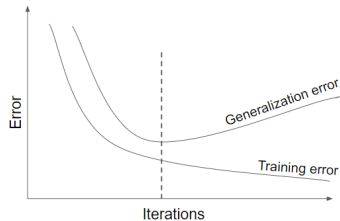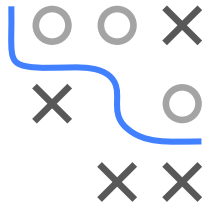# Introduction to Machine Learning
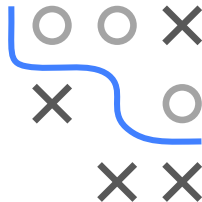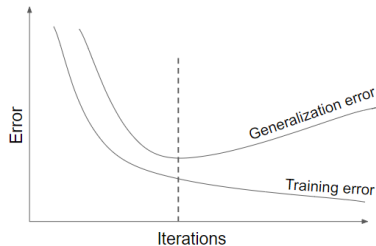
# Regularization
# Early Stopping



**Learning goals**

- Know how early stopping works
- Understand how early stopping acts as a regularizer
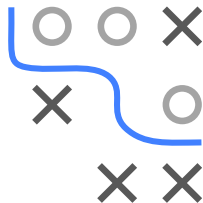
# EARLY STOPPING

- Especially for complex nonlinear models we can easily overfit
- In optimization: Often, after a certain number of iterations, generalization error begins to increase even though training error continues to decrease

## EARLY STOPPING AND *L2* ▸ Goodfellow, Bengio, and Courville 2016

| Strengths | Weaknesses |
|---|---|
| Effective and simple | Periodical evaluation of validation error |
| Applicable to almost any model without adjustment | Temporary copy of $\boldsymbol{\theta}$ (we have to save the whole model each time validation error improves) |
| Combinable with other regularization methods | Less data for training $\rightarrow$ include $\mathcal{D}_{\text{val}}$ afterwards |

- For simple case of LM with squared loss and GD optim initialized at $\boldsymbol{\theta} = 0$: Early stopping has exact correspondence with *L2* regularization/WD: optimal early-stopping iter $T_{\text{stop}}$ inversely proportional to $\lambda$ scaled by step-size $\alpha$
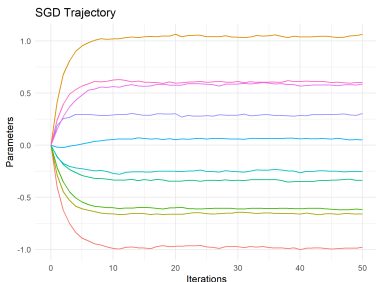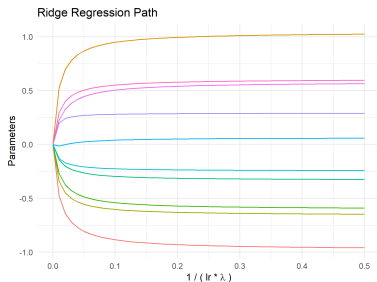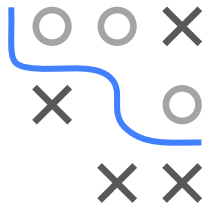
$$T_{\text{stop}} \approx \frac{1}{\alpha\lambda} \Leftrightarrow \lambda \approx \frac{1}{T_{\text{stop}}\alpha}$$

- Small $\lambda$ ( regu. $\downarrow$) $\Rightarrow$ large $T_{\text{stop}}$ (complexity $\uparrow$) and vice versa

# SGD TRAJECTORY AND *L2*  ▸ Ali, Dobriban, and Tibshirani 2020

Solution paths for *L2* regularized linear model closely matches SGD trajectory of unregularized LM initialized at $\boldsymbol{\theta} = 0$



**Caveat**: Initialization at the origin is crucial for this equivalence to hold, which is almost never exactly used in practice in ML/DL applications