



# Model Selection

---



# Criteria to choose a model

---

- Error
- Training computational cost
- Test (or production) computational cost
- Ability to explain predictions



# Model Complexity

---

- We are looking for a model with good generalization
  - Able to correctly predict values for new examples
- If the model is less complex than the function that actually generated the data it will not represent it well and underfit it
  - Using a line to fit a quadratic function
- If the model is more complex then it will overfit and also “learn” the noise in the data
  - A 4th degree polynomial to approximate a linear process
- The complexity of a model depends on
  - Its degrees of freedom
  - The number of examples used to train it



# Symptoms of over and under fitting

---

- Underfitting

- Large error in the training data.
- This means that the model has a large bias.

- Overfitting

- Low error in the training set
- Large error in the validation or test set
- This means that the model has a large variance



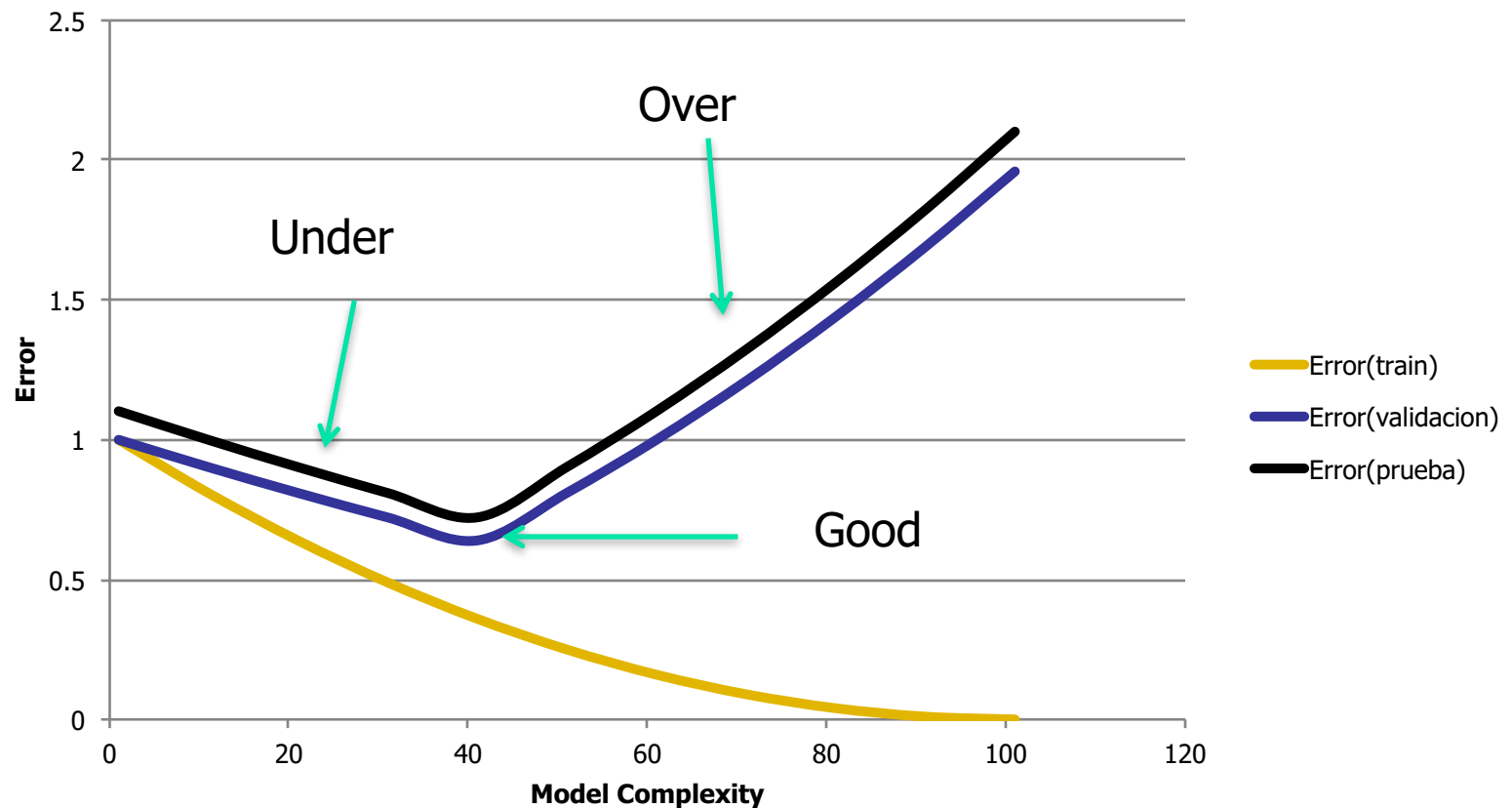
# ¿How can you determine the adequate complexity?

---

- Evaluate the training and validation errors for different model complexities
  - For each value of the complexity parameter estimate the generalization error
    - Use cross-validation or bootstrapping

# The Adequate Complexity

## Over and underfitting diagnosis





# ¿What to do?

---

- Overfitting (High variance)
  - Use more training examples
  - Reduce the complexity of the model
    - Reduce the number of attributes
    - Simplify the model by tuning its parameters
      - Increase regularization constant
      - Reduce the number of neurons in a neural net
      - Decrease the height of a decision tree



# ¿What to do?

---

- Underfitting (High bias)
  - Add attributes
    - Device new variables ( $x*x$ ,  $x*y$ ,  $\sin(x)$  etc). Use PCA, Factor analysis, etc.
    - Collect new and better variables
  - Increase the models complexity
    - Reduce the regularization constante
    - Use more neurons in a neural net
    - Change learning algorithm
  - Change data distribution





# Error estimation and hyper-parameter selection

---



# Model Creation

## Basic cycle

---

- Partition data into two groups
  - Training Set
    - Execute learning-validation cycle to choose model parameters (lambda, number of neurons, etc)
      - Details in following slides
  - Test set
    - The test set should not be used in any part of the learning and validation cycle
- Data should be partition following the original distribution and respecting temporal ordering
  - On occasion stratified sampling is in order
  - In general 25% of the data is left for testing
- Train with all the training set
- Report the error statistics from the learning-validation cycle for the selected hyperparameters and the test set error



# Model Selection

## Learning-validation cycle

---

- Most learning algorithms have hyper-parameters that need to be tuned
  - Lambda, neural net topology, etc.
- Hyper parameters need to be learnt too
  - The learning algorithm is basically trial and error
    - For each value  $h$  to explore of the hyper parameter
      - Train the model using  $h$  and estimate its error
    - Select the  $h$  with that yields the lowest error
- We cant use the test set to select  $h$  (why?)
  - We need an extra cycle, the learning-validation cycle



# Model Selection

## Learning-validation cycle

---

- We should execute learning-validation multiple times and gather error statistics
  - Randomize initialization parameters
  - Use different data subsets
- The different data subsets are generated using cross-validation or bootstrapping



# Learning-validation cycle

---

- Cross-validation

- Partition data into  $k$  subsets (normally between 5 and 15)
- Use one subset for validation and the remaining  $k-1$  for learning
- Repeat for each of the  $k$  subsets
- Compute statistics of the  $k$  results (mean and standard deviation of the error)

- Bootstrapping

- Suppose we have  $n$  data points in the training set
- Take  $k$  samples with replacement of size  $n$
- Divide each sample into learning and validation
  - Remove from validation data that appears in learning
- Compute statistics of the  $k$  results (mean and standard deviation of the error)

- Repeat process for each value of the hyper parameter and choose the one with the best statistics



# Note

---

- Attribute selection, standarization etc., mut go inside the learning-validation cycle



# Exercise

---

- Implement cross-validation to select the appropriate value of  $\lambda$  for the iterative linear regression model
  - Use the `regLinPoli2.csv` data file
  - You can use, if you wish, the `Regularization4class.ipynb` starter code



# Error Decomposition

---





# Error

---

- The learning error can be divided into three components
  - Irreducible error due to noisy data
  - Error due to high bias (unable to capture the phenomenon)
  - Error due to high variance (captures irrelevant aspects of the training)
- High bias is seen as underfitting and high variance as overfitting



# Error Decomposition

---

- Suppose the real function is of the form:
  - $y = f(x) + \varepsilon$ , where  $\varepsilon$  is noise normally distributed with mean zero and  $\sigma^2$  variance
- Our model produces a prediction
  - $\hat{V}(x)$ , for every  $x$
- And we measure error (in the case of regression) as
  - $\sum (y - \hat{V}(x))^2$ ,



# Error Decomposition

---

- We wish to estimate the error for a new data point  $x^*$

$$\text{Err}(x^*) = E[(y - \hat{V}(x^*))^2]$$

$$= E[(f(x^*) + \varepsilon - \hat{V}(x^*))^2]$$

$$= \sigma^2 + [E(\hat{V}(x^*)) - f(x^*)]^2 + E[\hat{V}(x^*) - E(\hat{V}(x^*))]^2$$

$$= \sigma^2 + \text{Bias}^2(\hat{V}(x^*)) + \text{Var}(\hat{V}(x^*))$$

$$= \text{Irreducible\_error} + \text{Bias}^2 + \text{Variance}$$

- Normally there is a tradeoff between bias and variance



# Error Decomposition

---

- Note that expectations are taken over everything that can be randomized
  - Initial weights ( $w$ 's)
  - Training set (expectation over all possible training sets)
    - For example in the iterative linear regression:  $E(\hat{V}(x^*))$  is the expected output of the model over all possible training sets and all possible initial  $w$ 's

Derivación

Descomposición de Error

---



# Derivación Versión 1

---

- Una propiedad importante (truco para derivar)
  - $\text{Var}(X) = E(X^2) - [E(X)]^2$
- Sustituimos la variable aleatoria  $X$  por la discrepancia de nuestro modelo
  - $\text{Var}(V^{\wedge}(x) - f(x) - \varepsilon) = \text{Var}(V^{\wedge}(x)) + \sigma^2$ 
    - Porque la varianza de  $f(x)$  es cero pues no es una variable aleatoria y la covarianza entre el ruido y  $V^{\wedge}(x)$  es cero



# Derivación Versión 1

---

- De la fórmula de la varianza sustituyendo:
- $\text{Var}(\hat{V}(x)) + \sigma^2 = E[(\hat{V}(x) - f(x) - \varepsilon)^2] - (E[\hat{V}(x) - f(x) - \varepsilon])^2$ 
  - $E[(\hat{V}(x) - f(x) - \varepsilon)^2] = \text{MSE}$  (error cuadrático medio)
  - $(E[\hat{V}(x) - f(x) - \varepsilon])^2 = (E(\hat{V}(x)) - E(f(x)) - E(\varepsilon))^2$   
 $= [E(\hat{V}(x)) - f(x)]^2 = \text{Bias}^2$ 
    - Porque  $E(f(x)) = f(x)$  y  $E(\varepsilon) = 0$
- Sustituyendo en la primer formula
- $\text{Var}(\hat{V}(x)) + \sigma^2 = \text{MSE} - \text{Bias}^2$
- $\text{MSE} = \text{Var}(\hat{V}(x)) + \sigma^2 + \text{Bias}^2$



# Version 2

## Derivación

---

- Algunas propiedades importantes:

- 1.  $E(E(x))=E(x)$

- 2.  $E((x-E(x))^2)=E[x^2-2xE(x)-E(x)^2]$   
 $=E(x^2)-2E(xE(x))+E[E(x)^2]$   
 $=E(x^2)-2E(x)E(x)+E(x)^2$   
 $=E(x^2)-E(x)^2$

- 3.  $E(x^2)=E((x-E(x))^2) + E(x)^2$  (fórmula varianza)

- 4.  $E((c+N(0,\sigma))x)$   
 $=E(cx+xN(0,\sigma))=cE(x)$  (la covarianza es cero)





# Derivación

---

- Regresando al error esperado:

$$\begin{aligned} E[(y - \hat{V}(x^*))^2] &= E[y^2 - 2y\hat{V}(x^*) + \hat{V}(x^*)^2] \\ &= E[y^2] - 2E(y\hat{V}(x^*)) + E[\hat{V}(x^*)^2] \\ &= E((y - E(y))^2) + E(y)^2 \text{ (propiedad 3)} \\ &\quad - 2E(\hat{V}(x^*))f(x^*) \text{ (propiedad 4)} \\ &\quad + E[(\hat{V}(x^*) - E(\hat{V}(x^*)))^2] + E(\hat{V}(x^*))^2 \text{ (propiedad 3)} \\ &= E((y - E(y))^2) \text{ (ruido. El desarrollo da } \sigma^2 \text{ usando prop.4 y 1)} \\ &\quad + E(y)^2 - 2E(\hat{V}(x^*))f(x^*) + E(\hat{V}(x^*))^2 \text{ (sesgo}^2 \text{ esto se} \\ &\quad \text{reduce a } (y - E(\hat{V}(x^*)))^2 \text{ note que } E(y) = y) \\ &\quad + E[(\hat{V}(x^*) - E(\hat{V}(x^*)))^2] \text{ (varianza)} \end{aligned}$$