

Training loss:

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N L(y(x^{(i)}), t^{(i)})$$

parameters of the model \uparrow θ \uparrow $x^{(i)}$ input \uparrow $t^{(i)}$ target
 \uparrow $y(x^{(i)})$ model prediction

Generalizes to new data?

Introduce a split between:

- training data \rightarrow used to find the parameters of our model
- validation data \rightarrow used to find the hyperparameters
- test data \rightarrow used to measure the generalization performance

\uparrow
 ability of the model to classify / predict on inputs it hasn't seen during training.

For each of these 3 sets, we can compute a loss:

- training loss
- validation loss
- test loss

Why do we need a validation set?

\rightarrow we can't tune hyperparameters on training data

\hookrightarrow we want hyperparameters that generalize.

example: picked a large number of neurons when designing a neural network could memorize training data

\rightarrow we can't tune hyperparameters on test data

\hookrightarrow a form of "cheating"

\hookrightarrow we should only measure test performance once

\hookrightarrow this would make it difficult to realistically evaluate/measure generalization.

Strategies for hyperparameter tuning

- Grid search: \rightarrow for each parameter, choose candidate values
 \rightarrow train one model for each possible combination of hyperparameter values.
 \rightarrow choose the combination that achieves minimal validation loss / error

- Random search: \rightarrow for each parameter, pick random values
 \rightarrow train a lot of models with different combinations of random hyperparameters values.
 \rightarrow choose the combination that achieves minimal validation loss.

Advantage of random search (over grid search) is that we have to try less combinations when the number of hyperparameters is large.

example: difficult to perform a grid search over 10 hyperparameters

but a random search could still be viable.

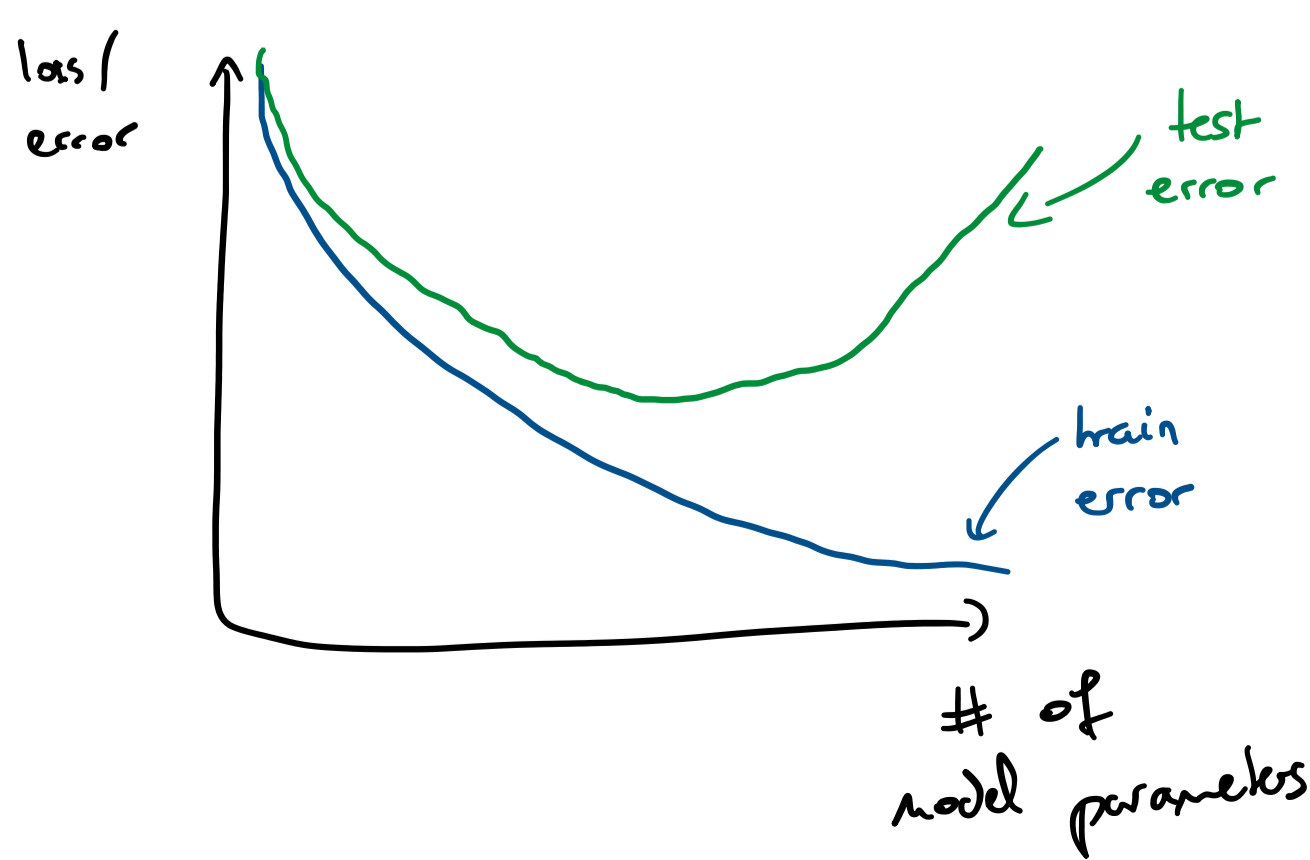
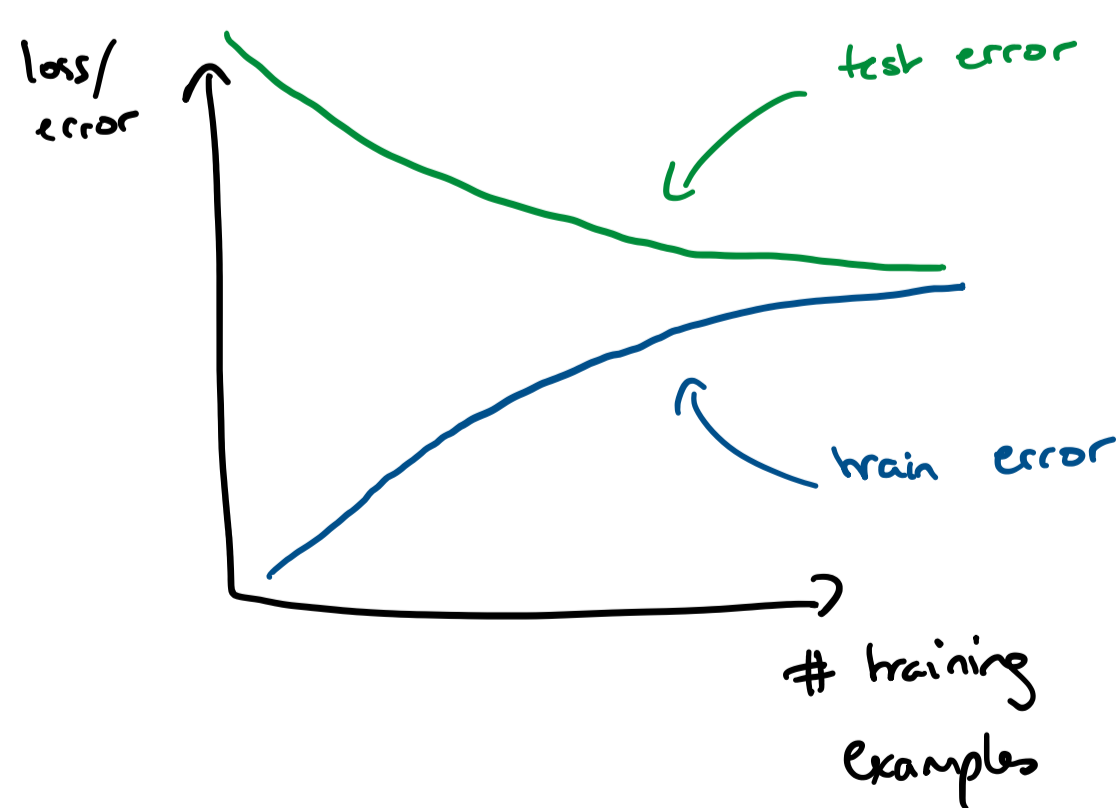
Reasoning about generalization

Avoid overfitting: avoid the model from recognizing data based on accidental regularities in the training set

Recall the example of a deep neural network which we covered in the video on the universal approximation theorem.

\hookrightarrow Pb: model has too much capacity (i.e., has ability to memorize training data)

Visualize generalization in practice:



(similar figure for loss wrt # epochs)