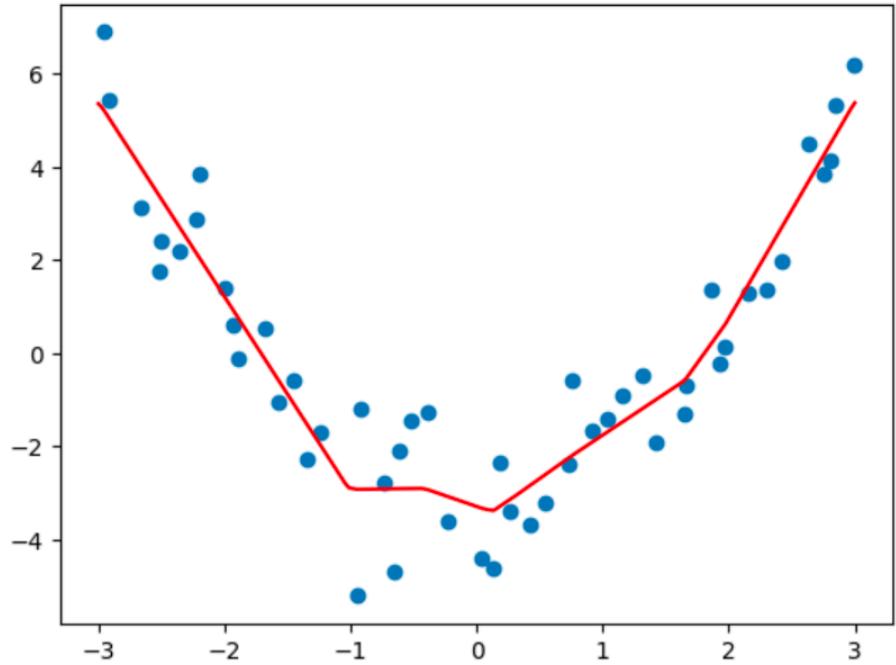
for i in range(steps):
    ...

    valid_loss = compute_loss(model, training_data)

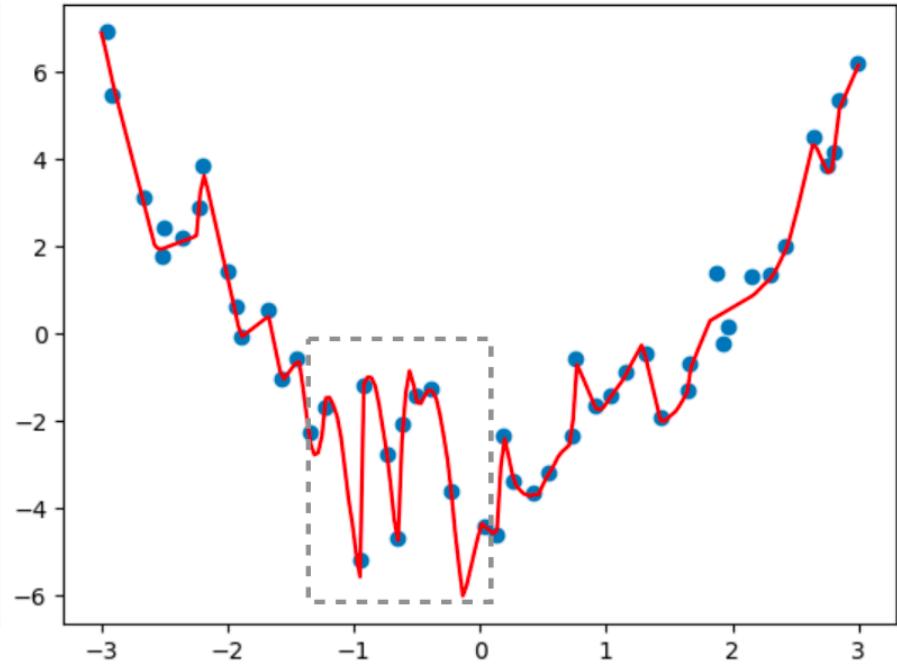
    # If the validation loss improves reset the counter
    if valid_loss < min_loss:
        steps_since_improvement = 0
        min_loss = valid_loss

    # Otherwise increment the counter
    else:
        steps_since_improvement += 1

    # If its been patience steps since the last improvement, stop
    if steps_since_improvement == patience:
        break
```

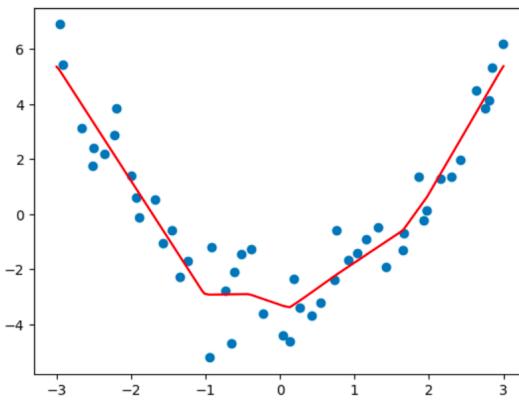


Well-fit

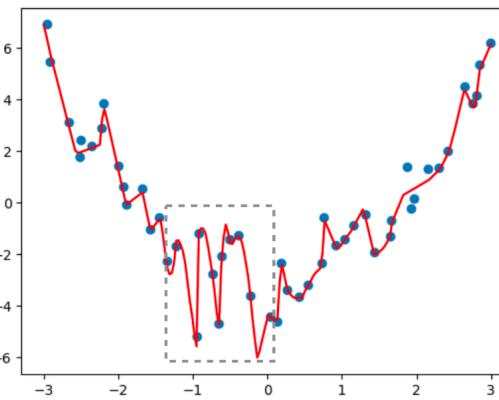


Overfit

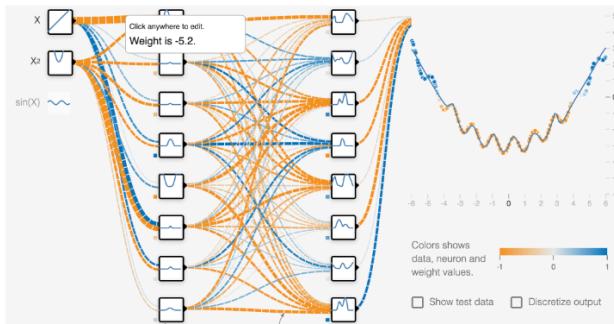
$$\text{MSE}(\mathbf{X}, \mathbf{y}, \mathbf{w}) = \frac{1}{N} \sum_{i=1}^N ((f(\mathbf{x}_i, \mathbf{w}) - y_i)^2)$$



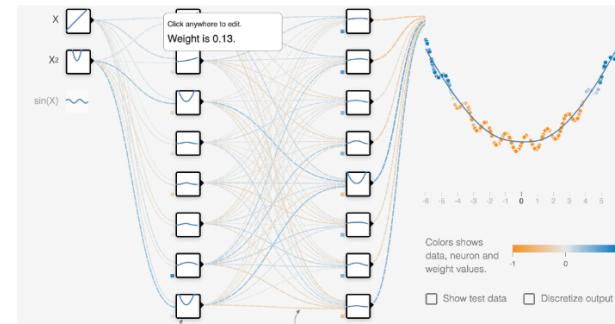
Well-fit



Overfit



An overfit network will have large weights to encode large slopes.



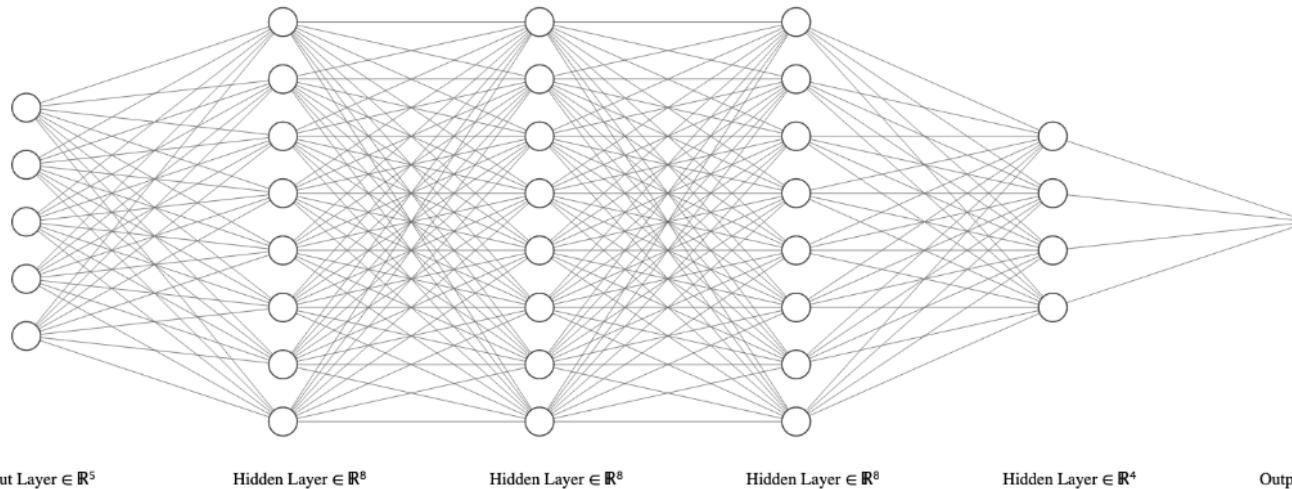
A regularized network will have smaller weights encoding a smooth function.

$$\mathbf{L}_2(\mathbf{w}) = \|\mathbf{w}\|_2^2 = \sum_{i=1}^d w_i^2$$

$$\mathbf{L}_2(\mathbf{W}) = \|\mathbf{W}\|_2^2 = \sum_{i=1}^d\sum_{j=1}^e w_{ij}^2$$

$$\textbf{Loss}(\mathbf{X}, \mathbf{y}, \mathbf{w}) = \textbf{MSE}(\mathbf{X}, \mathbf{y}, \mathbf{w}) + \lambda \mathbf{L}_2(\mathbf{w})$$

$$f(\mathbf{x}, \mathbf{w}_0, \dots) = \sigma(\sigma(\sigma(\sigma(\mathbf{x}^T \mathbf{W}_4)^T \mathbf{W}_3)^T \mathbf{W}_2)^T \mathbf{W}_1)^T \mathbf{w}_0$$



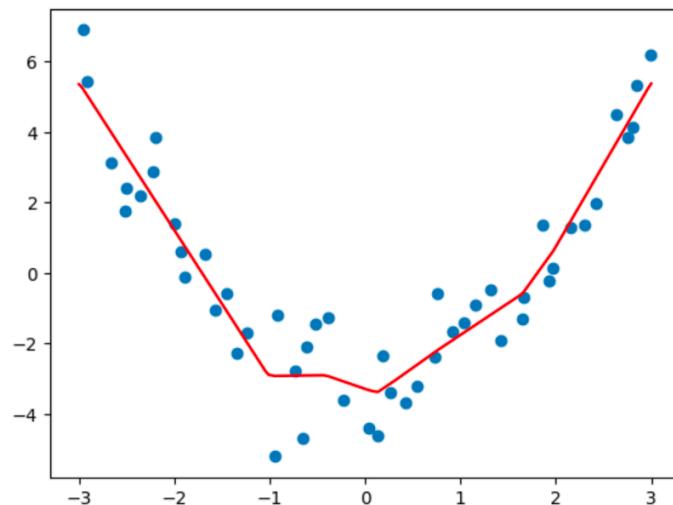
$$\mathbf{L}_2(\mathbf{w}_0, \mathbf{W}_1, \dots, \mathbf{W}_L) = \sum_{l=0}^L \|\mathbf{W}_l\|_2^2$$

In practice most networks also incorporate bias terms, so each linear function in our network can be written as:

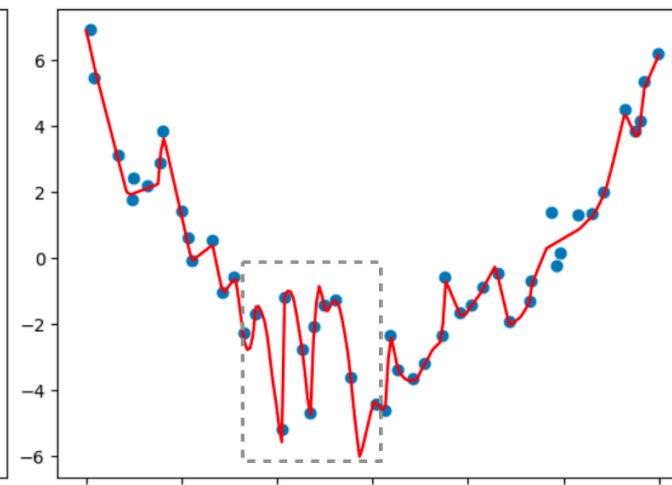
$$\mathbf{x}^T \mathbf{W} + \mathbf{b}$$

And the full prediction function for a sigmoid-activation network might be:

$$f(\mathbf{x}, \mathbf{w}_0, \dots) = \sigma(\sigma(\sigma(\sigma(\mathbf{x}^T \mathbf{W}_4 + \mathbf{b}_4)^T \mathbf{W}_3 + \mathbf{b}_3)^T \mathbf{W}_2 + \mathbf{b}_2)^T \mathbf{W}_1 + \mathbf{b}_1)^T \mathbf{w}_0 + \mathbf{b}_0$$

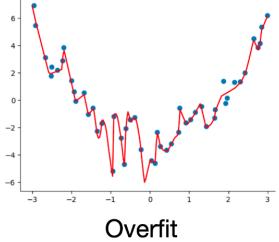
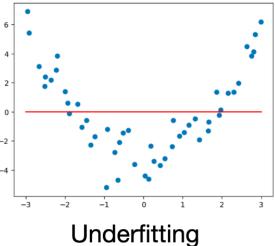
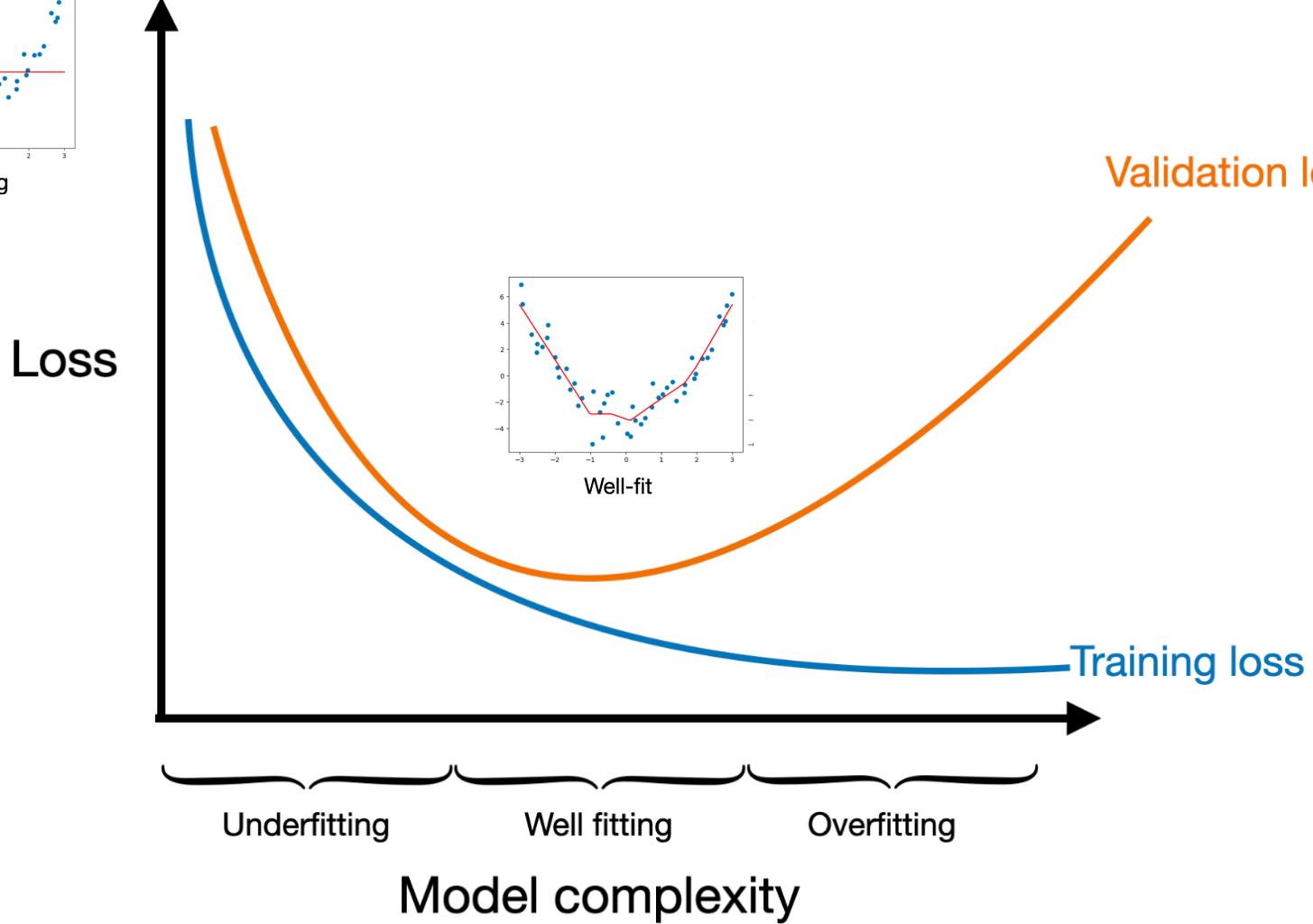


Well-fit

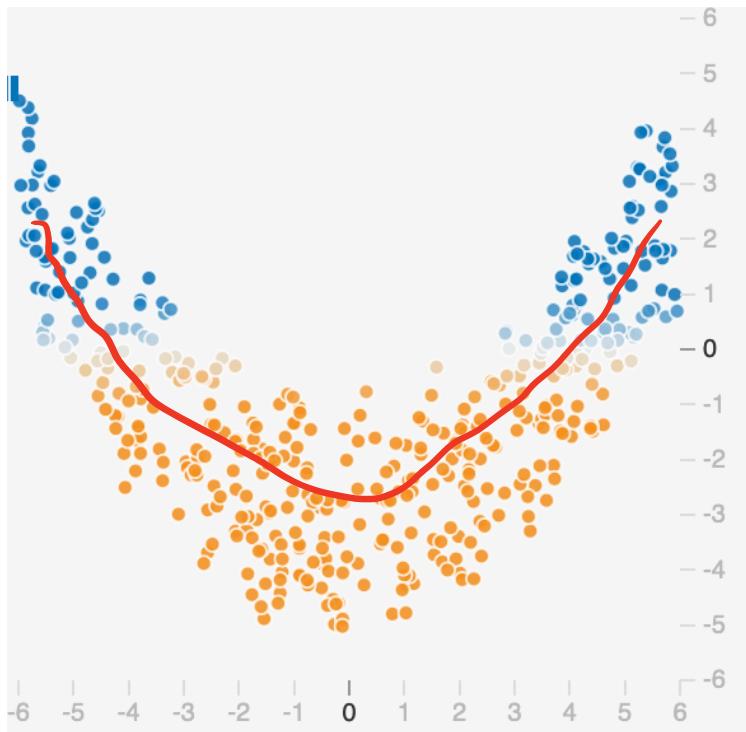


Overfit

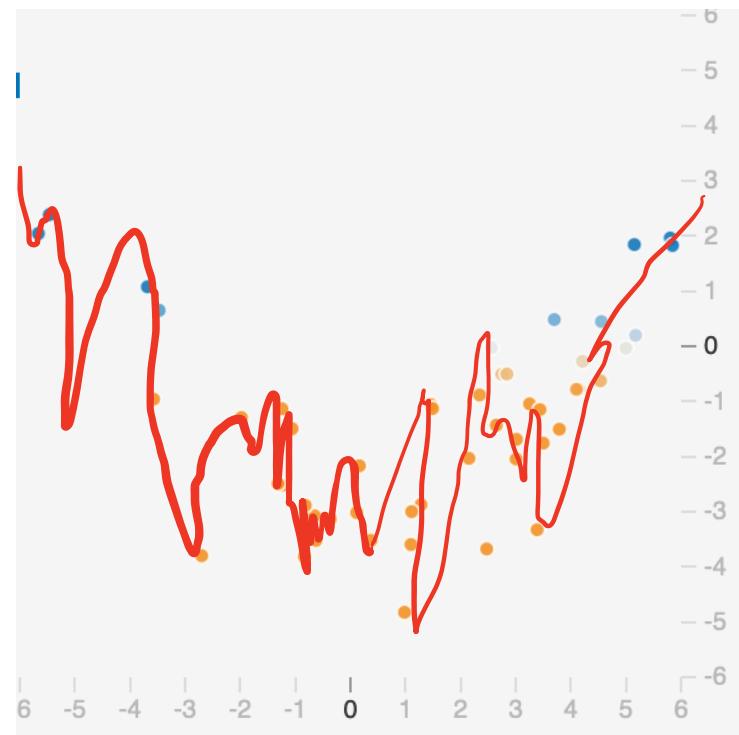
Does the bias contribute to overfitting?



# Sampling data

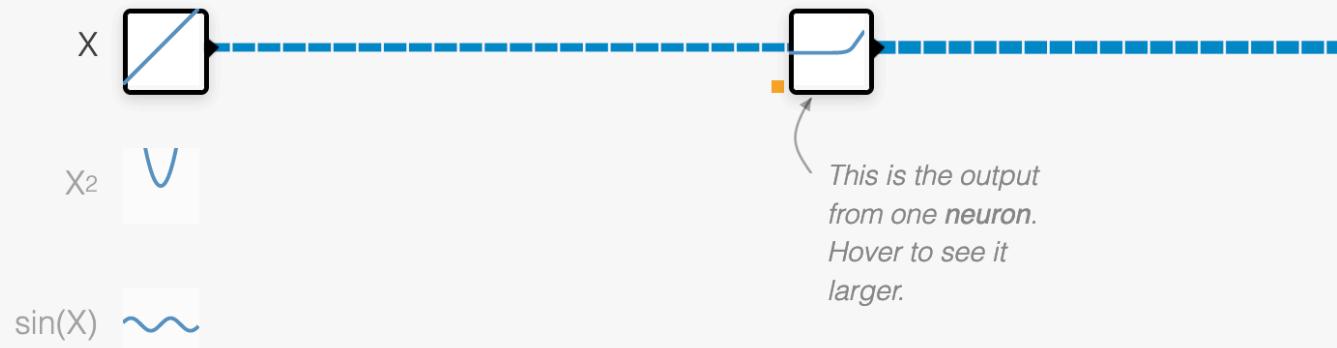


Data distribution

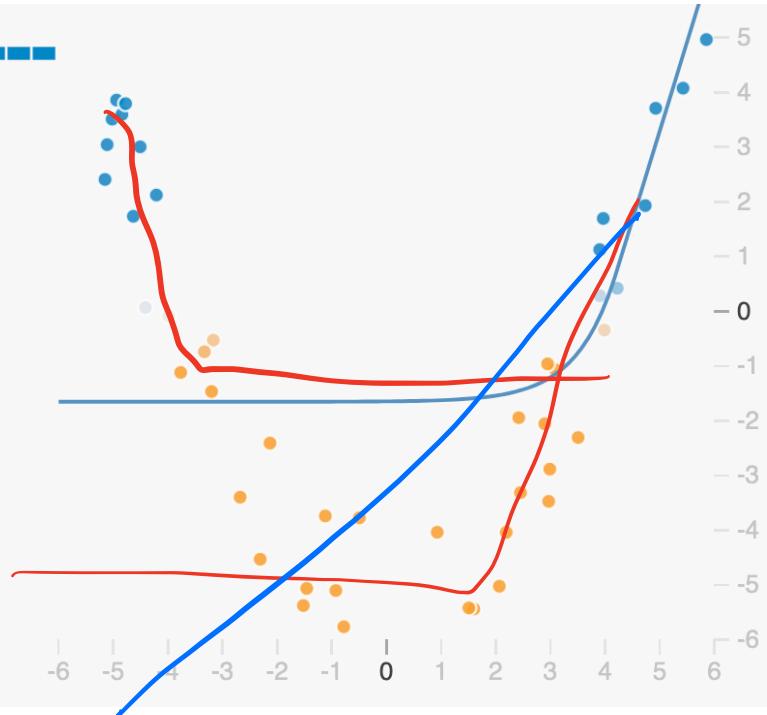


Training data

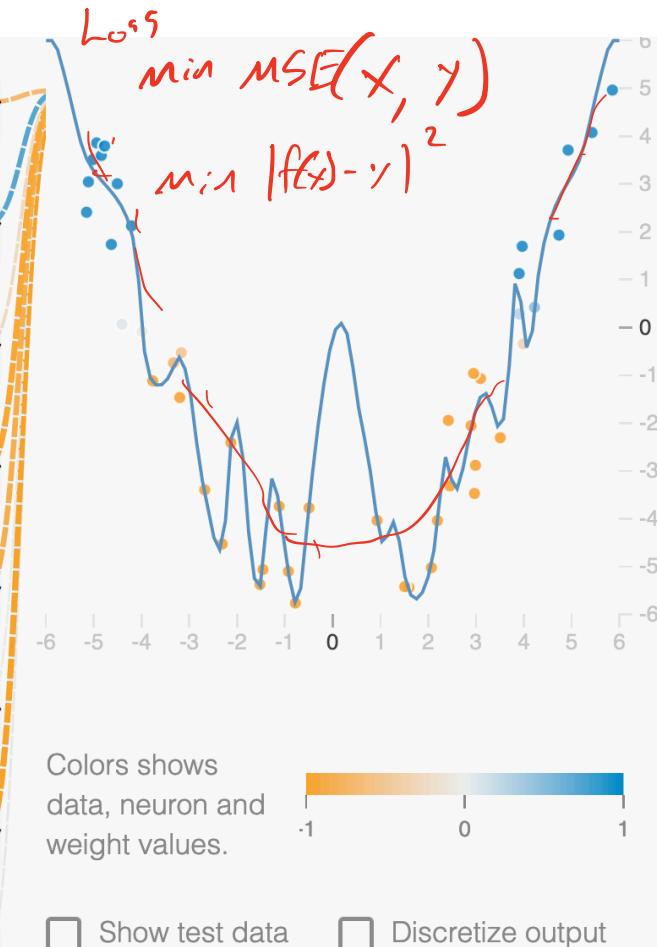
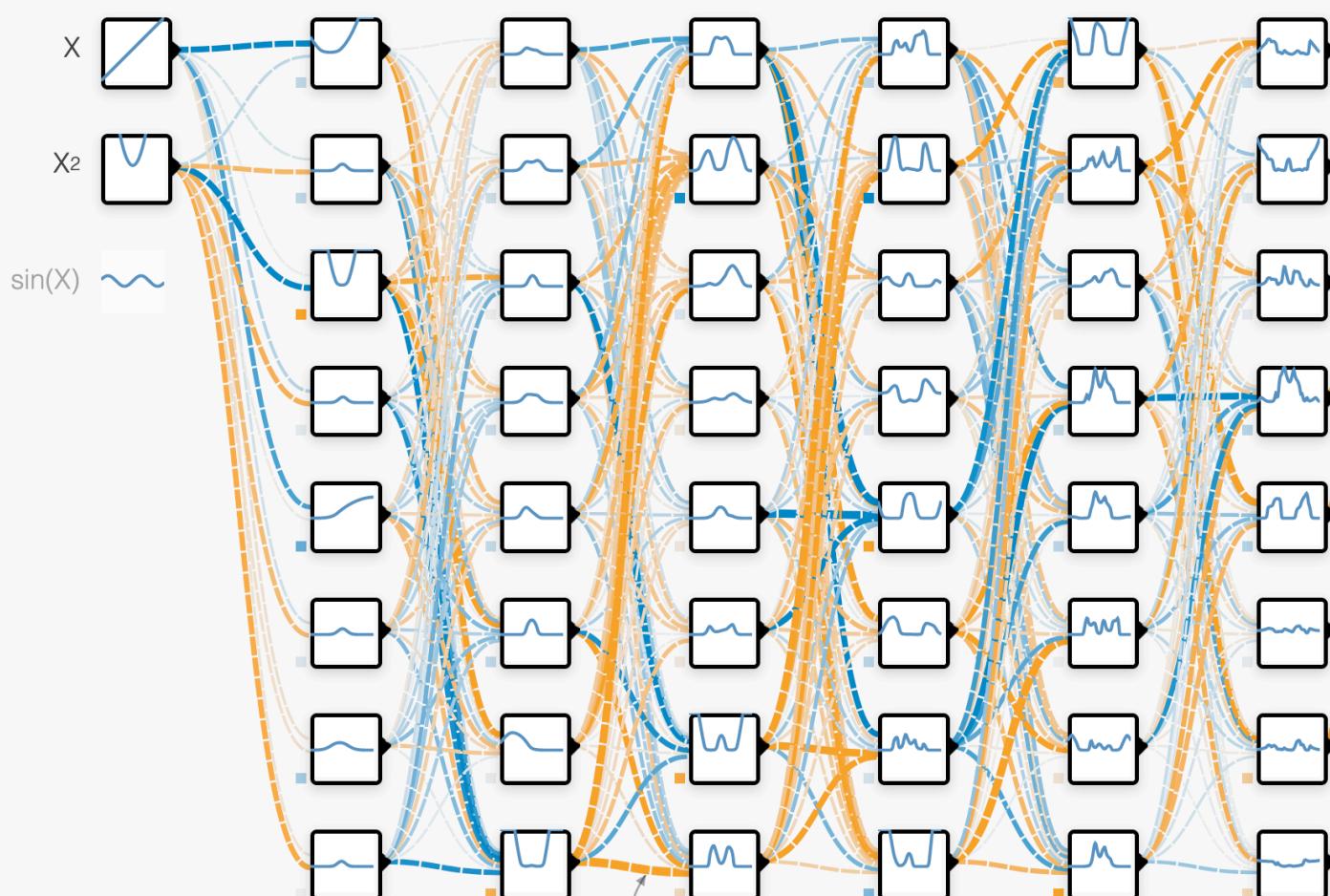
# Underfitting



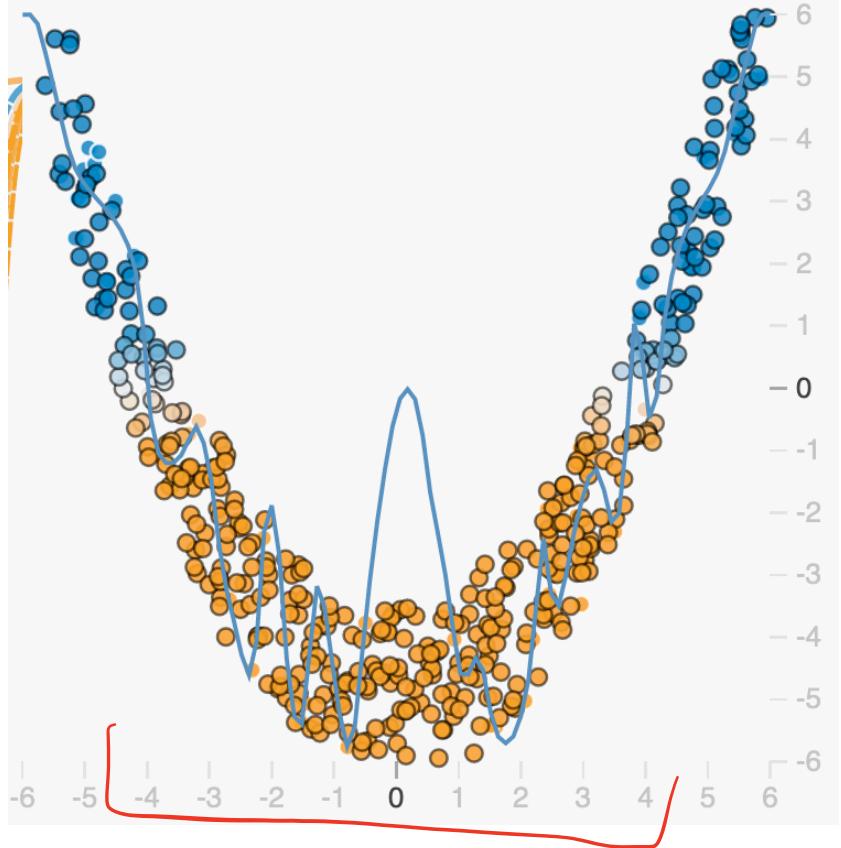
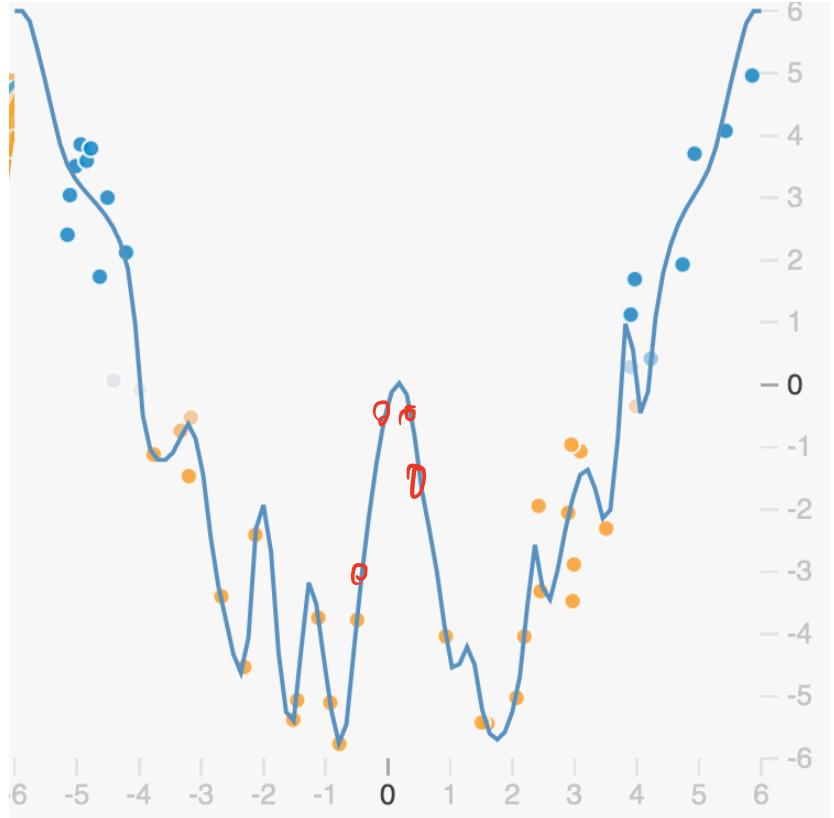
Model is too simple  
to capture the actual  
variations in data



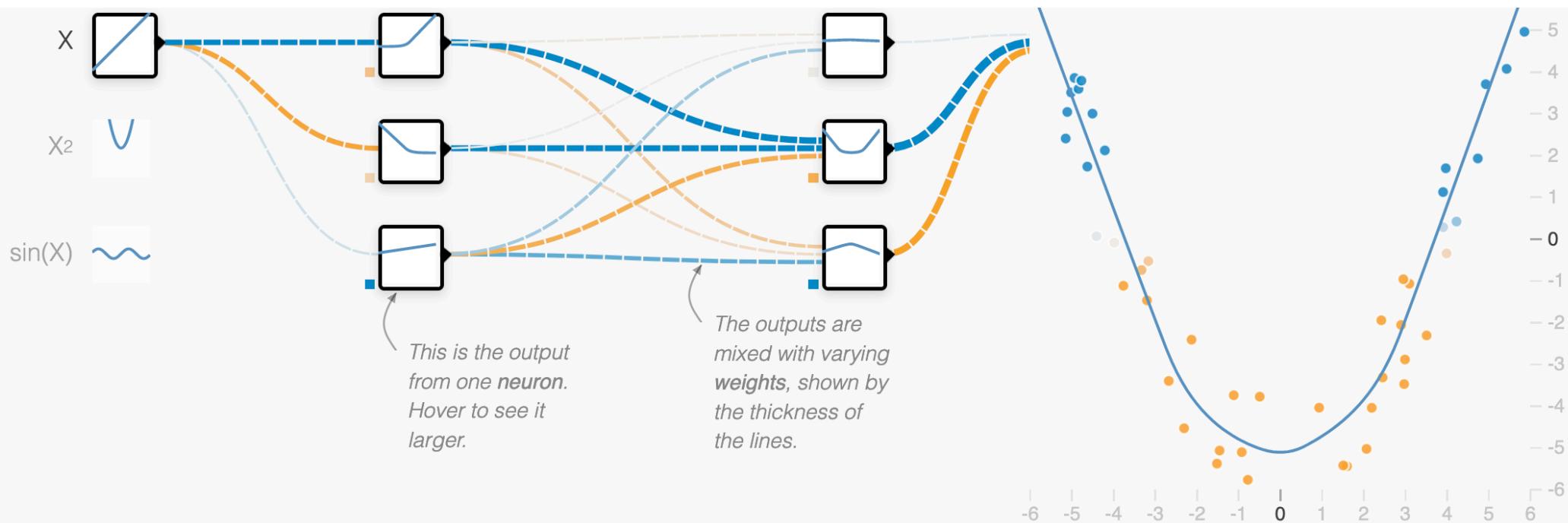
# Overfitting



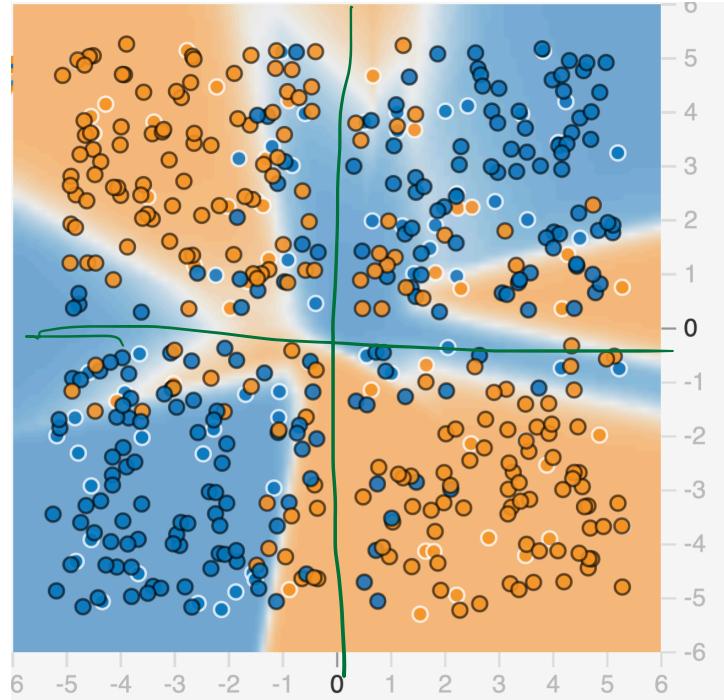
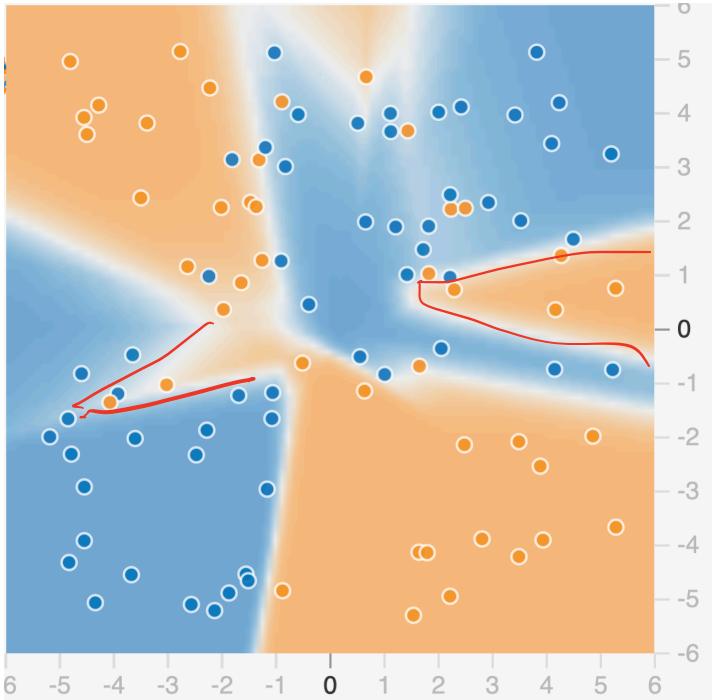
# Overfitting



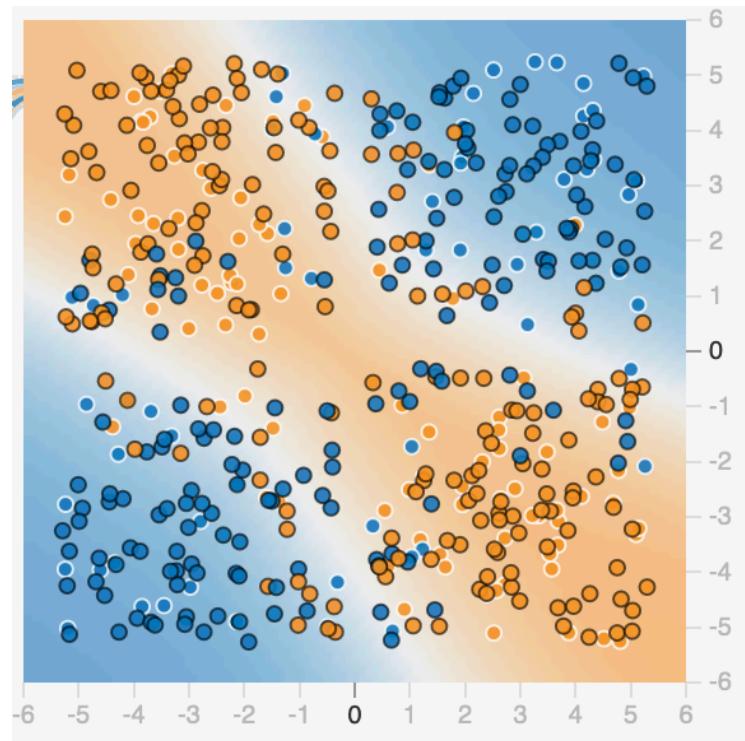
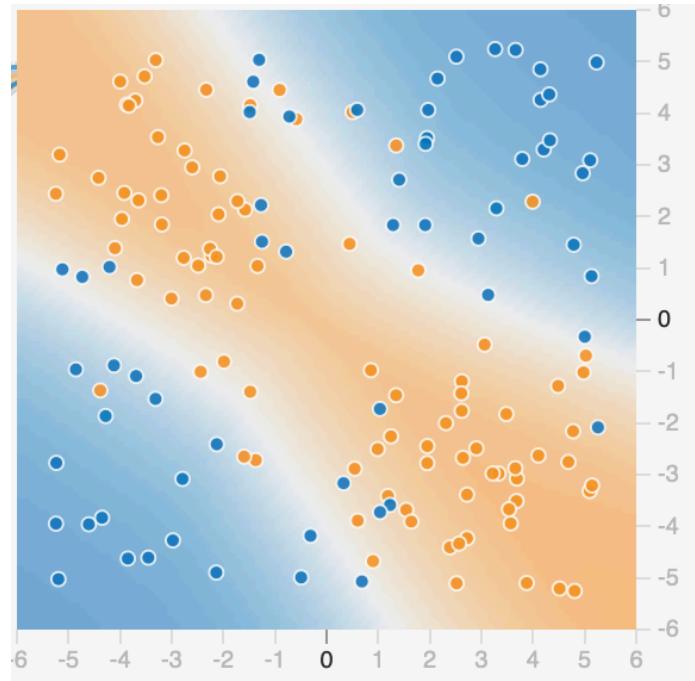
# Good fit

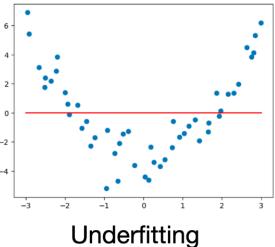
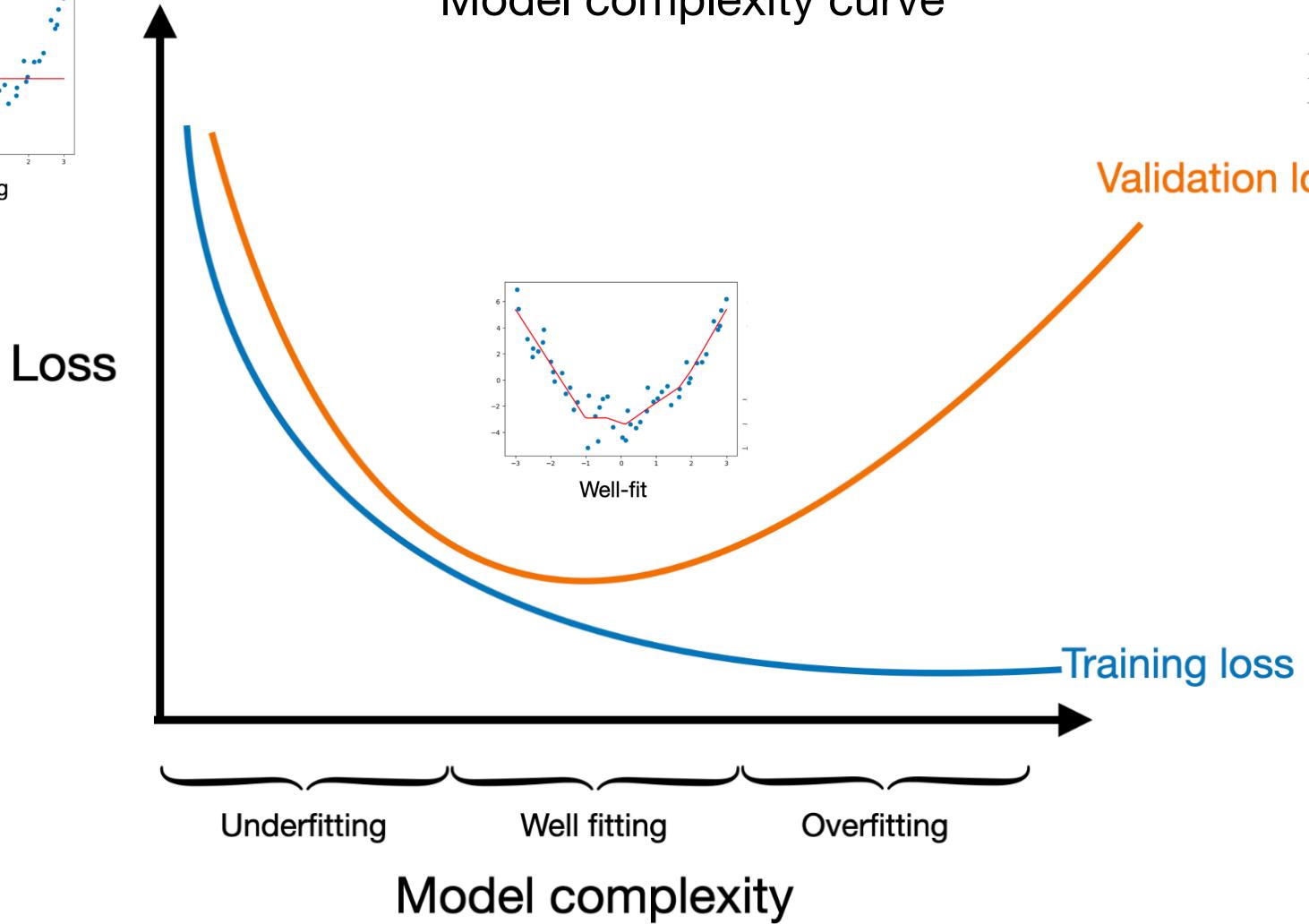


# Overfitting

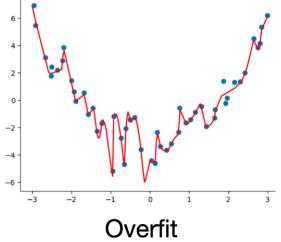


## Better fit



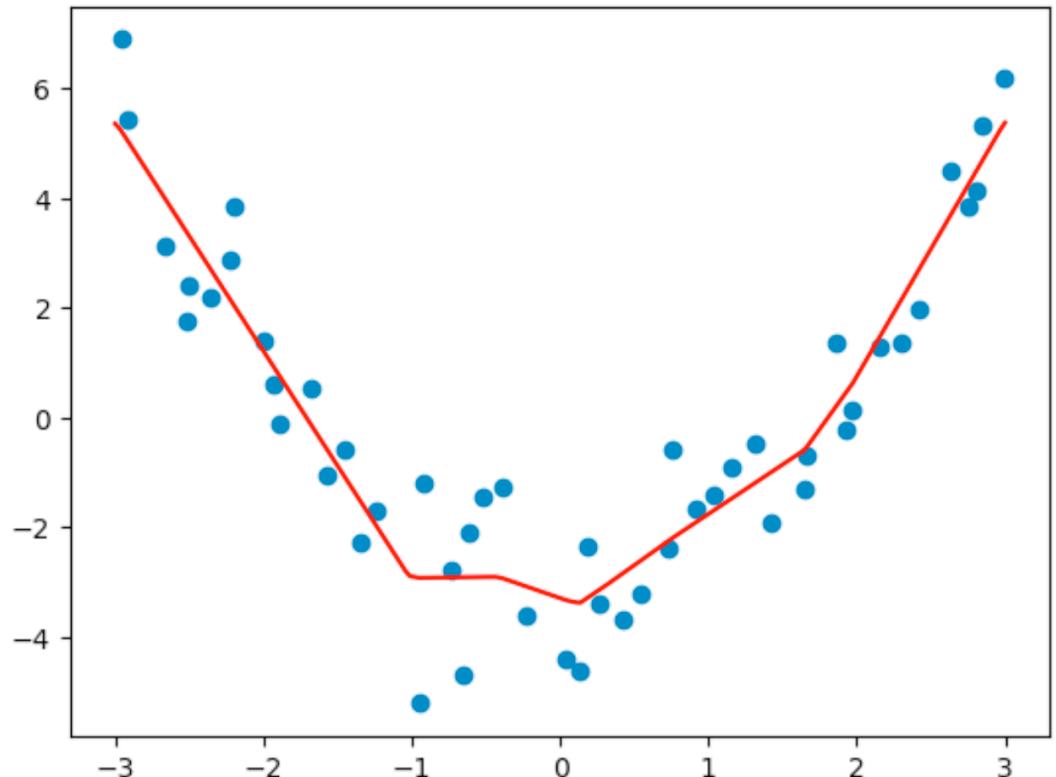


Model complexity curve

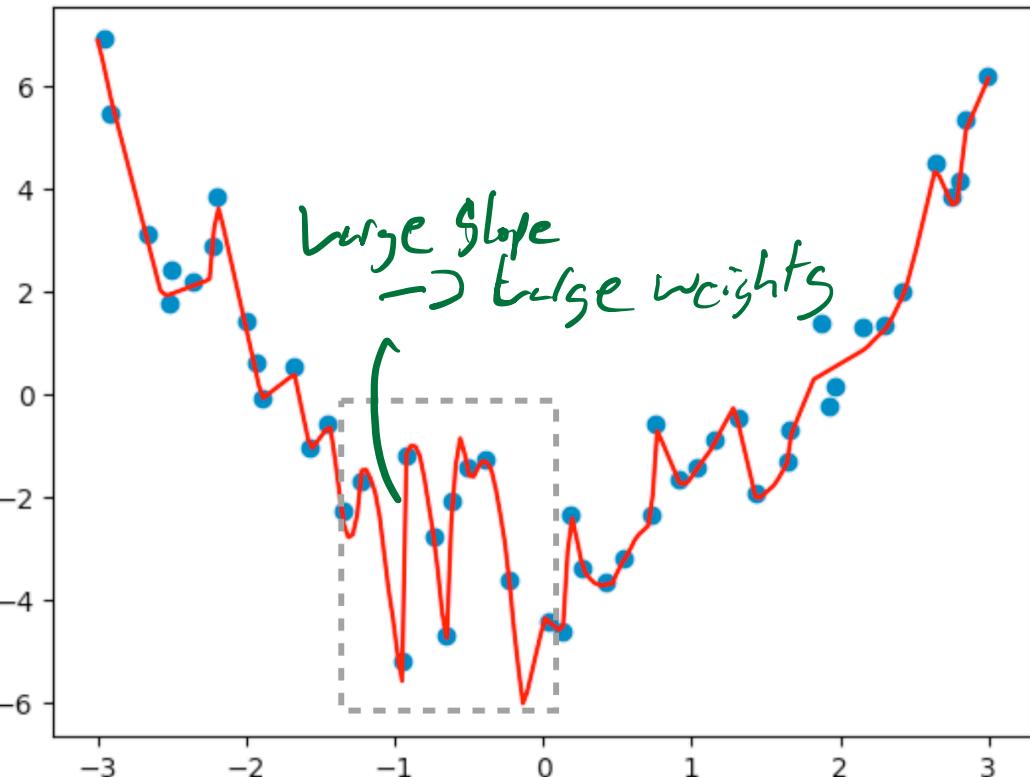


	Underfitting	Overfitting	Good fit
Training loss:	high	Low	Low
Test loss:	high	high	Low
→ Make our network larger or more complex			

## Overfitting with high weights

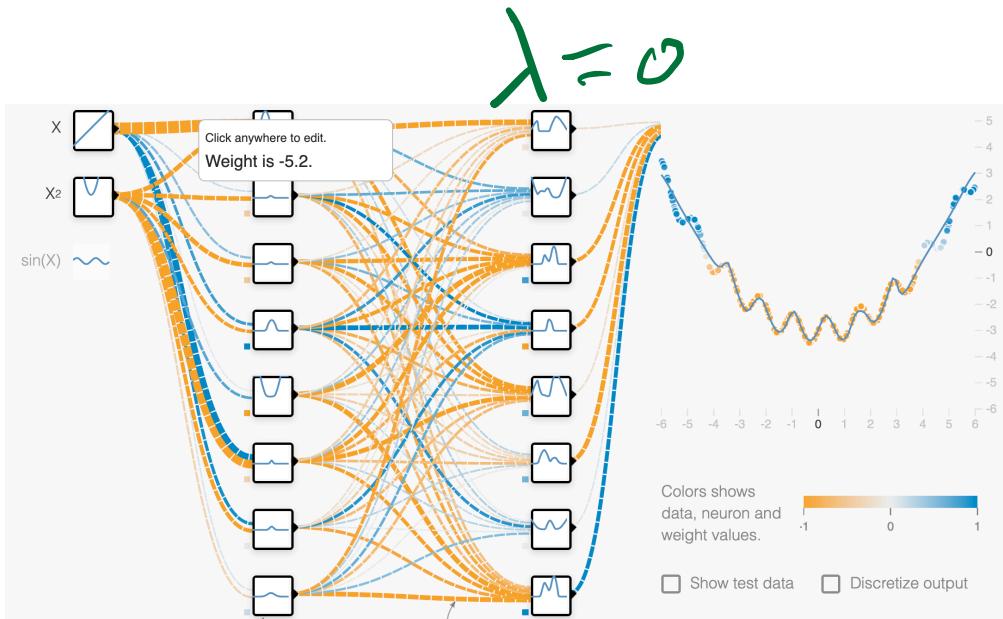


Well-fit



Overfit

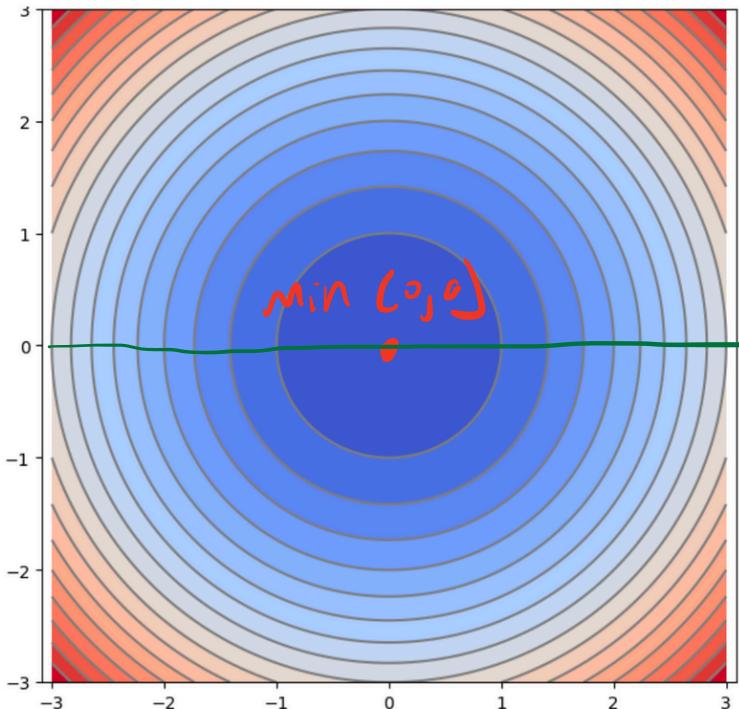
# Overfitting with high weights



## L2 Regularization

L2-Loss

$$\ell^2 = w_1^2 + w_2^2$$

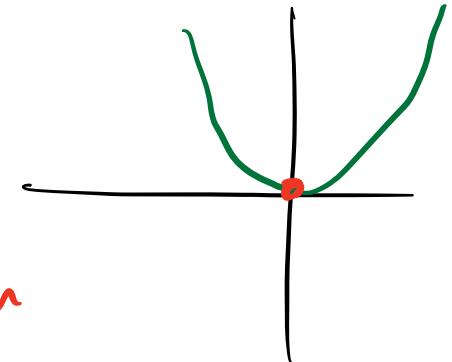


Squared L2-Norm

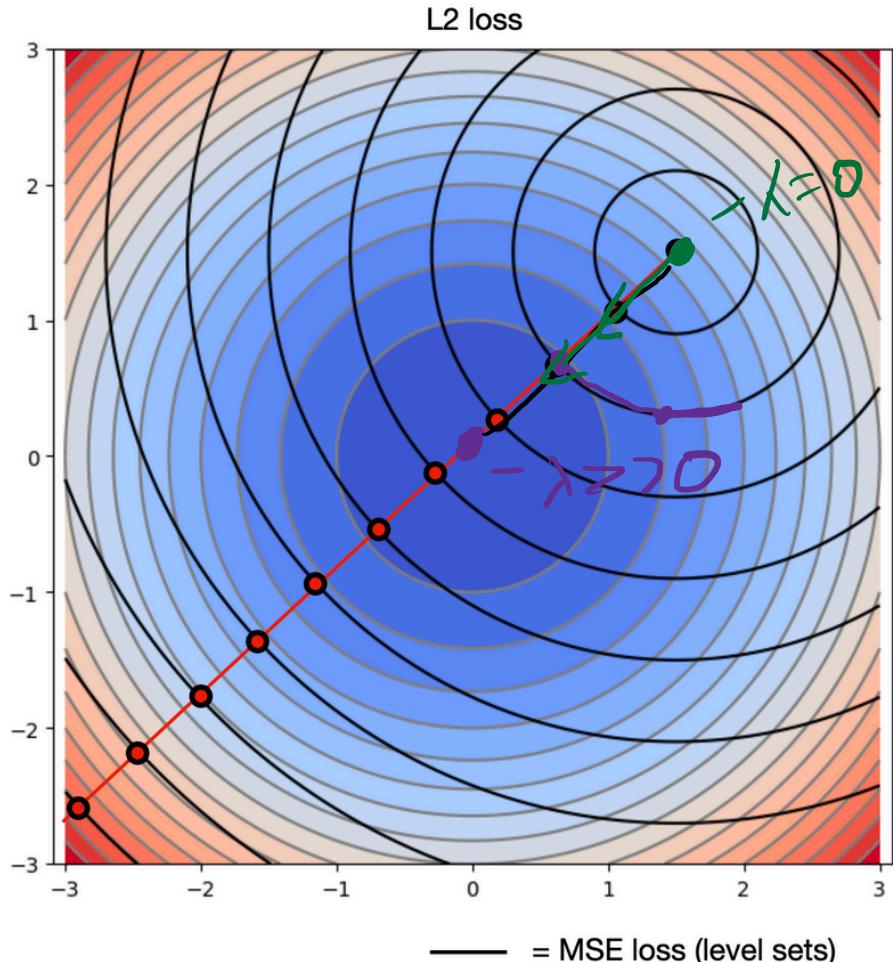
Min  $\sum_{i=1}^d \sum_{j=1}^e w_{ij}^2$

minimized where  $w_{ij} = 0$

$$\forall i, j$$



## L2 Regularization



$$L_2(\mathbf{W}) = \|\mathbf{W}\|_2^2 = \sum_{i=1}^d \sum_{j=1}^e w_{ij}^2$$

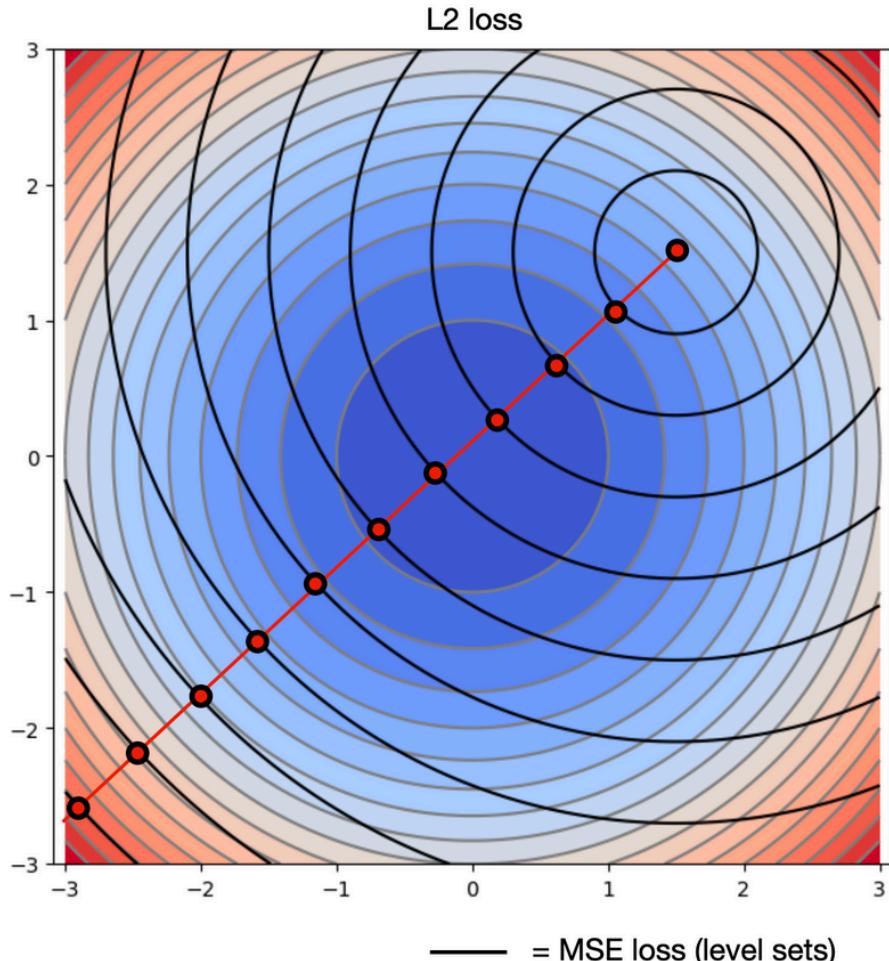
Loss( $\mathbf{X}, \mathbf{y}, \mathbf{w}$ ) = MSE( $\mathbf{X}, \mathbf{y}, \mathbf{w}$ ) +  $\lambda L_2(\mathbf{w})$

Scaling parameter  $\lambda$

λ makes good pred

has small weights

# L2 Regularization



$$\mathbf{L}_2(\mathbf{W}) = \|\mathbf{W}\|_2^2 = \sum_{i=1}^d \sum_{j=1}^e w_{ij}^2$$

$$\text{Loss}(\mathbf{X}, \mathbf{y}, \mathbf{w}) = \text{MSE}(\mathbf{X}, \mathbf{y}, \mathbf{w}) + \lambda \mathbf{L}_2(\mathbf{w})$$

```
from torch import optim  
optimizer = optim.SGD(model.parameters(), lr=0.1, weight_decay=0.01)
```

λ

# L1 Regularization

L2

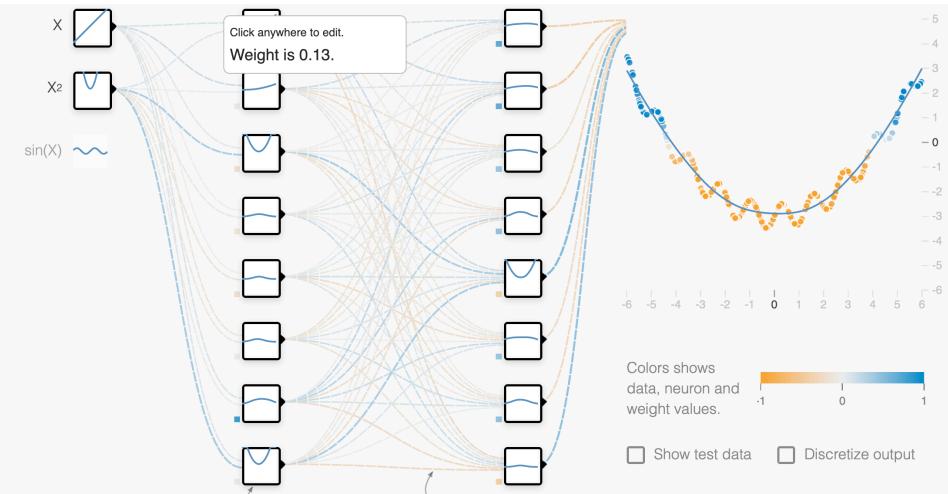
Vector:  $\mathbf{L}_2(\mathbf{w}) = \|\mathbf{w}\|_2^2 = \sum_{i=1}^d w_i^2$       Matrix:  $\mathbf{L}_2(\mathbf{W}) = \|\mathbf{W}\|_2^2 = \sum_{i=1}^d \sum_{j=1}^e w_{ij}^2$

L1

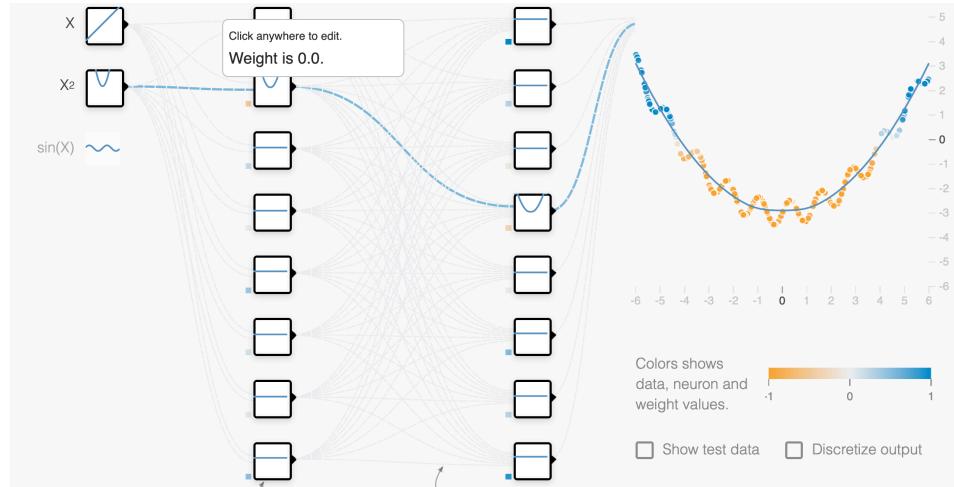
Vector:  $\mathbf{L}_1(\mathbf{w}) = \|\mathbf{w}\|_1 = \sum_{i=1}^d |w_i|$ ,      Matrix:  $\mathbf{L}_1(\mathbf{W}) = \|\mathbf{W}\|_1 = \sum_{i=1}^d \sum_{j=1}^e |w_{ij}|$

# L1 Regularization

L2



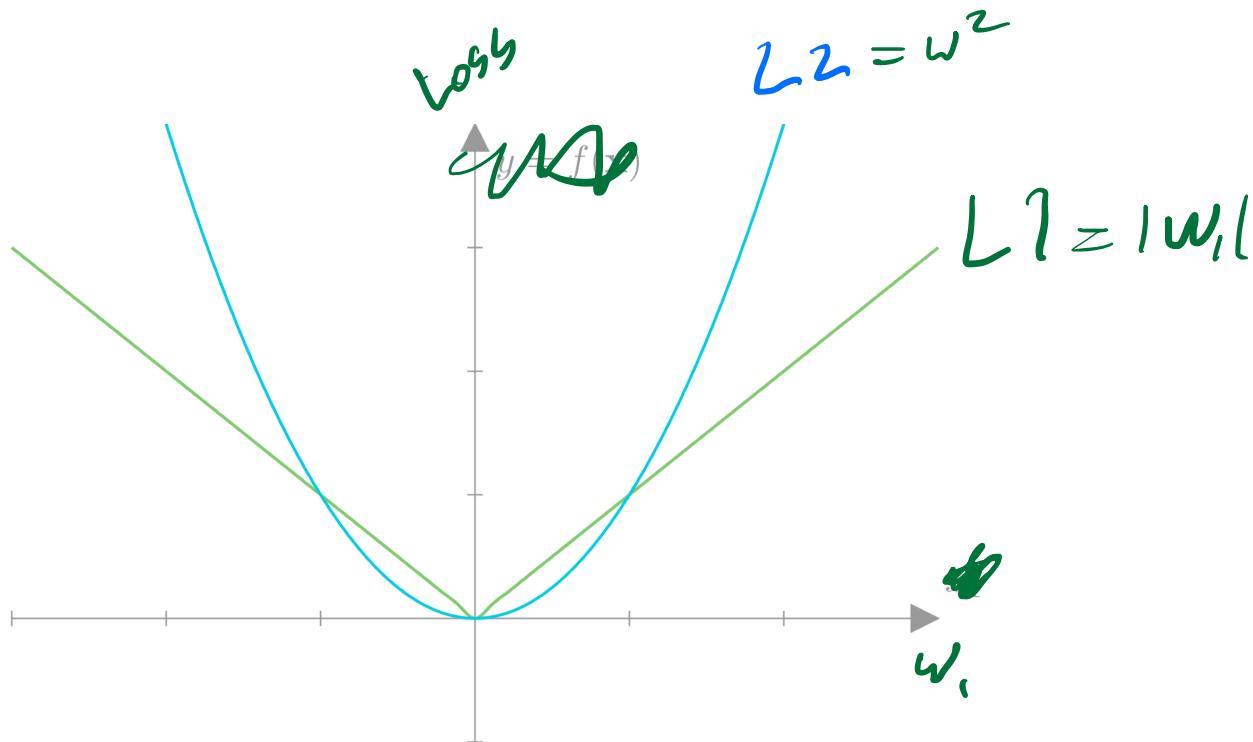
L1



Sparsity

most weights are 0

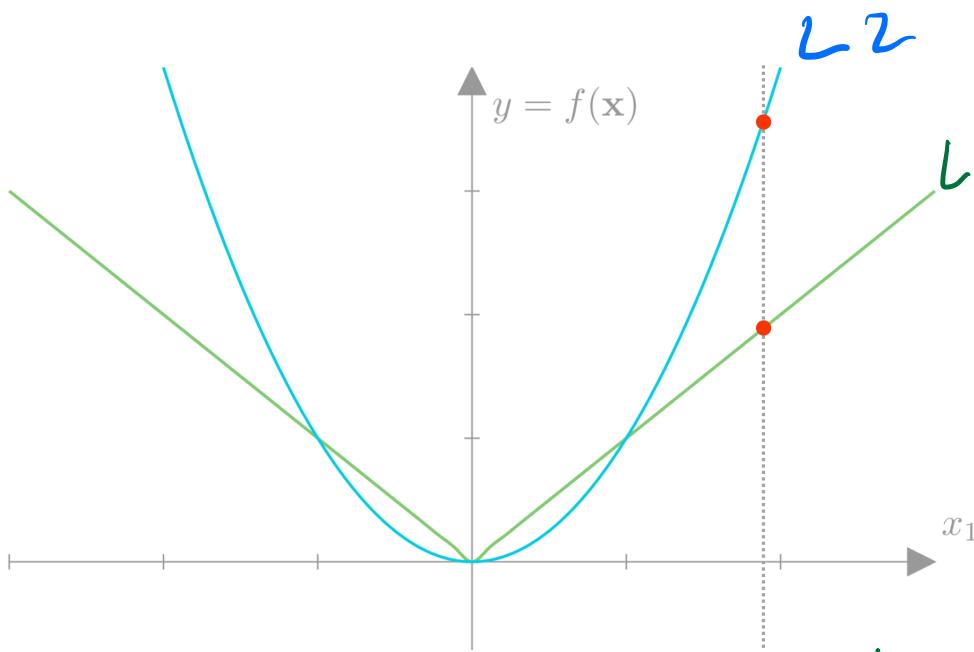
# L1 Vs. L2 Regularization



Vector:  $\mathbf{L}_2(\mathbf{w}) = \|\mathbf{w}\|_2^2 = \sum_{i=1}^d w_i^2$

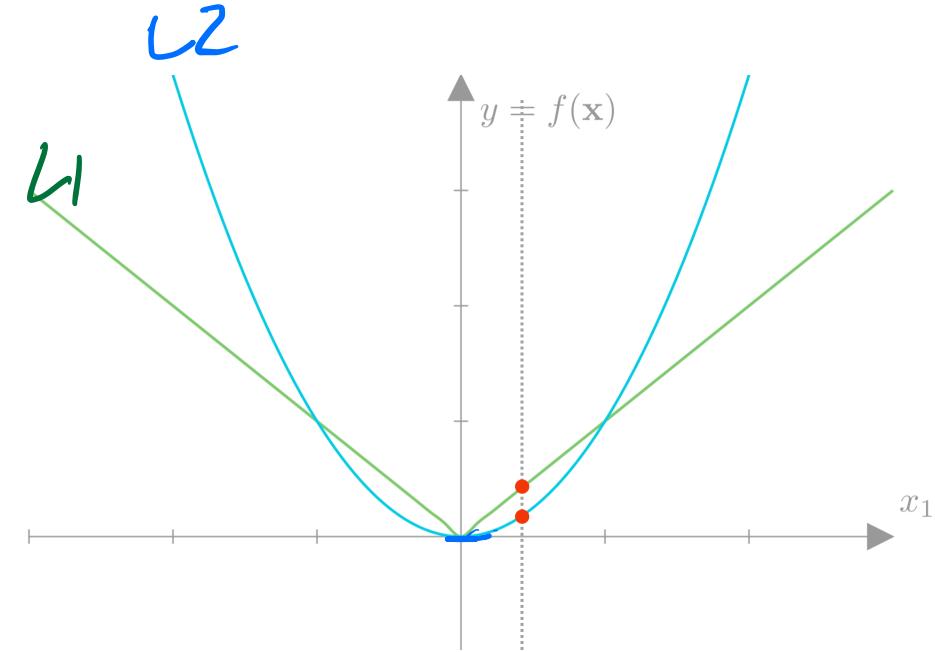
Vector:  $\mathbf{L}_1(\mathbf{w}) = \|\mathbf{w}\|_1 = \sum_{i=1}^d |w_i|$ ,

# L1 Vs. L2 Regularization

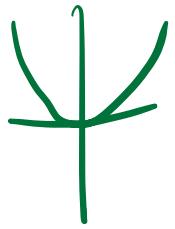


L2 Loss grows quickly  
away from 0

L1 loss grows consistently

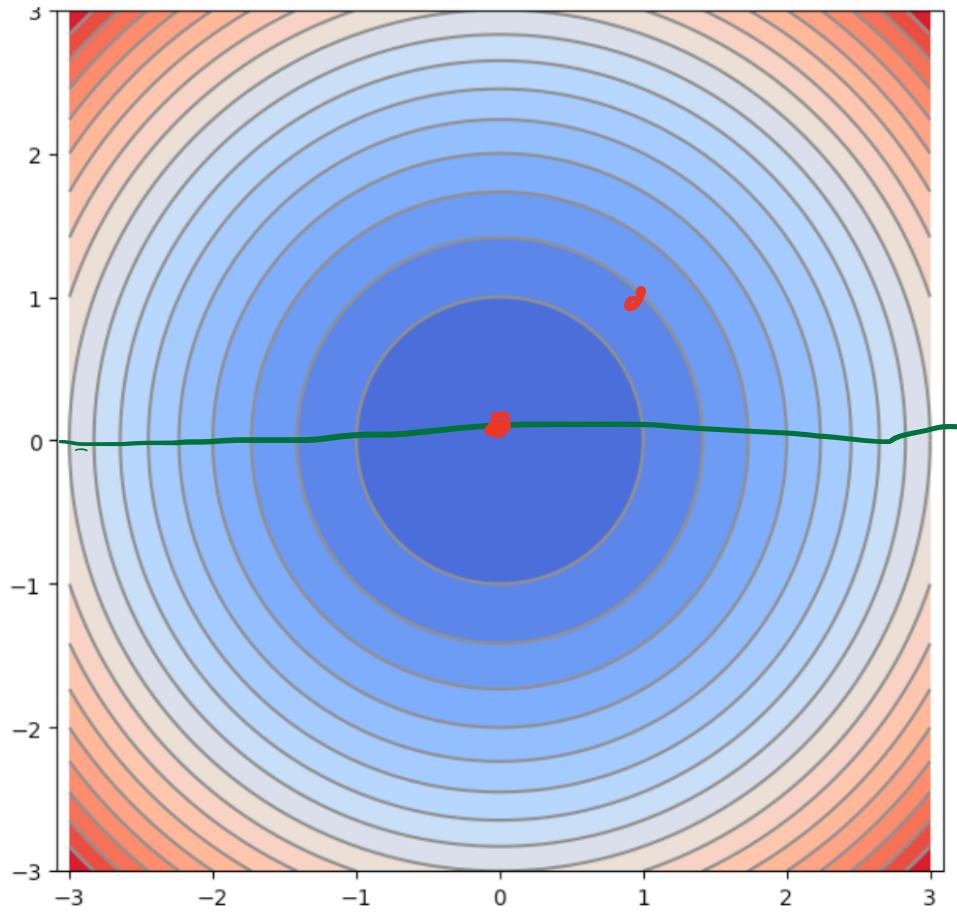


L2 Loss flat near 0  
L1 Loss remains constant



## L2-Loss

$$\ell^2 = w_1^2 + w_2^2$$

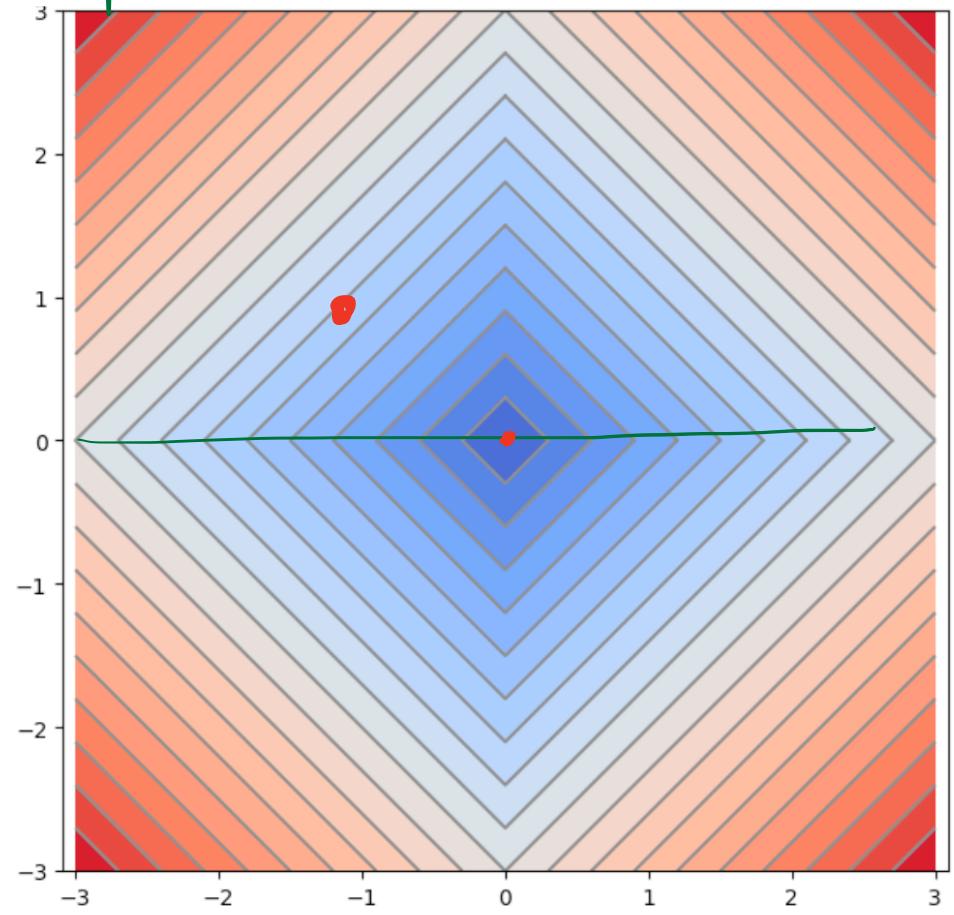


## L1 Vs. L2 Regularization

## L1-Loss



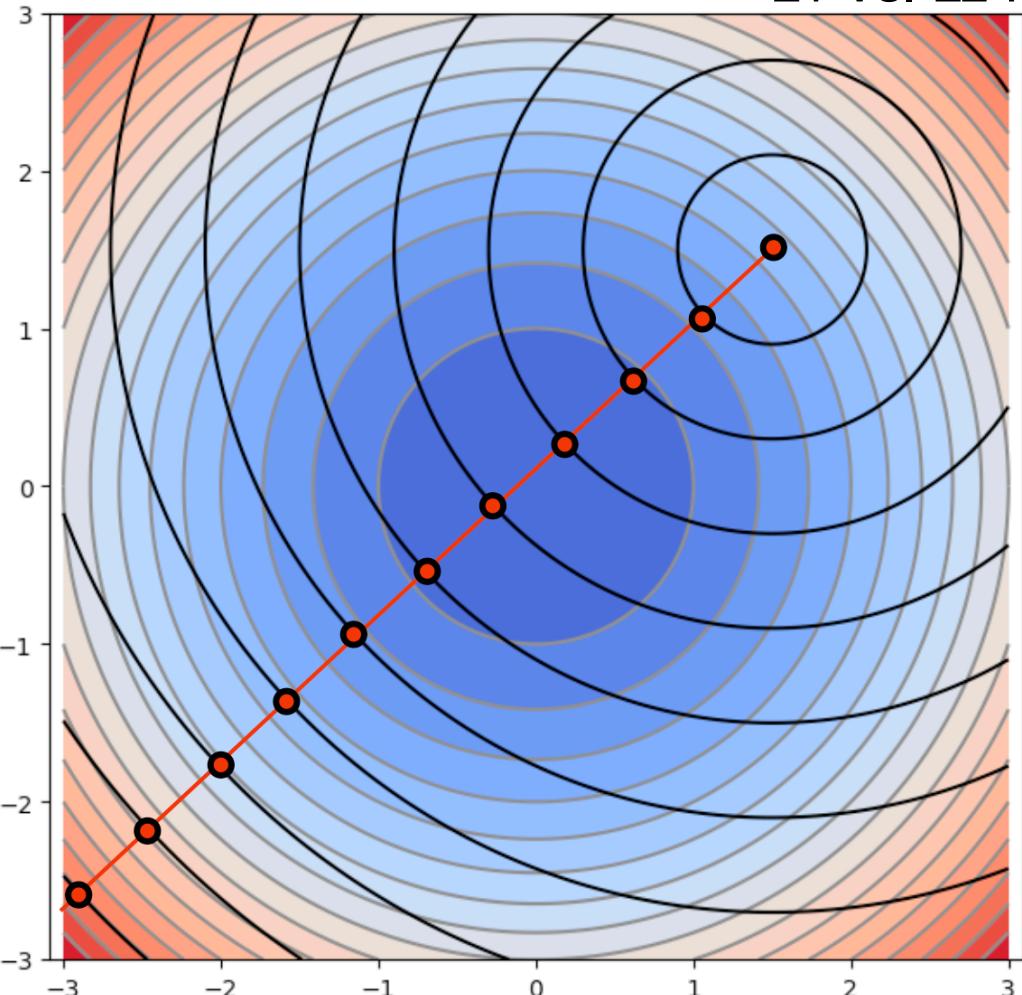
$$\ell^1 = |w_1| + |w_2|$$



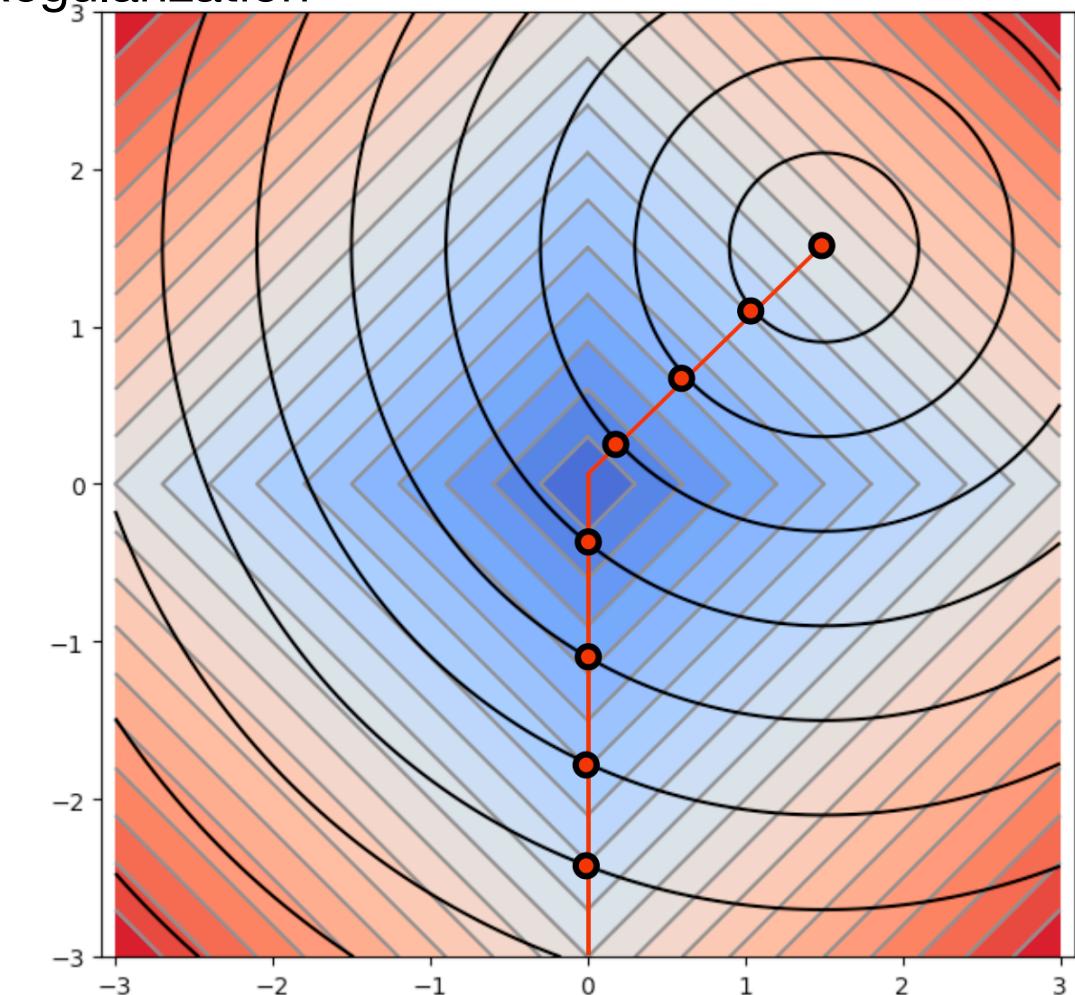
L2 loss

# L1 Vs. L2 Regularization

L1 loss

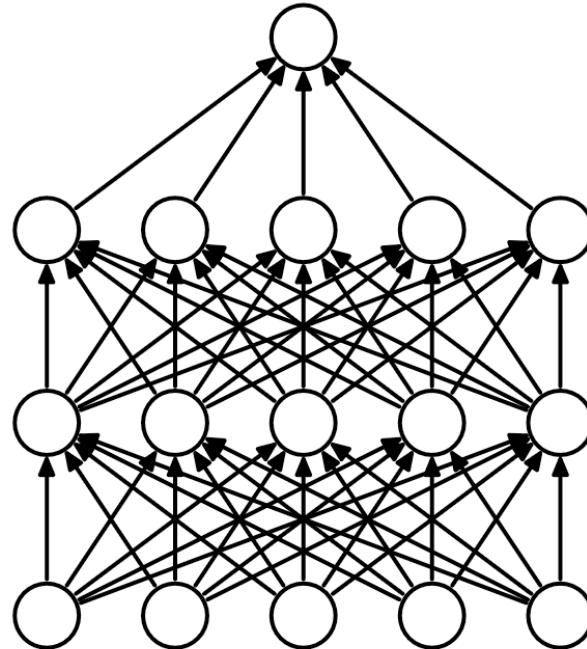


— = MSE loss (level sets)

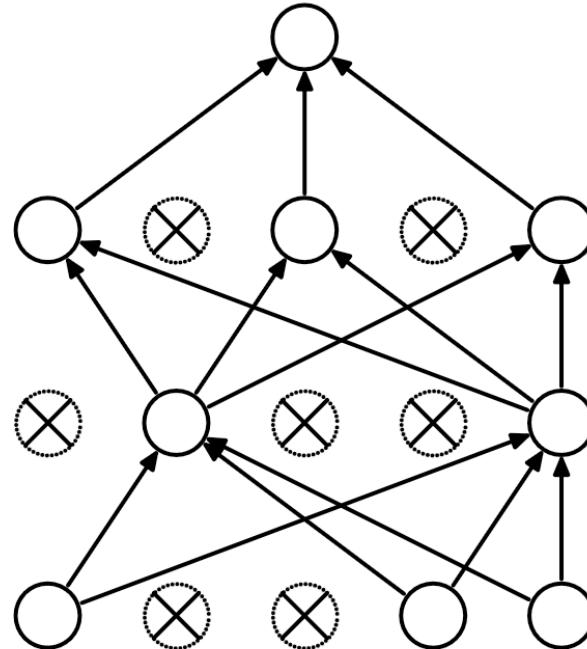


● = Minimum L2 or L1 loss for constant MSE

# Dropout



(a) Standard Neural Net



(b) After applying dropout.

Figure 1: Dropout Neural Net Model. **Left:** A standard neural net with 2 hidden layers. **Right:** An example of a thinned net produced by applying dropout to the network on the left. Crossed units have been dropped.

## Dropout

$$\text{Dropout}(\mathbf{X}, r) = \mathbf{D} \odot \mathbf{X}, \quad \mathbf{D} = \begin{bmatrix} d_{11} & d_{12} & \dots & d_{1n} \\ d_{21} & d_{22} & \dots & d_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ d_{m1} & d_{m2} & \dots & d_{mn} \end{bmatrix}, \quad d_{ij} \sim \text{Bernoulli}(1 - r)$$

$r$  = dropout rate      prob that we drop out layer  
on feature

In i-th layer

$$\phi(\mathbf{x}) = \sigma(\underline{\text{DO}_r(\mathbf{x})^T \mathbf{W}} + \mathbf{b})$$

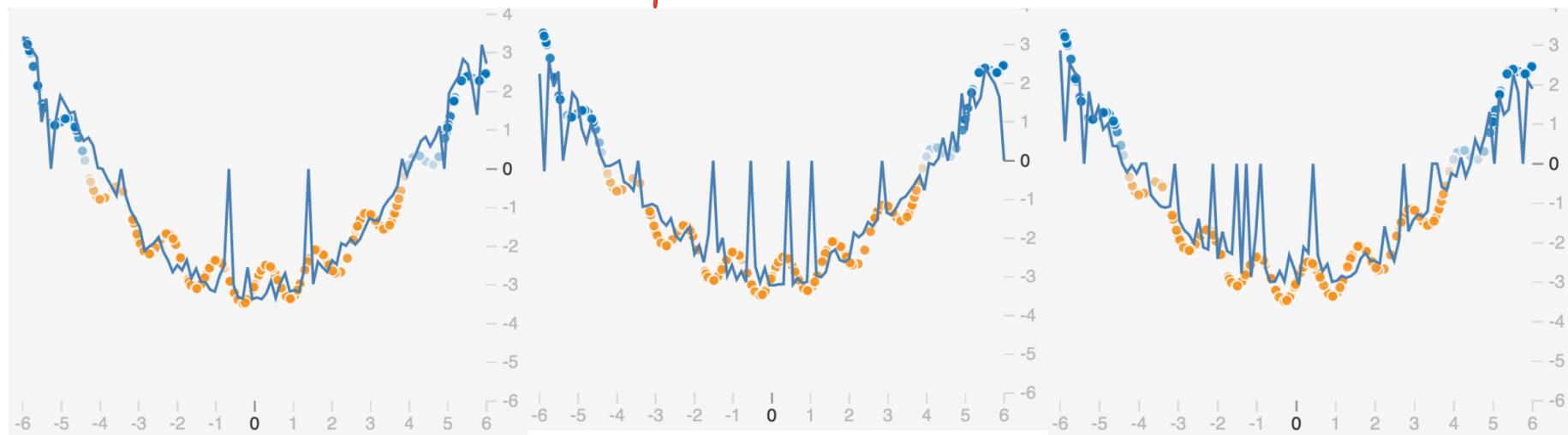
Or,

$$\phi(x) = \sigma(x^T w + b)$$

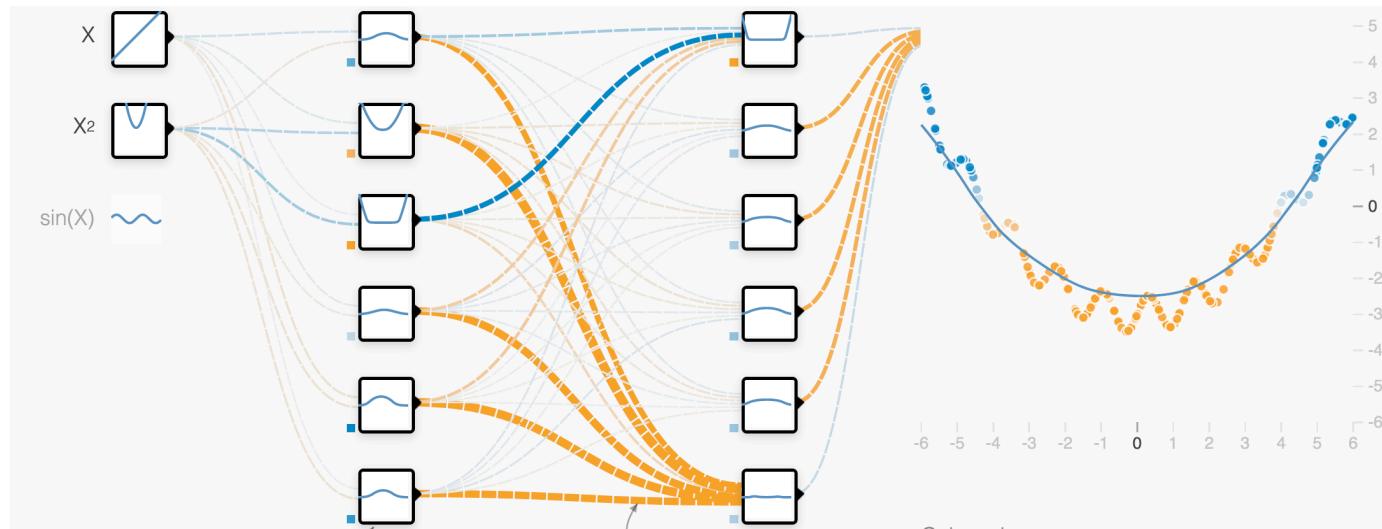
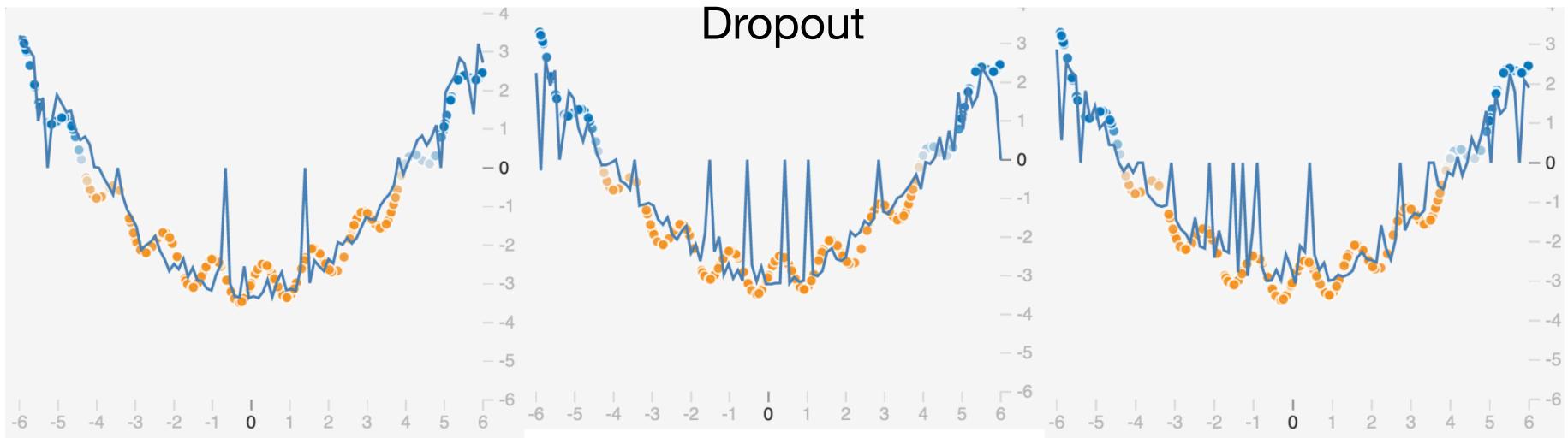
## Dropout

$$\text{Dropout}(\mathbf{X}, r) = \mathbf{D} \odot \mathbf{X}, \quad \mathbf{D} = \begin{bmatrix} d_{11} & d_{12} & \dots & d_{1n} \\ d_{21} & d_{22} & \dots & d_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ d_{m1} & d_{m2} & \dots & d_{mn} \end{bmatrix}, \quad d_{ij} \sim \text{Bernoulli}(1 - r)$$

Three random dropouts

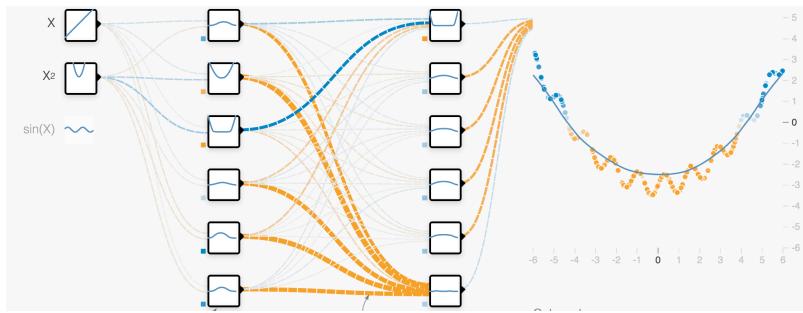
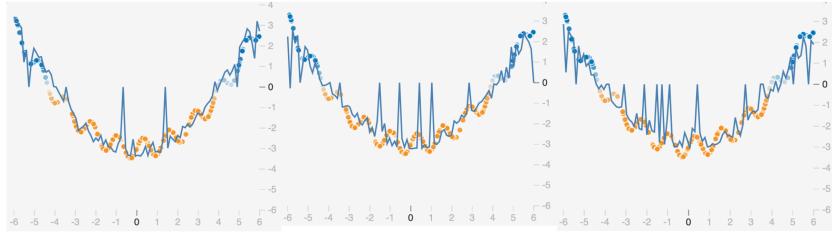


# Dropout



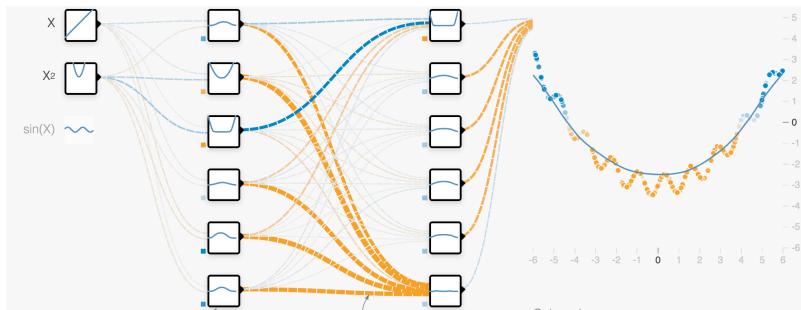
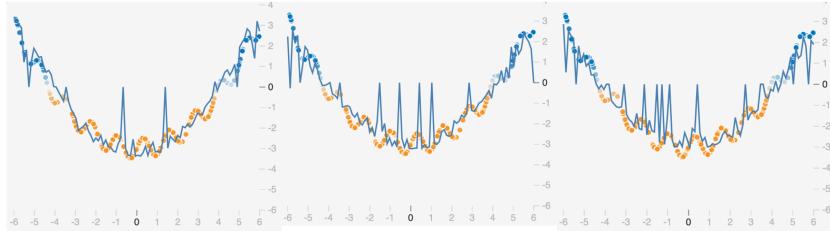
# Dropout

$$\phi(\mathbf{x})_{train} = \sigma(\text{DO}_r(\mathbf{x})^T \mathbf{W} + \mathbf{b}) \quad \rightarrow \quad \phi(\mathbf{x})_{eval} = \sigma(\mathbf{x}^T \mathbf{W} + \mathbf{b})$$



# Dropout

$$\phi(\mathbf{x})_{train} = \sigma(\text{DO}_r(\mathbf{x})^T \mathbf{W} + \mathbf{b}) \quad \rightarrow \quad \phi(\mathbf{x})_{eval} = \sigma(\mathbf{x}^T \mathbf{W} + \mathbf{b})$$



$$\mathbb{E}[\text{DO}_r(\mathbf{x})^T \mathbf{w}] = \sum_i d_i x_i w_i, \quad d_i \sim \text{Bernoulli}(1 - r)$$

$$= \sum_i p(d_i = 1) x_i w_i = (1 - r) \sum_i x_i w_i < \sum_i x_i w_i$$

## Dropout