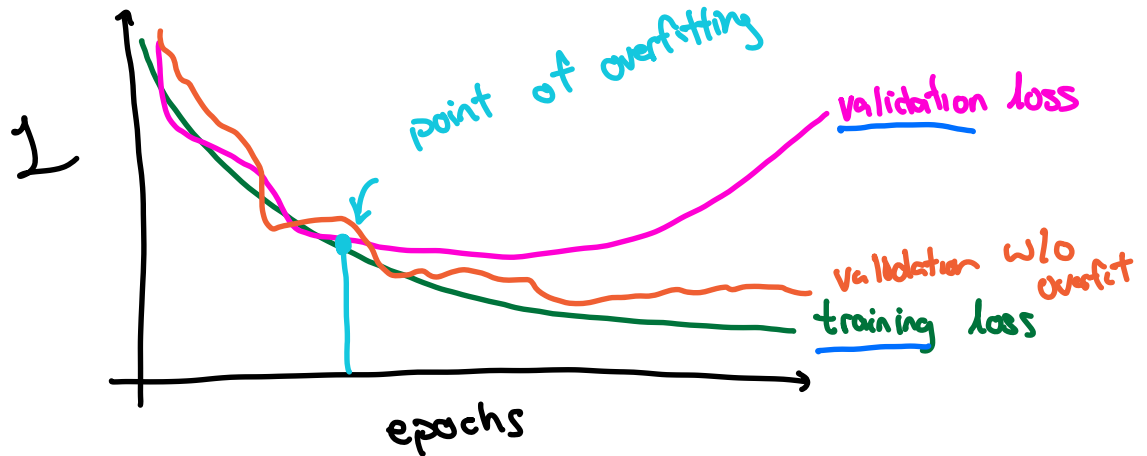


# Overfitting

When your model memorizes the training data.  
(Your model does not generalize.)



Take MNIST as an example.



60,000 training images → each input image (or subset) is recognized by a handful of parameters

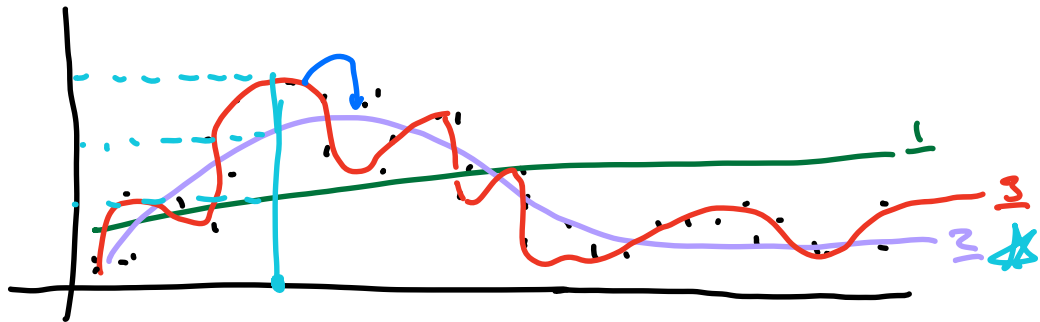


10,000 validation images → we perform poorly, because the network hasn't "memorized" them

# Polynomial Regression

$$y = x \theta_1 + x^2 \theta_2 + x^3 \theta_3 \dots$$

↖ degree

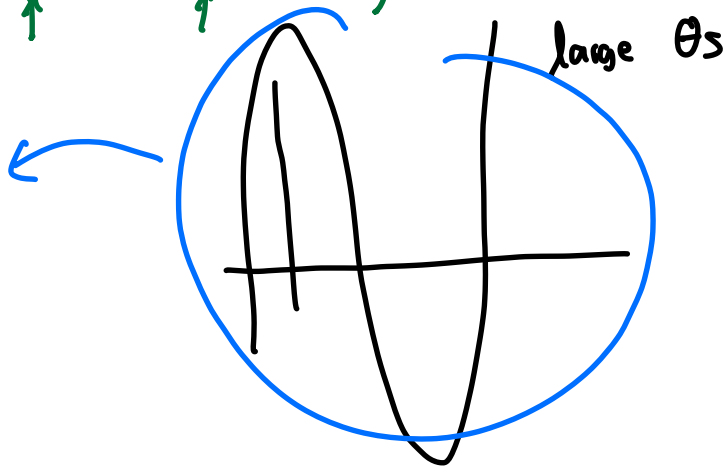
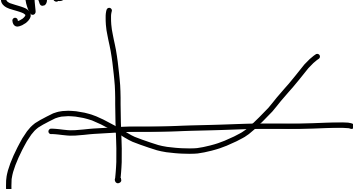


What is a symptom of the overfit model?

↳ High values for parameters.

$$y = \theta_0 + x \theta_1 + x^2 \theta_2 + x^3 \theta_3$$

Small  $\theta$ s

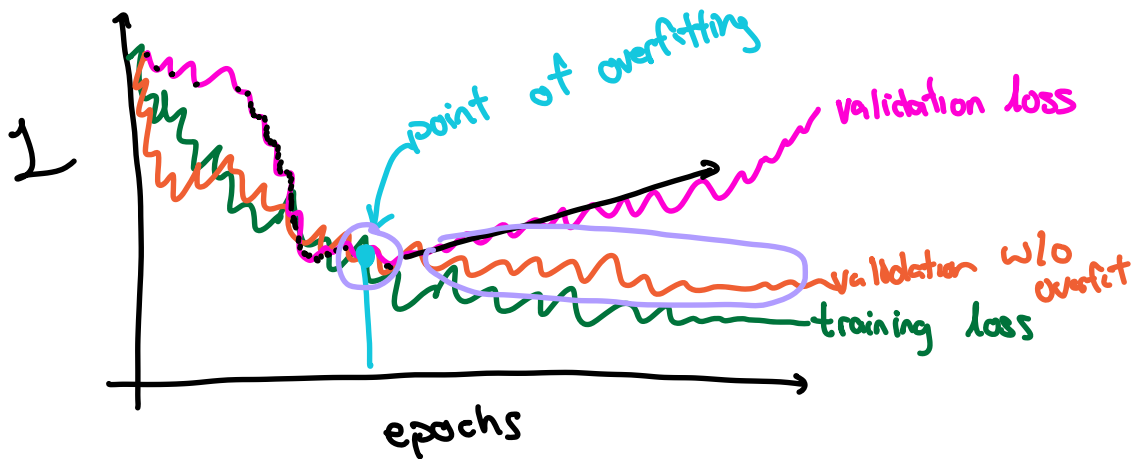


## Solutions

- early stopping
- parameter norm penalty (weight decay)
- dropout
- data augmentation

Cross validation  
ensemble methods

# Early Stopping



1. Track validation loss
2. Save model checkpoints
3. Take model just prior to validation loss not beating the best value  $\times$  updates in a row

## Parameter Norm Penalization

- Weight Decay
- Regularization
- L-1, L-2 Norm Penalty
- Ridge regression
- Tikhonov regularization

$$\mathcal{L}(\hat{y}, y)_{\text{HMSE}} = \frac{1}{2} \|(\hat{y} - y)^2\|_1$$

$$\mathcal{L}(\hat{y}, y) = \mathcal{L}(\hat{y}, y)_{\text{HMSE}} + \underbrace{\frac{\lambda}{2} \theta^T \theta}$$

What does this do to loss?

- When params are large?  $\rightarrow$  larger loss
- small?  $\rightarrow$  increase loss by smaller amount

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \theta} &= \frac{\partial \mathcal{L}_{\text{HMSE}}}{\partial \theta} + \frac{\partial}{\partial \theta} \underbrace{\frac{\lambda}{2} \theta^T \theta}_{\sum \theta_i^2} \\ &= \frac{\partial \mathcal{L}_{\text{HMSE}}}{\partial \theta} + \frac{\lambda}{2} 2\theta \\ &= \frac{\partial \mathcal{L}_{\text{HMSE}}}{\partial \theta} + \lambda \theta \end{aligned}$$

*weight decay param.*

How does regularization impact the parameter updates?

What if  $\theta_i$  is a large # (high positive)  
small # (high negative)

$$\Theta_{t+1} = \Theta_t - \eta \frac{\partial J}{\partial \Theta}$$

$$= \Theta_t - \eta \left( \frac{\partial J_{\text{HMSE}}}{\partial \Theta} + \lambda \Theta \right)$$

$$= \Theta_t - \eta \frac{\partial J_{\text{HMSE}}}{\partial \Theta} - \lambda \eta \Theta$$

1.  $1,756 - \eta \nabla J - \lambda \eta 1,756$  sub

2.  $-845 - \eta \nabla J - \lambda \eta (-845)$

$-845 - \eta \nabla J + \lambda \eta 845$  add

Drive to zero

Fixing the symptom, not the problem.