

# Machine Learning for Finance

## UCEMA 2020

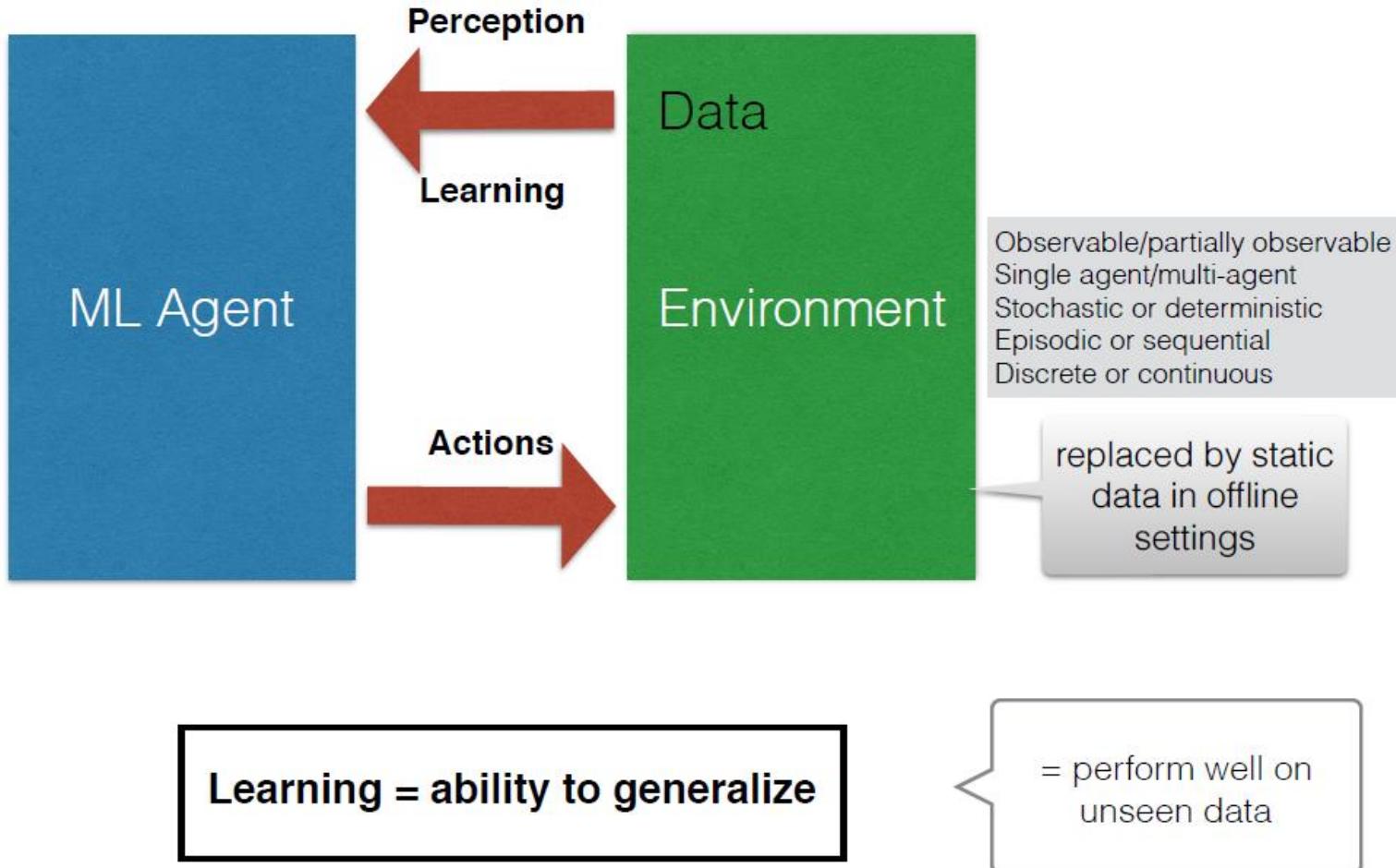
# Lecture 2

- Bias-variance tradeoff
- Training and test sets
- Model capacity and overfitting
- Learning in linear regressions
- Regularizations
- Hyperparameters and cross validation
- Gradient descent-like algorithms
- Probabilistic models, MLE, MAP
- Probabilistic classifications
- Logistic regressions
- Scikit-Learn demo

# Lectures Structure

- 1<sup>st</sup> block
  - ~1h: Presentation of subjects + code explanation.
  - ~15': Discussion + Q&A.
- 2<sup>nd</sup> block
  - ~1h: Presentation of subjects + code explanation.
  - ~15': Discussion + Q&A.

# Generalization as a Goal of Learning



# Generalization Error in Regression

# Generalization Error in Regression

- Consider a **regression problem**

$$y = f(\mathbf{x}) + \varepsilon$$

- Here  $\mathbb{E}[\varepsilon] = 0, \mathbb{E}[\varepsilon^2] = \sigma^2$

# Generalization Error in Regression

- Consider a **regression problem**

$$y = f(\mathbf{x}) + \varepsilon$$

- Here  $\mathbb{E}[\varepsilon] = 0, \mathbb{E}[\varepsilon^2] = \sigma^2$
- This produces  $\mathbb{E}[y] = f(x), \text{Var}(y) = \sigma^2$

# Generalization Error in Regression

- Consider a regression problem

$$y = f(\mathbf{x}) + \varepsilon$$

- Here  $\mathbb{E}[\varepsilon] = 0, \mathbb{E}[\varepsilon^2] = \sigma^2$
- This produces  $\mathbb{E}[y] = f(x), \text{Var}(y) = \sigma^2$
- We look for approximation  $f(\mathbf{x}) \approx \hat{f}(\mathbf{x})$  such that the **generalization error** (expected loss)  $\mathbb{E}[(y - \hat{f}(x))^2]$  is minimized over all possible data, both seen and unseen!

# Generalization Error in Regression

- Consider a regression problem

$$y = f(\mathbf{x}) + \varepsilon$$

- Here  $\mathbb{E}[\varepsilon] = 0, \mathbb{E}[\varepsilon^2] = \sigma^2$
- This produces  $\mathbb{E}[y] = f(x), \text{Var}(y) = \sigma^2$
- We look for approximation  $f(\mathbf{x}) \approx \hat{f}(\mathbf{x})$  such that the **generalization error** (expected loss)  $\mathbb{E}\left[\left(y - \hat{f}(x)\right)^2\right]$  is minimized over all possible data, both seen and unseen!
- Evaluate:  $\mathbb{E}\left[\left(y - \hat{f}(x)\right)^2\right] = \mathbb{E}[y^2] + \mathbb{E}[\hat{f}^2] - 2\mathbb{E}[y\hat{f}]$

# Generalization Error in Regression

- Consider a regression problem

$$y = f(\mathbf{x}) + \varepsilon$$

- Here  $\mathbb{E}[\varepsilon] = 0, \mathbb{E}[\varepsilon^2] = \sigma^2$

- This produces  $\mathbb{E}[y] = f(x), \text{Var}(y) = \sigma^2$

- We look for approximation  $f(\mathbf{x}) \approx \hat{f}(\mathbf{x})$  such that the **generalization error** (expected loss)  $\mathbb{E}\left[\left(y - \hat{f}(x)\right)^2\right]$  is minimized over all possible data, both seen and unseen!

- Evaluate:  $\mathbb{E}\left[\left(y - \hat{f}(x)\right)^2\right] = \mathbb{E}[y^2] + \mathbb{E}[\hat{f}^2] - 2\mathbb{E}[y\hat{f}] =$

$$\text{Var}(y) + \mathbb{E}[y]^2 + \text{Var}(\hat{f}) + \mathbb{E}[\hat{f}]^2 - 2\mathbb{E}[y\hat{f}] =$$

$$(y = f(\mathbf{x}) + \varepsilon)$$

# Generalization Error in Regression

- Consider a regression problem

$$y = f(\mathbf{x}) + \varepsilon$$

- Here  $\mathbb{E}[\varepsilon] = 0, \mathbb{E}[\varepsilon^2] = \sigma^2$
- This produces  $\mathbb{E}[y] = f(x), \text{Var}(y) = \sigma^2$
- We look for approximation  $f(\mathbf{x}) \approx \hat{f}(\mathbf{x})$  such that the **generalization error** (expected loss)  $\mathbb{E}\left[(y - \hat{f}(x))^2\right]$  is minimized over all possible data, both seen and unseen!

- Evaluate:  $\mathbb{E}\left[(y - \hat{f}(x))^2\right] = \mathbb{E}[y^2] + \mathbb{E}[\hat{f}^2] - 2\mathbb{E}[y\hat{f}] =$

$$\text{Var}(y) + \mathbb{E}[y]^2 + \text{Var}(\hat{f}) + \mathbb{E}[\hat{f}]^2 - 2\mathbb{E}[y\hat{f}] =$$

$$\text{Var}(y) + \text{Var}(\hat{f}) + \mathbb{E}[f - \hat{f}]^2$$

$y = f(\mathbf{x}) + \varepsilon$

# Generalization Error in Regression

- Consider a regression problem

$$y = f(\mathbf{x}) + \varepsilon$$

- Here  $\mathbb{E}[\varepsilon] = 0, \mathbb{E}[\varepsilon^2] = \sigma^2$
- This produces  $\mathbb{E}[y] = f(x), \text{Var}(y) = \sigma^2$
- We look for approximation  $f(\mathbf{x}) \approx \hat{f}(\mathbf{x})$  such that the **generalization error** (expected loss)  $\mathbb{E}\left[\left(y - \hat{f}(x)\right)^2\right]$  is minimized over all possible data, both seen and unseen!
- Evaluate:  $\mathbb{E}\left[\left(y - \hat{f}(x)\right)^2\right] = \mathbb{E}[y^2] + \mathbb{E}[\hat{f}^2] - 2\mathbb{E}[y\hat{f}] =$   
 $Var(y) + \mathbb{E}[y]^2 + Var(\hat{f}) + \mathbb{E}[\hat{f}]^2 - 2\mathbb{E}[y\hat{f}] =$   
 $Var(y) + Var(\hat{f}) + \mathbb{E}\left[f - \hat{f}\right]^2$   
noise ( $\sigma^2$ )

$y = f(\mathbf{x}) + \varepsilon$

# Generalization Error in Regression

- Consider a regression problem

$$y = f(\mathbf{x}) + \varepsilon$$

- Here  $\mathbb{E}[\varepsilon] = 0, \mathbb{E}[\varepsilon^2] = \sigma^2$
- This produces  $\mathbb{E}[y] = f(x), \text{Var}(y) = \sigma^2$
- We look for approximation  $f(\mathbf{x}) \approx \hat{f}(\mathbf{x})$  such that the **generalization error** (expected loss)  $\mathbb{E}\left[\left(y - \hat{f}(x)\right)^2\right]$  is minimized over all possible data, both seen and unseen!
- Evaluate:  $\mathbb{E}\left[\left(y - \hat{f}(x)\right)^2\right] = \mathbb{E}[y^2] + \mathbb{E}[\hat{f}^2] - 2\mathbb{E}[y\hat{f}] =$   
 $\text{Var}(y) + \mathbb{E}[y]^2 + \text{Var}(\hat{f}) + \mathbb{E}[\hat{f}]^2 - 2\mathbb{E}[y\hat{f}] =$   
 $\text{Var}(y) + \text{Var}(\hat{f}) + \mathbb{E}\left[f - \hat{f}\right]^2$   
noise ( $\sigma^2$ ) variance

$y = f(\mathbf{x}) + \varepsilon$

# Generalization Error in Regression

- Consider a regression problem

$$y = f(\mathbf{x}) + \varepsilon$$

- Here  $\mathbb{E}[\varepsilon] = 0, \mathbb{E}[\varepsilon^2] = \sigma^2$
- This produces  $\mathbb{E}[y] = f(x), \text{Var}(y) = \sigma^2$
- We look for approximation  $f(\mathbf{x}) \approx \hat{f}(\mathbf{x})$  such that the **generalization error** (expected loss)  $\mathbb{E}\left[(y - \hat{f}(x))^2\right]$  is minimized over all possible data, both seen and unseen!

- Evaluate:  $\mathbb{E}\left[(y - \hat{f}(x))^2\right] = \mathbb{E}[y^2] + \mathbb{E}[\hat{f}^2] - 2\mathbb{E}[y\hat{f}] =$

$$\text{Var}(y) + \mathbb{E}[y]^2 + \text{Var}(\hat{f}) + \mathbb{E}[\hat{f}]^2 - 2\mathbb{E}[y\hat{f}] =$$

$$\text{noise } (\sigma^2) \quad \text{variance} \quad \text{bias}$$

$y = f(\mathbf{x}) + \varepsilon$

# The Bias-Variance Decomposition

# The Bias-Variance Decomposition

This is bias-variance decomposition:

$$\mathbb{E}\left[\left(y - \hat{f}(\mathbf{x})\right)^2\right] = (\text{bias})^2 + \text{variance} + \text{noise}$$

Here:

$$(\text{bias})^2 = \mathbb{E}\left[f - \hat{f}\right]^2$$

$$\text{variance} = \text{Var}(\hat{f})$$

$$\text{noise} = \text{Var}(y) = \sigma^2$$

# The Bias-Variance Decomposition

This is bias-variance decomposition:

$$\mathbb{E}\left[\left(y - \hat{f}(\mathbf{x})\right)^2\right] = (\text{bias})^2 + \text{variance} + \text{noise}$$

Here:

$$(\text{bias})^2 = \mathbb{E}\left[f - \hat{f}\right]^2$$

$$\text{variance} = \text{Var}(\hat{f})$$

$$\text{noise} = \text{Var}(y) = \sigma^2$$

**Bias:** the (square of) expected difference of approximate predictor  $\hat{f}(\mathbf{x})$  from the “true” predictor  $f(\mathbf{x})$

# The Bias-Variance Decomposition

This is bias-variance decomposition:

$$\mathbb{E}\left[\left(y - \hat{f}(\mathbf{x})\right)^2\right] = (\text{bias})^2 + \text{variance} + \text{noise}$$

Here:

$$(\text{bias})^2 = \mathbb{E}\left[f - \hat{f}\right]^2$$

$$\text{variance} = \text{Var}(\hat{f})$$

$$\text{noise} = \text{Var}(y) = \sigma^2$$

**Bias:** the (square of) expected difference of approximate predictor  $\hat{f}(\mathbf{x})$  from the “true” predictor  $f(\mathbf{x})$

**Variance:** sensitivity of  $\hat{f}(\mathbf{x})$  to the choice of data set

# The Bias-Variance Decomposition

This is bias-variance decomposition:

$$\mathbb{E}\left[\left(y - \hat{f}(\mathbf{x})\right)^2\right] = (\text{bias})^2 + \text{variance} + \text{noise}$$

Here:

$$(\text{bias})^2 = \mathbb{E}\left[f - \hat{f}\right]^2$$

$$\text{variance} = \text{Var}(\hat{f})$$

$$\text{noise} = \text{Var}(y) = \sigma^2$$

**Bias:** the (square of) expected difference of approximate predictor  $\hat{f}(\mathbf{x})$  from the “true” predictor  $f(\mathbf{x})$

**Variance:** sensitivity of  $\hat{f}(\mathbf{x})$  to the choice of data set

**Noise:** a property of data; beyond our control

# The Bias-Variance Tradeoff

# The Bias-Variance Tradeoff

$$\mathbb{E}\left[\left(y - \hat{f}(\mathbf{x})\right)^2\right] = (\text{bias})^2 + \text{variance} + \text{noise}$$

- This relation is mostly of a theoretical value: we do not know how to compute the bias and variance!
- But it shows a general decomposition of expected loss of a generic regression model.

# The Bias-Variance Tradeoff

$$\mathbb{E}\left[\left(y - \hat{f}(\mathbf{x})\right)^2\right] = (\text{bias})^2 + \text{variance} + \text{noise}$$

- This relation is mostly of a theoretical value: we do not know how to compute the bias and variance!
- But it shows a general decomposition of expected loss of a generic regression model.
- **Bias-variance tradeoff:**
  - Complex models (more features) tend to have low bias and high variance

# The Bias-Variance Tradeoff

$$\mathbb{E}\left[\left(y - \hat{f}(\mathbf{x})\right)^2\right] = (\text{bias})^2 + \text{variance} + \text{noise}$$

- This relation is mostly of a theoretical value: we do not know how to compute the bias and variance!
- But it shows a general decomposition of expected loss of a generic regression model.
- **Bias-variance tradeoff:**
  - Complex models (more features) tend to have low bias and high variance
  - Simple models (less features) tend to have high bias and low variance

# The Bias-Variance Tradeoff

$$\mathbb{E}\left[\left(y - \hat{f}(\mathbf{x})\right)^2\right] = (\text{bias})^2 + \text{variance} + \text{noise}$$

- This relation is mostly of a theoretical value: we do not know how to compute the bias and variance!
- But it shows a general decomposition of expected loss of a generic regression model.
- **Bias-variance tradeoff:**
  - Complex models (more features) tend to have low bias and high variance
  - Simple models (less features) tend to have high bias and low variance
  - Important to control **model complexity** to have an optimal tradeoff!

# Data Representation

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

Feature vector

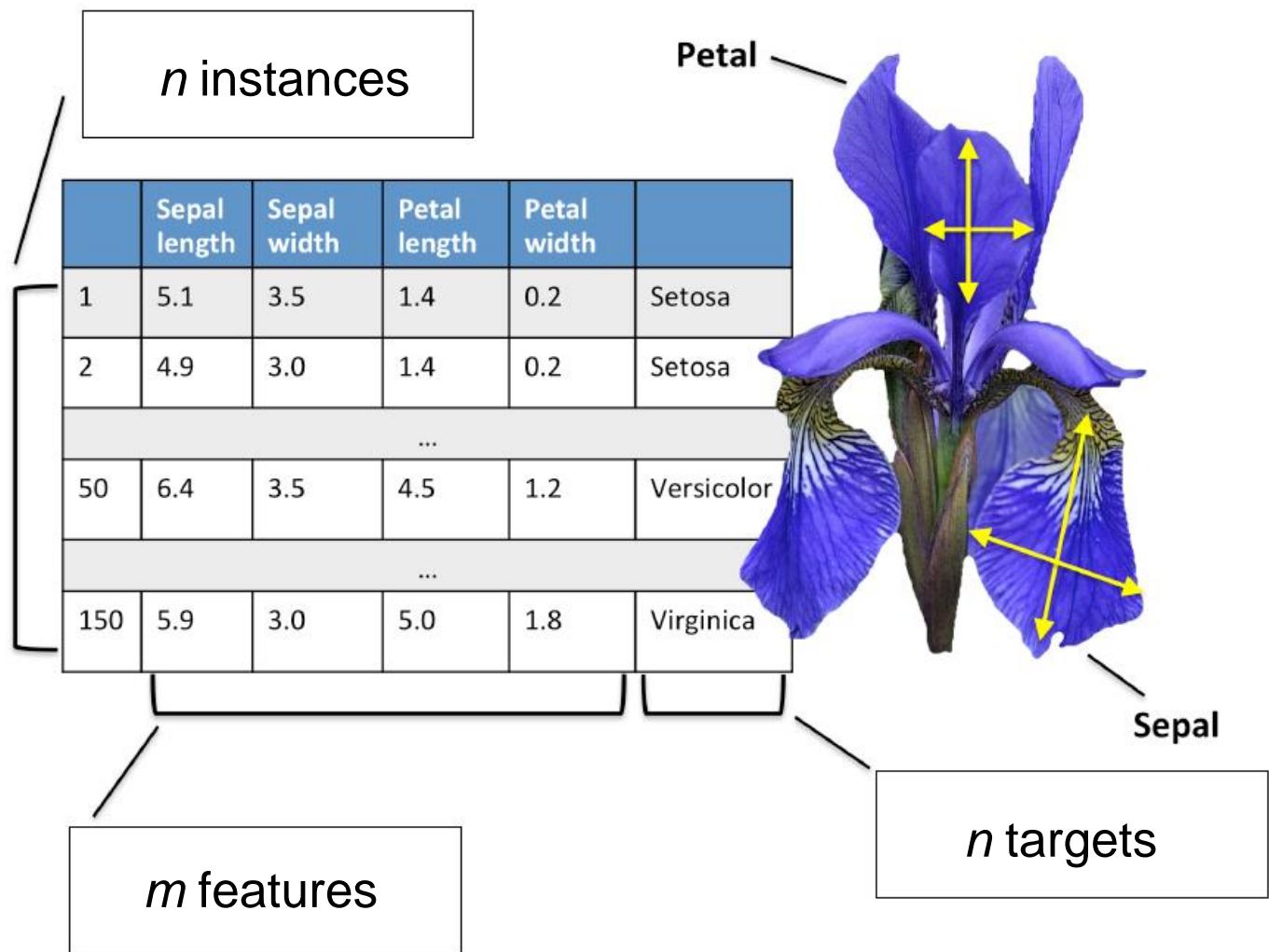
$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}$$

Design Matrix

$$\mathbf{X} = \begin{bmatrix} x_1^{[1]} & x_2^{[1]} & \cdots & x_m^{[1]} \\ x_1^{[2]} & x_2^{[2]} & \cdots & x_m^{[2]} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{[n]} & x_2^{[n]} & \cdots & x_m^{[n]} \end{bmatrix}$$

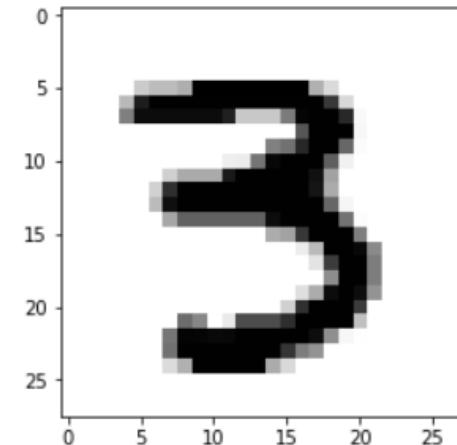
Design Matrix

# Structured Data



# Unstructured Data (Images)

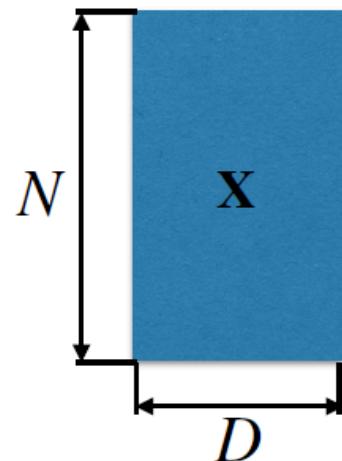
## "traditional methods"



# Training and Test Sets

# Training and Test Sets

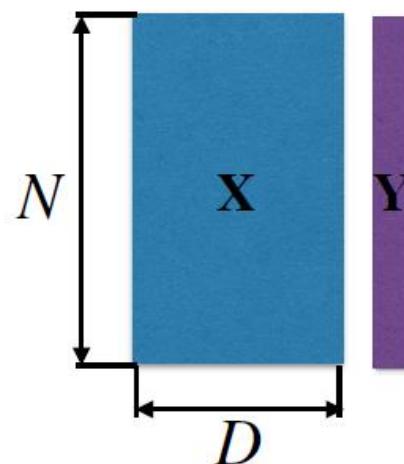
**Design matrix** (dimension  $N \times D$ ):



*i.i.d.* samples  
from a data-  
generating  
distribution  
 $p_{data}$

# Training and Test Sets

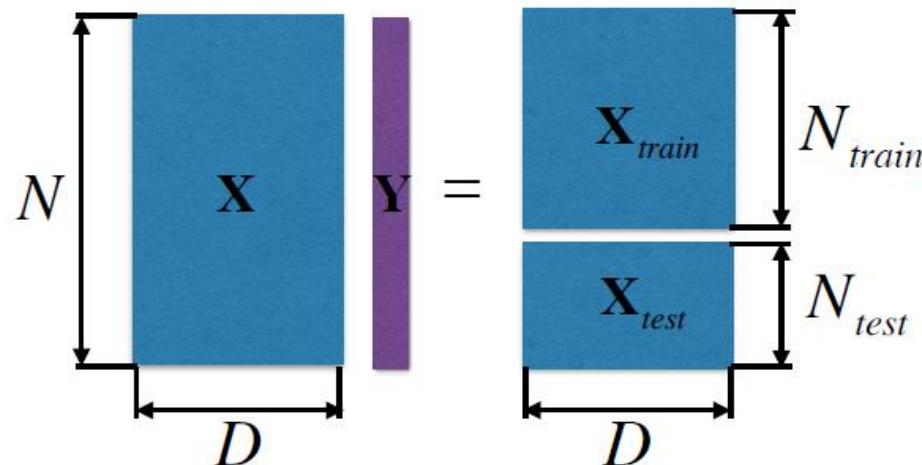
**Design matrix** (dimension  $N \times D$ ):



*i.i.d.* samples  
from a data-  
generating  
distribution  
 $p_{data}$

# Training and Test Sets

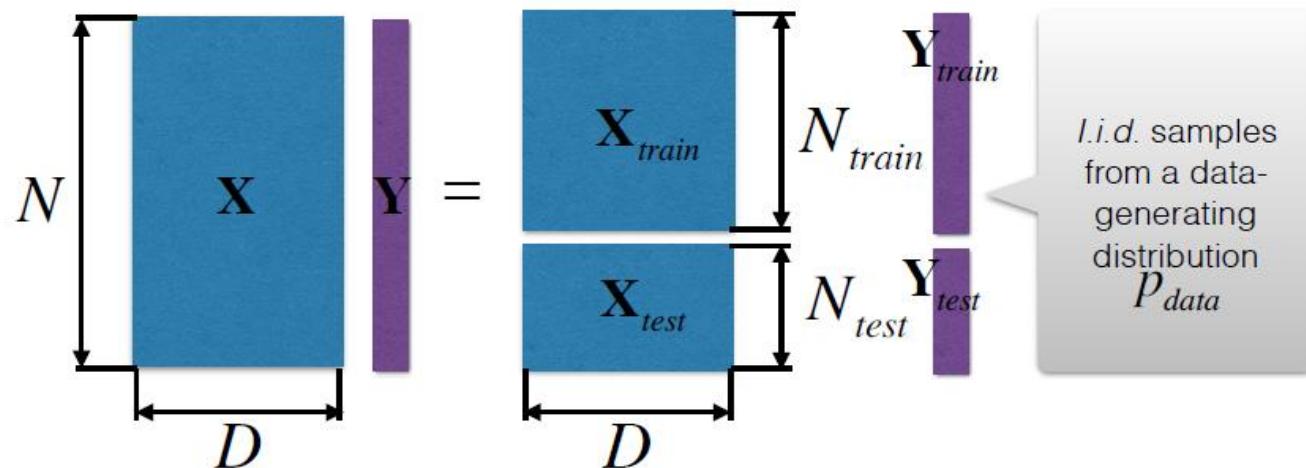
**Design matrix** (dimension  $N \times D$ ):



*i.i.d.* samples  
from a data-  
generating  
distribution  
 $p_{data}$

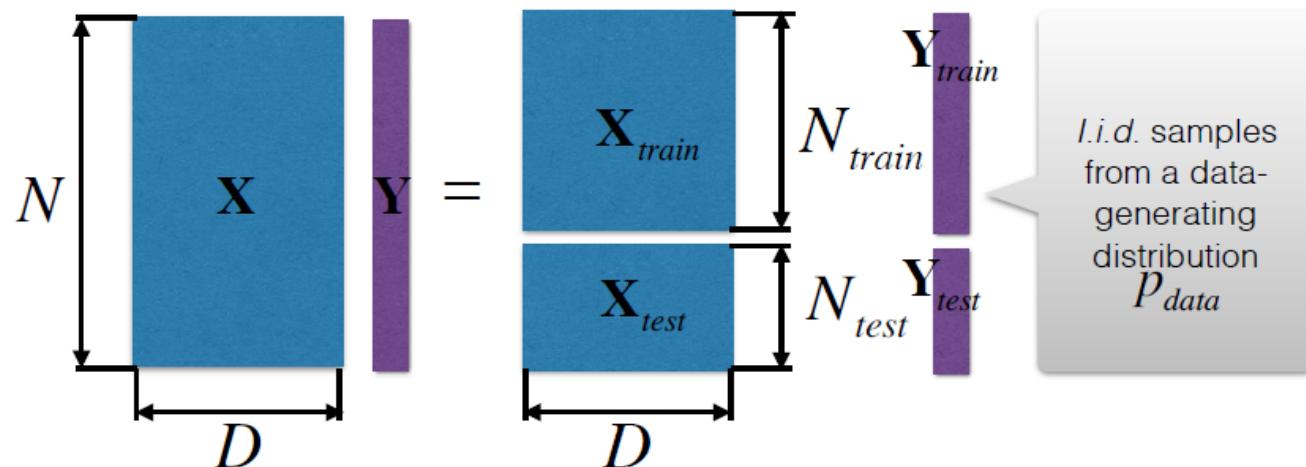
# Training and Test Sets

**Design matrix** (dimension  $N \times D$ ):



# Training and Test Sets

**Design matrix** (dimension  $N \times D$ ):

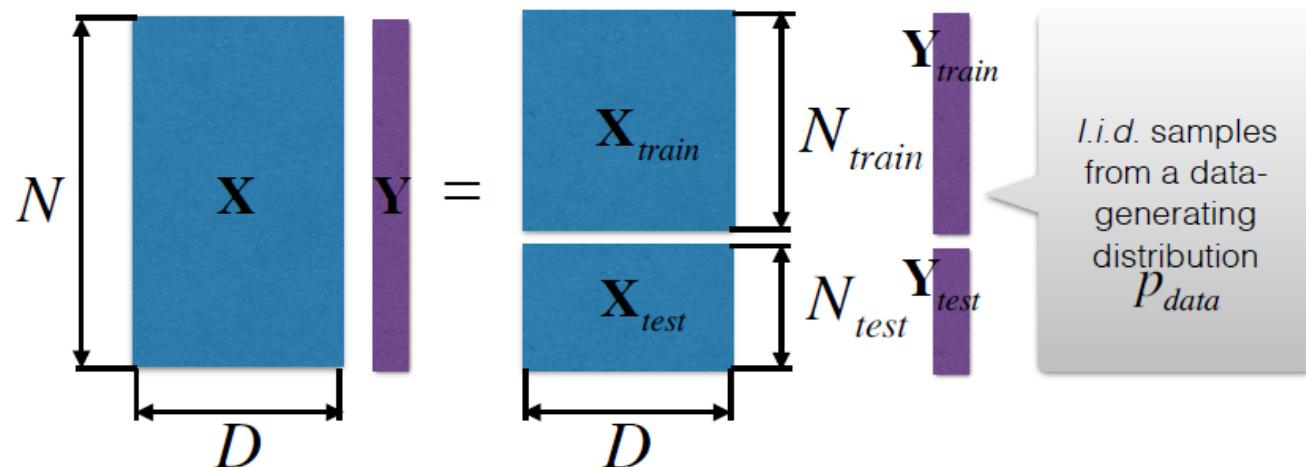


A model is trained using only a **training set**  $(\mathbf{X}_{train}, y_{train}) \sim p_{data}$

A **test set**  $(\mathbf{X}_{test}, y_{test}) \sim p_{data}$  is used to estimate algorithm's ability to **generalize**, i.e. perform well on unseen data.

# Training and Test Sets

**Design matrix** (dimension  $N \times D$ ):

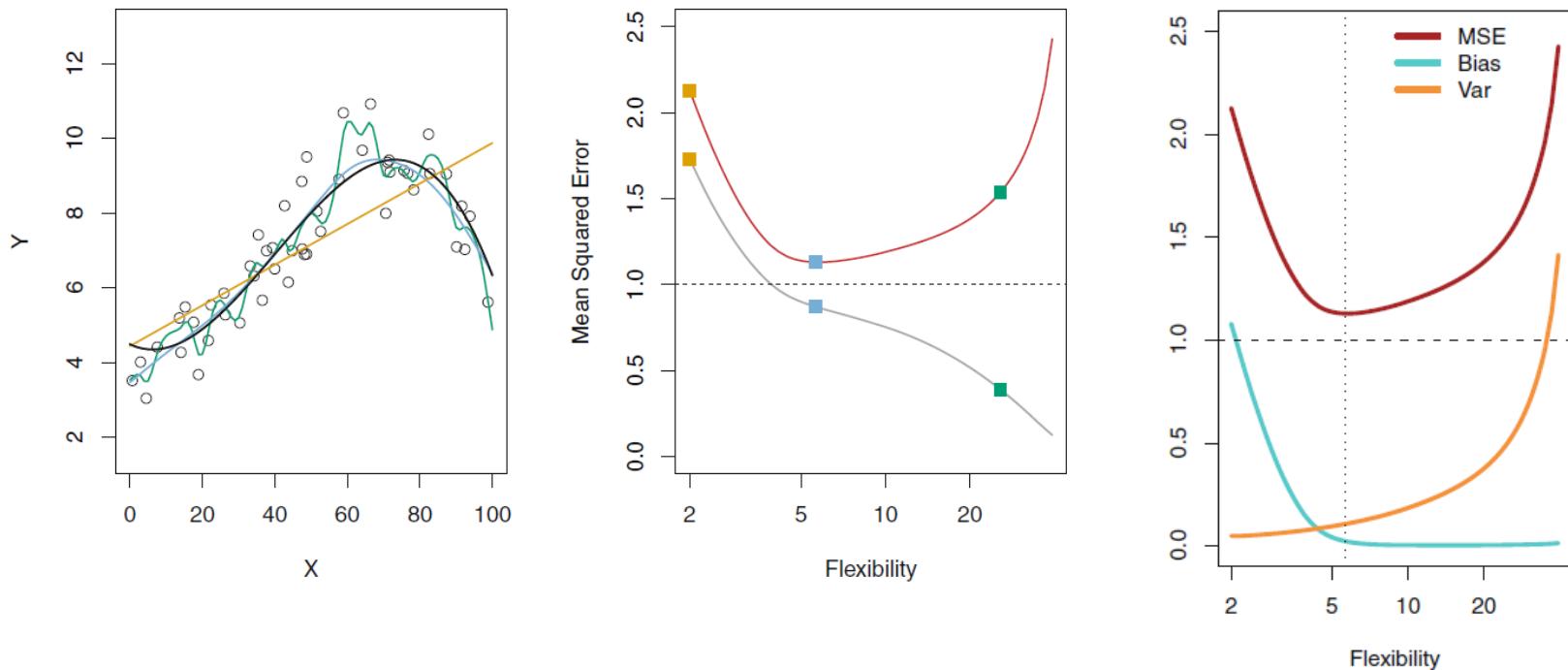


A model is trained using only a **training set**  $(\mathbf{X}_{train}, y_{train}) \sim p_{data}$

A **test set**  $(\mathbf{X}_{test}, y_{test}) \sim p_{data}$  is used to estimate algorithm's ability to **generalize**, i.e. perform well on unseen data.

More specifically, a test set is used to detect when a ML algorithm starts to **overfit** data, by estimating a generalization error by a test error.

# The Bias-Variance Tradeoff

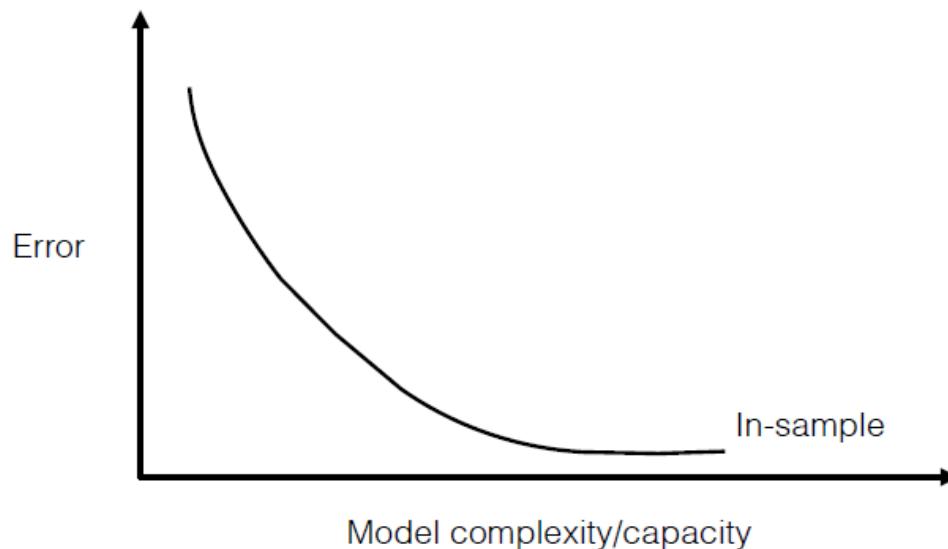


*Data simulated from  $f$ , shown in black. Three estimates of  $f$  are shown: the linear regression line (orange curve), and two smoothing spline fits (blue and green curves). Right: Training MSE (grey curve), test MSE (red curve), and minimum possible test MSE over all methods (dashed line). Squares represent the training and test MSEs for the three fits shown in the left-hand panel.*

# Overfitting

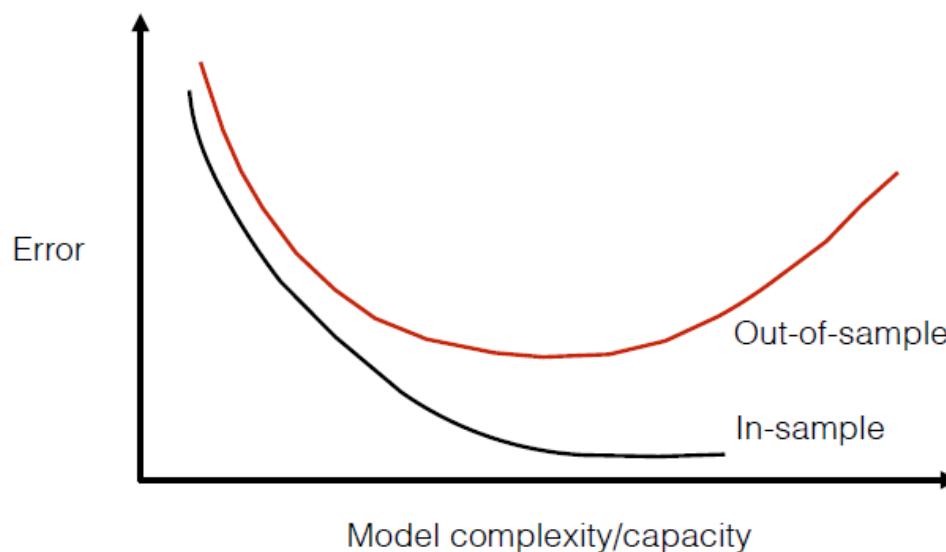
# Overfitting

Trying to exactly match all available data is almost always a bad idea



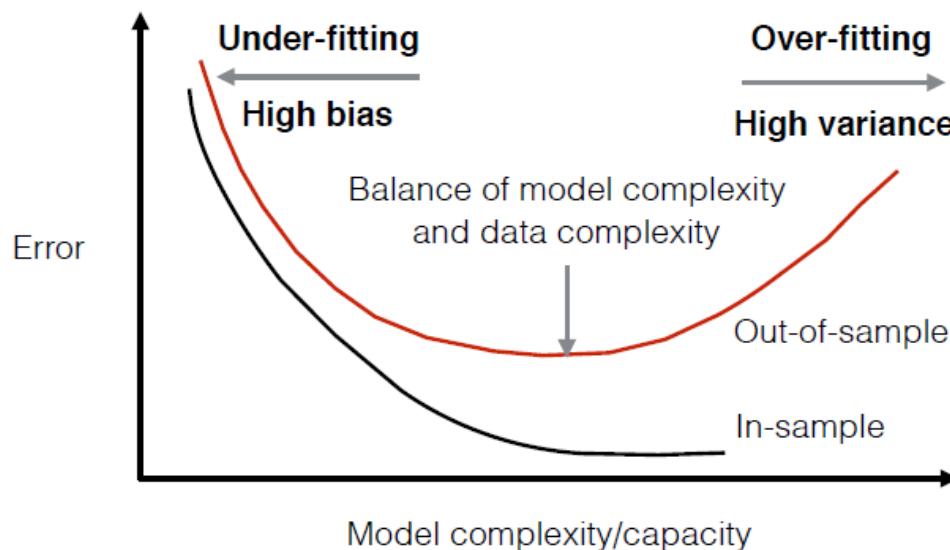
# Overfitting

Trying to exactly match all available data is almost always a bad idea



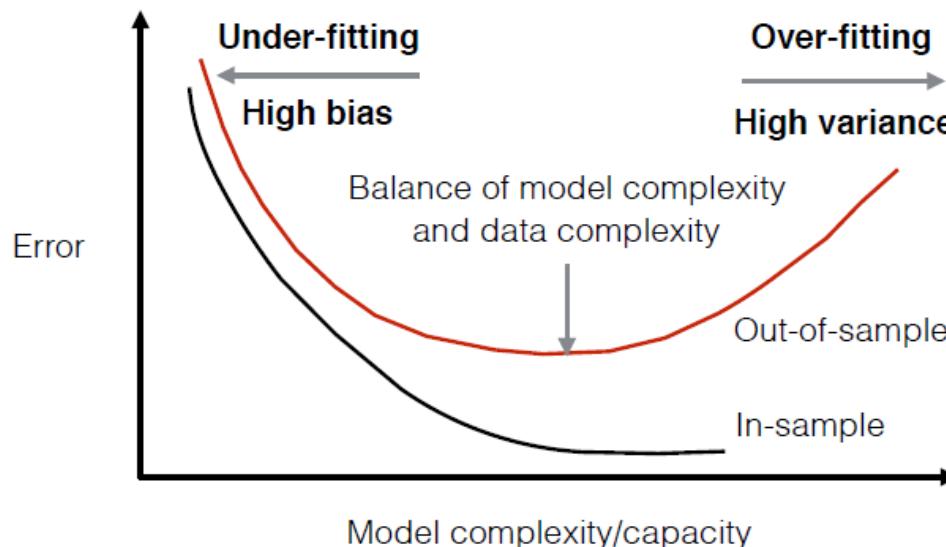
# Overfitting

Trying to exactly match all available data is almost always a bad idea



# Overfitting

Trying to exactly match all available data is almost always a bad idea



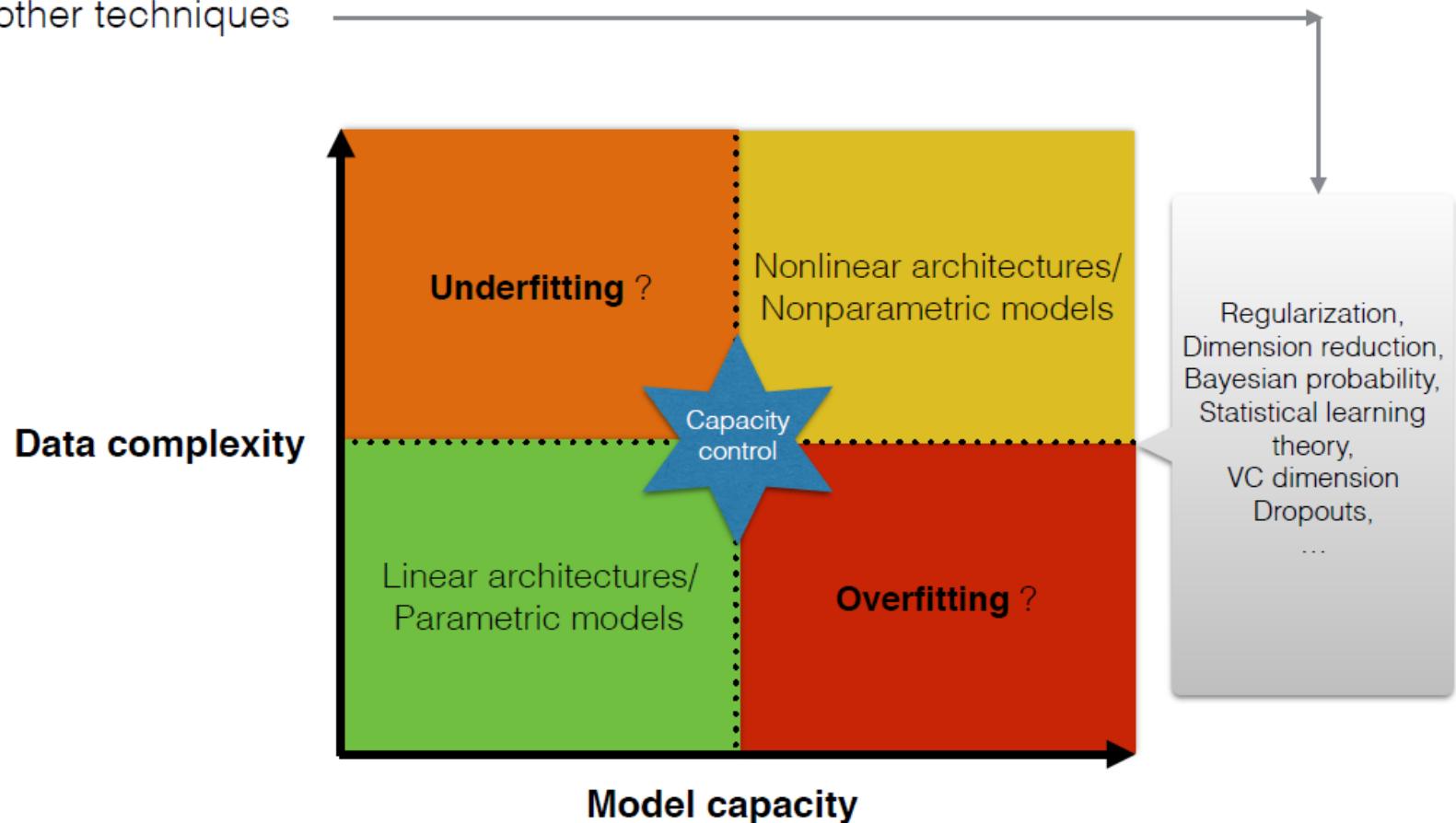
A good ML algorithm should achieve **two goals**:

1. Make the **training error** small (avoid under-fitting)
2. Make the gap between training and **test errors** small (avoid over-fitting)

Key ingredients: 1) data is i.i.d.  $\sim p_{data}$ , and 2) model **capacity control**

# Model Capacity and Overfitting

- Model **capacity** controls model's ability to fit a wide variety of functions.
- Models with low capacity can **under-fit**, but models with high capacity can **over-fit!**
- Capacity is controlled by the choice of a **hypothesis space** (architecture), and other techniques



# Linear Regression as a ML Task

**Task:** predict a scalar value  $y \in \mathbb{R}$  from a vector of predictors (“**features**”)  $\mathbf{X} = (X_0, X_1, \dots, X_{D-1}) \in \mathbb{R}^D$  ( $D \geq 1$  is the dimension of the feature space, or the number of predictors).

# Linear Regression as a ML Task

**Task:** predict a scalar value  $y \in \mathbb{R}$  from a vector of predictors (“**features**”)  $\mathbf{X} = (X_0, X_1, \dots, X_{D-1}) \in \mathbb{R}^D$  ( $D \geq 1$  is the dimension of the feature space, or the number of predictors).

**Given:** a dataset  $(\mathbf{X}, y)_{data} = [(\mathbf{X}_{train}, y_{train}), (\mathbf{X}_{test}, y_{test})] \sim p_{data}$

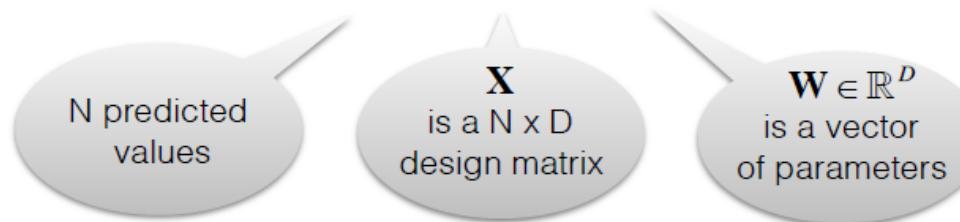
# Linear Regression as a ML Task

**Task:** predict a scalar value  $y \in \mathbb{R}$  from a vector of predictors (“**features**”)  $\mathbf{X} = (X_0, X_1, \dots, X_{D-1}) \in \mathbb{R}^D$  ( $D \geq 1$  is the dimension of the feature space, or the number of predictors).

**Given:** a dataset  $(\mathbf{X}, y)_{data} = [(\mathbf{X}_{train}, y_{train}), (\mathbf{X}_{test}, y_{test})] \sim p_{data}$

**Architecture: Linear**

$$\hat{\mathbf{y}} = \mathbf{X} \mathbf{W}$$



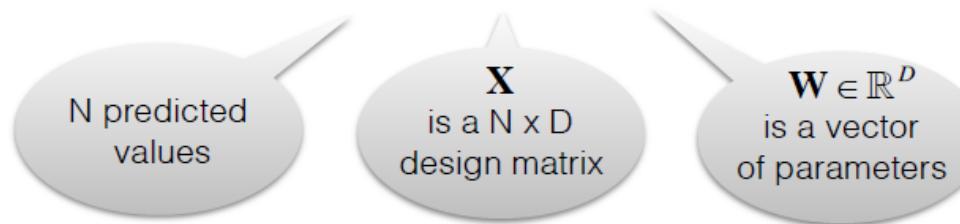
# Linear Regression as a ML Task

**Task:** predict a scalar value  $y \in \mathbb{R}$  from a vector of predictors (“**features**”)  $\mathbf{X} = (X_0, X_1, \dots, X_{D-1}) \in \mathbb{R}^D$  ( $D \geq 1$  is the dimension of the feature space, or the number of predictors).

**Given:** a dataset  $(\mathbf{X}, y)_{data} = [(\mathbf{X}_{train}, y_{train}), (\mathbf{X}_{test}, y_{test})] \sim p_{data}$

**Architecture: Linear**

$$\hat{\mathbf{y}} = \mathbf{X} \mathbf{W}$$



**Example:**  $\mathbf{y}$  is a vector of  $N$  daily returns of AMZN, and  $\mathbf{X}$  is a  $N \times D$  design matrix made of  $D - 1$  daily market index returns (S&P 500, NASDAQ, VIX, etc.)

# Learning in Linear Regression

# Learning in Linear Regression

**Performance measure P :** mean square error (**MSE**) on the **test set**:

$$\text{MSE}_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} (\hat{y}_i^{test} - y_i^{test})^2$$

# Learning in Linear Regression

**Performance measure P :** mean square error (**MSE**) on the **test set**:

$$\text{MSE}_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} (\hat{y}_i^{test} - y_i^{test})^2 = \frac{1}{N_{test}} \left\| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right\|_2^2$$

# Learning in Linear Regression

**Performance measure P :** mean square error (**MSE**) on the **test set**:

$$\text{MSE}_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} (\hat{y}_i^{test} - y_i^{test})^2 = \frac{1}{N_{test}} \left\| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right\|_2^2 = \frac{1}{N_{test}} \left\| \mathbf{X}^{test} \mathbf{W} - \mathbf{Y}^{test} \right\|_2^2$$

# Learning in Linear Regression

**Performance measure P :** mean square error (**MSE**) on the **test set**:

$$\text{MSE}_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} (\hat{y}_i^{test} - y_i^{test})^2 = \frac{1}{N_{test}} \left\| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right\|_2^2 = \frac{1}{N_{test}} \left\| \mathbf{X}^{test} \mathbf{W} - \mathbf{Y}^{test} \right\|_2^2$$

How to optimize parameters  $\mathbf{W} \in \mathbb{R}^D$  ?

# Learning in Linear Regression

**Performance measure P :** mean square error (**MSE**) on the **test set**:

$$\text{MSE}_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} (\hat{y}_i^{test} - y_i^{test})^2 = \frac{1}{N_{test}} \left\| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right\|_2^2 = \frac{1}{N_{test}} \left\| \mathbf{X}^{test} \mathbf{W} - \mathbf{Y}^{test} \right\|_2^2$$

How to optimize parameters  $\mathbf{W} \in \mathbb{R}^D$ ?

Minimize MSE on the **training set**!

Why? Because both  $\text{MSE}_{train}$  and  $\text{MSE}_{test}$  estimate the same generalization (expected) error  $\mathbb{E}[(\hat{\mathbf{y}} - \mathbf{y})^2]$  from the empirical distribution  $\sim p_{data}$ .

*Note:* In ML, we often minimize one function, while actually caring about minimization of another function. This makes ML different from optimization that typically focuses only on one function.

# Learning in Linear Regression

**Performance measure P :** mean square error (**MSE**) on the **test set**:

$$\text{MSE}_{\text{test}} = \frac{1}{N_{\text{test}}} \sum_{n=1}^{N_{\text{test}}} (\hat{y}_i^{\text{test}} - y_i^{\text{test}})^2 = \frac{1}{N_{\text{test}}} \left\| \hat{\mathbf{Y}}^{\text{test}} - \mathbf{Y}^{\text{test}} \right\|_2^2 = \frac{1}{N_{\text{test}}} \left\| \mathbf{X}^{\text{test}} \mathbf{W} - \mathbf{Y}^{\text{test}} \right\|_2^2$$

How to optimize parameters  $\mathbf{W} \in \mathbb{R}^D$ ? Minimize MSE on the **training set**! Why? Because both  $\text{MSE}_{\text{test}}$  and  $\text{MSE}_{\text{train}}$  estimate the generalization (expected) error.

For optimal weights  $\mathbf{W}$ , the gradient of  $\text{MSE}_{\text{train}}$  should be 0:

$$\nabla_{\mathbf{W}} \text{MSE}_{\text{train}} = \nabla_{\mathbf{W}} \frac{1}{N_{\text{train}}} \left\| \hat{\mathbf{Y}}^{\text{train}} - \mathbf{Y}^{\text{train}} \right\|_2^2 = \frac{1}{N_{\text{train}}} \nabla_{\mathbf{W}} \left\| \mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{Y}^{\text{train}} \right\|_2^2 = 0$$

# Learning in Linear Regression

**Performance measure P :** mean square error (**MSE**) on the **test set**:

$$\text{MSE}_{\text{test}} = \frac{1}{N_{\text{test}}} \sum_{n=1}^{N_{\text{test}}} (\hat{y}_i^{\text{test}} - y_i^{\text{test}})^2 = \frac{1}{N_{\text{test}}} \left\| \hat{\mathbf{Y}}^{\text{test}} - \mathbf{Y}^{\text{test}} \right\|_2^2 = \frac{1}{N_{\text{test}}} \left\| \mathbf{X}^{\text{test}} \mathbf{W} - \mathbf{Y}^{\text{test}} \right\|_2^2$$

How to optimize parameters  $\mathbf{W} \in \mathbb{R}^D$ ? Minimize MSE on the **training set**! Why? Because both  $\text{MSE}_{\text{test}}$  and  $\text{MSE}_{\text{train}}$  estimate the generalization (expected) error.

For optimal weights  $\mathbf{W}$ , the gradient of  $\text{MSE}_{\text{train}}$  should be 0:

$$\begin{aligned} \nabla_{\mathbf{W}} \text{MSE}_{\text{train}} &= \nabla_{\mathbf{W}} \frac{1}{N_{\text{train}}} \left\| \hat{\mathbf{Y}}^{\text{train}} - \mathbf{Y}^{\text{train}} \right\|_2^2 = \frac{1}{N_{\text{train}}} \nabla_{\mathbf{W}} \left\| \mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{Y}^{\text{train}} \right\|_2^2 = 0 \\ \Rightarrow \nabla_{\mathbf{W}} (\mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{y}^{\text{train}})^T (\mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{y}^{\text{train}}) &= 0 \end{aligned}$$

# Learning in Linear Regression

**Performance measure P :** mean square error (**MSE**) on the **test set**:

$$\text{MSE}_{\text{test}} = \frac{1}{N_{\text{test}}} \sum_{n=1}^{N_{\text{test}}} (\hat{y}_i^{\text{test}} - y_i^{\text{test}})^2 = \frac{1}{N_{\text{test}}} \left\| \hat{\mathbf{Y}}^{\text{test}} - \mathbf{Y}^{\text{test}} \right\|_2^2 = \frac{1}{N_{\text{test}}} \left\| \mathbf{X}^{\text{test}} \mathbf{W} - \mathbf{Y}^{\text{test}} \right\|_2^2$$

How to optimize parameters  $\mathbf{W} \in \mathbb{R}^D$ ? Minimize MSE on the **training set**! Why? Because both  $\text{MSE}_{\text{test}}$  and  $\text{MSE}_{\text{train}}$  estimate the generalization (expected) error.

For optimal weights  $\mathbf{W}$ , the gradient of  $\text{MSE}_{\text{train}}$  should be 0:

$$\nabla_{\mathbf{W}} \text{MSE}_{\text{train}} = \nabla_{\mathbf{W}} \frac{1}{N_{\text{train}}} \left\| \hat{\mathbf{Y}}^{\text{train}} - \mathbf{Y}^{\text{train}} \right\|_2^2 = \frac{1}{N_{\text{train}}} \nabla_{\mathbf{W}} \left\| \mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{Y}^{\text{train}} \right\|_2^2 = 0$$

$$\Rightarrow \nabla_{\mathbf{W}} (\mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{y}^{\text{train}})^T (\mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{y}^{\text{train}}) = 0$$

Simplify notation:  
 $\mathbf{X}^{\text{train}} \rightarrow \mathbf{X}, \mathbf{y}^{\text{train}} \rightarrow \mathbf{y}$

# Learning in Linear Regression

**Performance measure P :** mean square error (**MSE**) on the **test set**:

$$\text{MSE}_{\text{test}} = \frac{1}{N_{\text{test}}} \sum_{n=1}^{N_{\text{test}}} (\hat{y}_i^{\text{test}} - y_i^{\text{test}})^2 = \frac{1}{N_{\text{test}}} \left\| \hat{\mathbf{Y}}^{\text{test}} - \mathbf{Y}^{\text{test}} \right\|_2^2 = \frac{1}{N_{\text{test}}} \left\| \mathbf{X}^{\text{test}} \mathbf{W} - \mathbf{Y}^{\text{test}} \right\|_2^2$$

How to optimize parameters  $\mathbf{W} \in \mathbb{R}^D$ ? Minimize MSE on the **training set**! Why? Because both  $\text{MSE}_{\text{test}}$  and  $\text{MSE}_{\text{train}}$  estimate the generalization (expected) error.

For optimal weights  $\mathbf{W}$ , the gradient of  $\text{MSE}_{\text{train}}$  should be 0:

$$\nabla_{\mathbf{W}} \text{MSE}_{\text{train}} = \nabla_{\mathbf{W}} \frac{1}{N_{\text{train}}} \left\| \hat{\mathbf{Y}}^{\text{train}} - \mathbf{Y}^{\text{train}} \right\|_2^2 = \frac{1}{N_{\text{train}}} \nabla_{\mathbf{W}} \left\| \mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{Y}^{\text{train}} \right\|_2^2 = 0$$

$$\Rightarrow \nabla_{\mathbf{W}} (\mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{y}^{\text{train}})^T (\mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{y}^{\text{train}}) = 0$$

$$\Rightarrow \nabla_{\mathbf{W}} (\mathbf{W}^T \mathbf{X}^T \mathbf{X} \mathbf{W} - 2\mathbf{W}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y})$$

Simplify notation:  
 $\mathbf{X}^{\text{train}} \rightarrow \mathbf{X}, \mathbf{y}^{\text{train}} \rightarrow \mathbf{y}$

# Learning in Linear Regression

**Performance measure P :** mean square error (**MSE**) on the **test set**:

$$\text{MSE}_{\text{test}} = \frac{1}{N_{\text{test}}} \sum_{n=1}^{N_{\text{test}}} (\hat{y}_i^{\text{test}} - y_i^{\text{test}})^2 = \frac{1}{N_{\text{test}}} \left\| \hat{\mathbf{Y}}^{\text{test}} - \mathbf{Y}^{\text{test}} \right\|_2^2 = \frac{1}{N_{\text{test}}} \left\| \mathbf{X}^{\text{test}} \mathbf{W} - \mathbf{Y}^{\text{test}} \right\|_2^2$$

How to optimize parameters  $\mathbf{W} \in \mathbb{R}^D$ ? Minimize MSE on the **training set**! Why? Because both  $\text{MSE}_{\text{test}}$  and  $\text{MSE}_{\text{train}}$  estimate the generalization (expected) error.

For optimal weights  $\mathbf{W}$ , the gradient of  $\text{MSE}_{\text{train}}$  should be 0:

$$\nabla_{\mathbf{W}} \text{MSE}_{\text{train}} = \nabla_{\mathbf{W}} \frac{1}{N_{\text{train}}} \left\| \hat{\mathbf{Y}}^{\text{train}} - \mathbf{Y}^{\text{train}} \right\|_2^2 = \frac{1}{N_{\text{train}}} \nabla_{\mathbf{W}} \left\| \mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{Y}^{\text{train}} \right\|_2^2 = 0$$

$$\Rightarrow \nabla_{\mathbf{W}} (\mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{y}^{\text{train}})^T (\mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{y}^{\text{train}}) = 0$$

$$\Rightarrow \nabla_{\mathbf{W}} (\mathbf{W}^T \mathbf{X}^T \mathbf{X} \mathbf{W} - 2\mathbf{W}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) = 2(\mathbf{X}^T \mathbf{X}) \mathbf{W} - 2(\mathbf{X}^T \mathbf{y}) = 0$$

Simplify notation:  
 $\mathbf{X}^{\text{train}} \rightarrow \mathbf{X}, \mathbf{y}^{\text{train}} \rightarrow \mathbf{y}$

# Learning in Linear Regression

**Performance measure P :** mean square error (**MSE**) on the **test set**:

$$\text{MSE}_{\text{test}} = \frac{1}{N_{\text{test}}} \sum_{n=1}^{N_{\text{test}}} (\hat{y}_i^{\text{test}} - y_i^{\text{test}})^2 = \frac{1}{N_{\text{test}}} \left\| \hat{\mathbf{Y}}^{\text{test}} - \mathbf{Y}^{\text{test}} \right\|_2^2 = \frac{1}{N_{\text{test}}} \left\| \mathbf{X}^{\text{test}} \mathbf{W} - \mathbf{Y}^{\text{test}} \right\|_2^2$$

How to optimize parameters  $\mathbf{W} \in \mathbb{R}^D$ ? Minimize MSE on the **training set**! Why? Because both  $\text{MSE}_{\text{test}}$  and  $\text{MSE}_{\text{train}}$  estimate the generalization (expected) error.

For optimal weights  $\mathbf{W}$ , the gradient of  $\text{MSE}_{\text{train}}$  should be 0:

$$\nabla_{\mathbf{W}} \text{MSE}_{\text{train}} = \nabla_{\mathbf{W}} \frac{1}{N_{\text{train}}} \left\| \hat{\mathbf{Y}}^{\text{train}} - \mathbf{Y}^{\text{train}} \right\|_2^2 = \frac{1}{N_{\text{train}}} \nabla_{\mathbf{W}} \left\| \mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{Y}^{\text{train}} \right\|_2^2 = 0$$

$$\Rightarrow \nabla_{\mathbf{W}} (\mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{y}^{\text{train}})^T (\mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{y}^{\text{train}}) = 0$$

$$\Rightarrow \nabla_{\mathbf{W}} (\mathbf{W}^T \mathbf{X}^T \mathbf{X} \mathbf{W} - 2\mathbf{W}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) = 2(\mathbf{X}^T \mathbf{X}) \mathbf{W} - 2(\mathbf{X}^T \mathbf{y}) = 0$$

Simplify notation:  
 $\mathbf{X}^{\text{train}} \rightarrow \mathbf{X}, \mathbf{y}^{\text{train}} \rightarrow \mathbf{y}$

$\Rightarrow$

$$\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Normal  
equation

# Learning in Linear Regression

**Performance measure P :** mean square error (**MSE**) on the **test set**:

$$\text{MSE}_{\text{test}} = \frac{1}{N_{\text{test}}} \sum_{n=1}^{N_{\text{test}}} (\hat{y}_i^{\text{test}} - y_i^{\text{test}})^2 = \frac{1}{N_{\text{test}}} \left\| \hat{\mathbf{Y}}^{\text{test}} - \mathbf{Y}^{\text{test}} \right\|_2^2 = \frac{1}{N_{\text{test}}} \left\| \mathbf{X}^{\text{test}} \mathbf{W} - \mathbf{Y}^{\text{test}} \right\|_2^2$$

How to optimize parameters  $\mathbf{W} \in \mathbb{R}^D$ ? Minimize MSE on the **training set**! Why? Because both  $\text{MSE}_{\text{test}}$  and  $\text{MSE}_{\text{train}}$  estimate the generalization (expected) error.

For optimal weights  $\mathbf{W}$ , the gradient of  $\text{MSE}_{\text{train}}$  should be 0:

$$\nabla_{\mathbf{W}} \text{MSE}_{\text{train}} = \nabla_{\mathbf{W}} \frac{1}{N_{\text{train}}} \left\| \hat{\mathbf{Y}}^{\text{train}} - \mathbf{Y}^{\text{train}} \right\|_2^2 = \frac{1}{N_{\text{train}}} \nabla_{\mathbf{W}} \left\| \mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{Y}^{\text{train}} \right\|_2^2 = 0$$

$$\Rightarrow \nabla_{\mathbf{W}} (\mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{y}^{\text{train}})^T (\mathbf{X}^{\text{train}} \mathbf{W} - \mathbf{y}^{\text{train}}) = 0$$

$$\Rightarrow \nabla_{\mathbf{W}} (\mathbf{W}^T \mathbf{X}^T \mathbf{X} \mathbf{W} - 2\mathbf{W}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) = 2(\mathbf{X}^T \mathbf{X}) \mathbf{W} - 2(\mathbf{X}^T \mathbf{y}) = 0$$

Simplify notation:  
 $\mathbf{X}^{\text{train}} \rightarrow \mathbf{X}, \mathbf{y}^{\text{train}} \rightarrow \mathbf{y}$

Here  
 $\mathbf{X} = \mathbf{X}^{\text{train}}, \mathbf{y} = \mathbf{y}^{\text{train}}$

$\Rightarrow$

$$\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Normal equation

# Learning in Linear Regression

# Learning in Linear Regression

Visualize the result for optimal weights  $\mathbf{W}$  in the vector form

$$\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$
$$\mathbf{W} = \begin{bmatrix} \mathbf{X}^T & \mathbf{X} \end{bmatrix}^{-1} \mathbf{X}^T \mathbf{y}$$

Matrix D x 1      Matrix N x 1

Matrix D x D      Matrix D x N

Prediction  
in-sample:

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{W} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \equiv \mathbf{H}\mathbf{y}$$

$$= \mathbf{H}$$

The “hat” (projection) matrix

# Learning in Linear Regression

Visualize the result for optimal weights  $\mathbf{W}$  in the vector form

$$\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\mathbf{W} = \left[ \begin{array}{c|c} \mathbf{X}^T & \mathbf{X} \end{array} \right]^{-1} \mathbf{X}^T \mathbf{y}$$

Matrix D x 1      Matrix N x 1

Matrix D x D      Matrix D x N

Sir R. Penrose

Nobel Prize in Physics 2020

Prediction  
in-sample



$$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \equiv \mathbf{H} \mathbf{y}$$

$$= \mathbf{H}$$

The “hat” (projection) matrix

# Learning in Linear Regression

Visualize the result for optimal weights  $\mathbf{W}$  in the vector form

$$\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Sir R. Penrose

# Nobel Prize in Physics 2020

## Prediction in-sample



$$\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \equiv \mathbf{H} \mathbf{y}$$

Symmetric and  
idempotent:  
 $\mathbf{H}^2 = \mathbf{H}$

## The “hat” (projection) matrix

# Learning in Linear Regression

Visualize the result for optimal weights  $\mathbf{W}$  in the vector form

$$\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

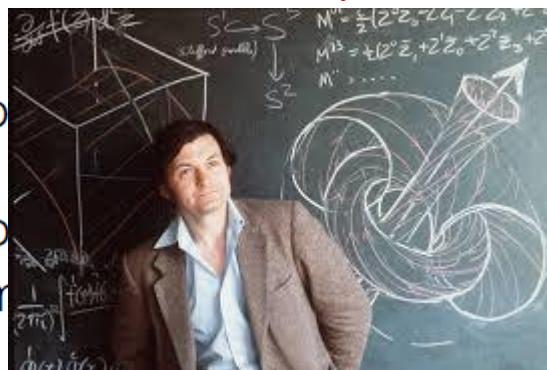
$$\mathbf{W} = \left[ \begin{array}{c|c} \mathbf{X}^T & \mathbf{X} \end{array} \right]^{-1} \mathbf{X}^T \mathbf{y}$$

Matrix D x 1      Matrix N x 1  
Matrix D x D      Matrix D x N

Sir R. Penrose

Nobel Prize in Physics 2020

Prediction



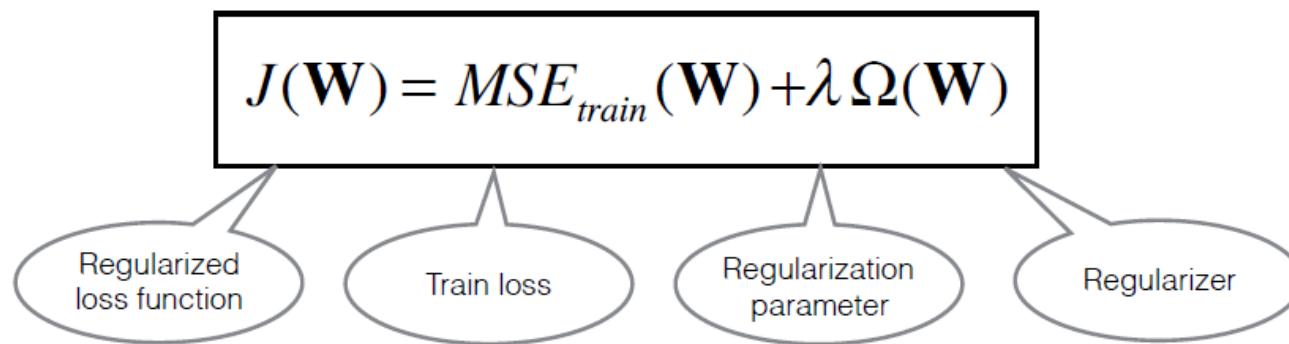
$$\mathbf{X}^{test} \mathbf{W}$$

How to obtain a better fit?  
design matrix: add more features to the control!

# Regularization

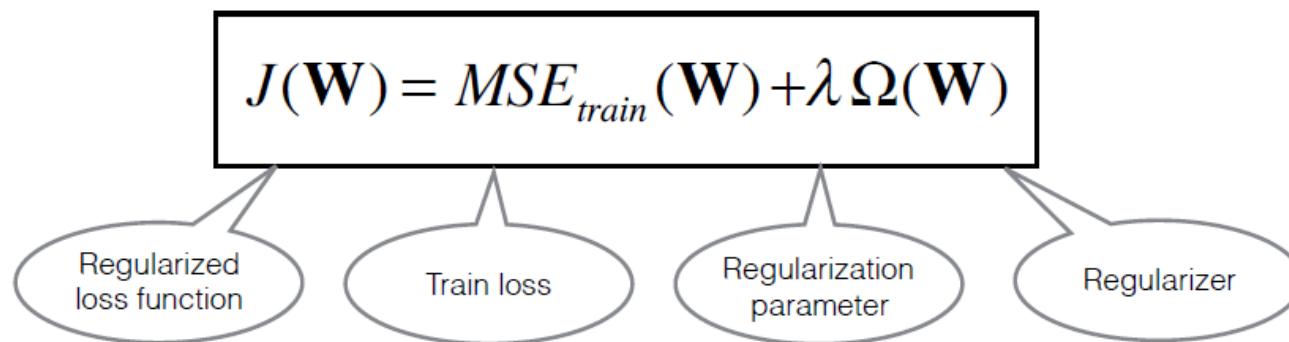
# Regularization

Recall that we minimize  $MSE_{train}$ , though want in fact to minimize  $MSE_{test}$ .  
The idea of regularization is to modify the objective function of  
minimization of  $MSE_{train}$  so that  $MSE_{test}$  has a smaller variance:



# Regularization

Recall that we minimize  $MSE_{train}$ , though want in fact to minimize  $MSE_{test}$ .  
The idea of regularization is to modify the objective function of minimization of  $MSE_{train}$  so that  $MSE_{test}$  has a smaller variance:



Popular choices for the **regularizer**:

- $L_2$  regularization     $\Omega(\mathbf{W}) = \mathbf{W}^T \mathbf{W} = \|\mathbf{W}\|_2$     Penalties large weights

- $L_1$  regularization     $\Omega(\mathbf{W}) = \sum_i |W_i| = \|\mathbf{W}\|_1$     Enforces a sparse solution

- Entropy regularization     $\Omega(\mathbf{W}) = \sum_i W_i \log W_i$     ( $W_i \geq 0, \sum_i W_i = 1$ )

Motivated by Bayesian statistics

# Hyperparameters and Validation Set

# Hyperparameters and Validation Set

- **Hyperparameters** are any quantitative features of ML models that are not directly optimized by minimizing in-sample loss such as  $\text{MSE}_{train}$
- **Hyperparameters control model capacity**

# Hyperparameters and Validation Set

- **Hyperparameters** are any quantitative features of ML models that are not directly optimized by minimizing in-sample loss such as  $\text{MSE}_{train}$
- **Hyperparameters control model capacity**
- Examples of hyperparameters:
  - ◆ Degree of a polynomial regression (linear, quadratic, cubic, etc.)
  - ◆ Regularization parameter  $\lambda$
  - ◆ Number of levels in a decision tree
  - ◆ Number of layers and nodes per layer in neural networks
  - ◆ Learning rate
  - ◆ ...

# Hyperparameters and Validation Set

- **Hyperparameters** are any quantitative features of ML models that are not directly optimized by minimizing in-sample loss such as  $\text{MSE}_{\text{train}}$
- **Hyperparameters control model capacity**
- Examples of hyperparameters:
  - ◆ Degree of a polynomial regression (linear, quadratic, cubic, etc.)
  - ◆ Regularization parameter  $\lambda$
  - ◆ Number of levels in a decision tree
  - ◆ Number of layers and nodes per layer in neural networks
  - ◆ Learning rate
  - ◆ ...
- How to choose hyperparameters:
  - ◆ Split a training set into training and **validation** sets (e.g. as 80:20)

# Hyperparameters and Validation Set

- **Hyperparameters** are any quantitative features of ML models that are not directly optimized by minimizing in-sample loss such as  $\text{MSE}_{\text{train}}$
- **Hyperparameters control model capacity**
- Examples of hyperparameters:
  - ◆ Degree of a polynomial regression (linear, quadratic, cubic, etc.)
  - ◆ Regularization parameter  $\lambda$
  - ◆ Number of levels in a decision tree
  - ◆ Number of layers and nodes per layer in neural networks
  - ◆ Learning rate
  - ◆ ...
- How to choose hyperparameters:
  - ◆ Split a training set into training and **validation** sets (e.g. as 80:20)
  - ◆ Use **cross-validation**

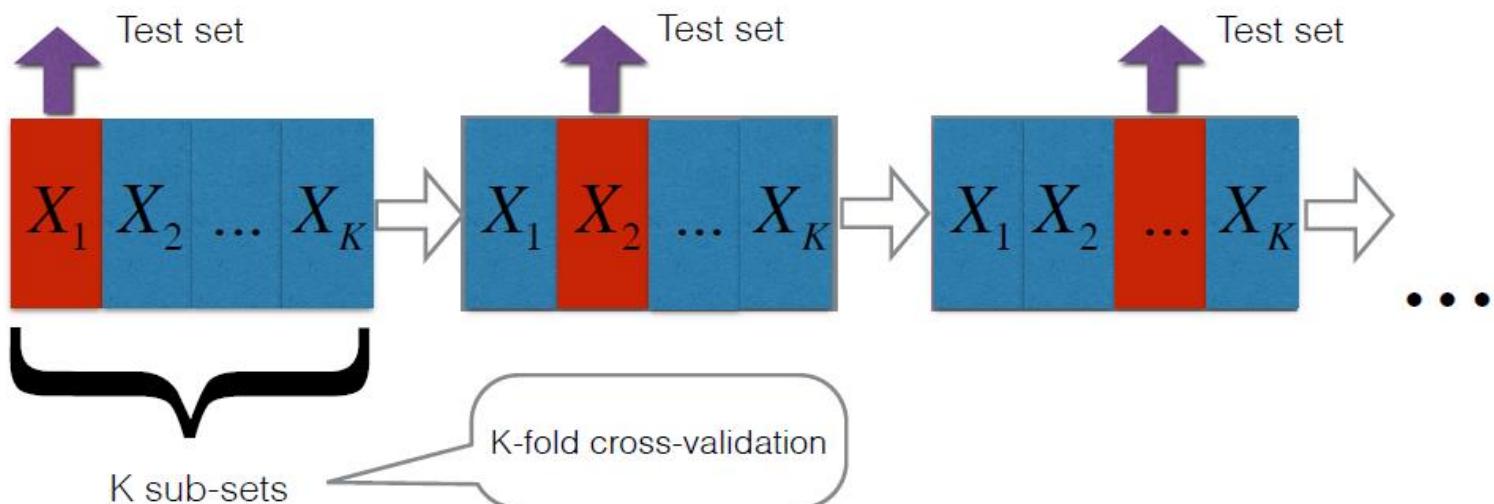
# Cross-Validation

# Cross-Validation

- Assume we are given  $N$  samples, but  $N$  is small, so setting aside a fixed test set is problematic.
- We want to use all samples for training!
- This is achieved using cross-validation:

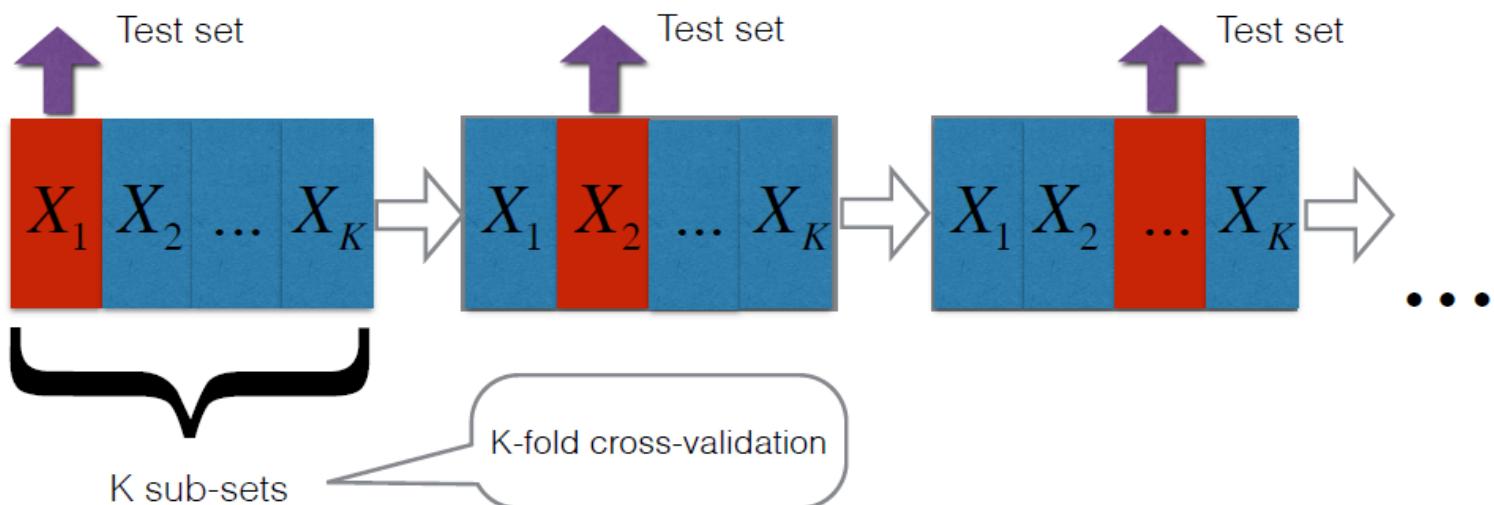
# Cross-Validation

- Assume we are given  $N$  samples, but  $N$  is small, so setting aside a fixed test set is problematic.
- We want to use all samples for training!
- This is achieved using cross-validation:



# Cross-Validation

- Assume we are given  $N$  samples, but  $N$  is small, so setting aside a fixed test set is problematic.
- We want to use all samples for training!
- This is achieved using cross-validation:

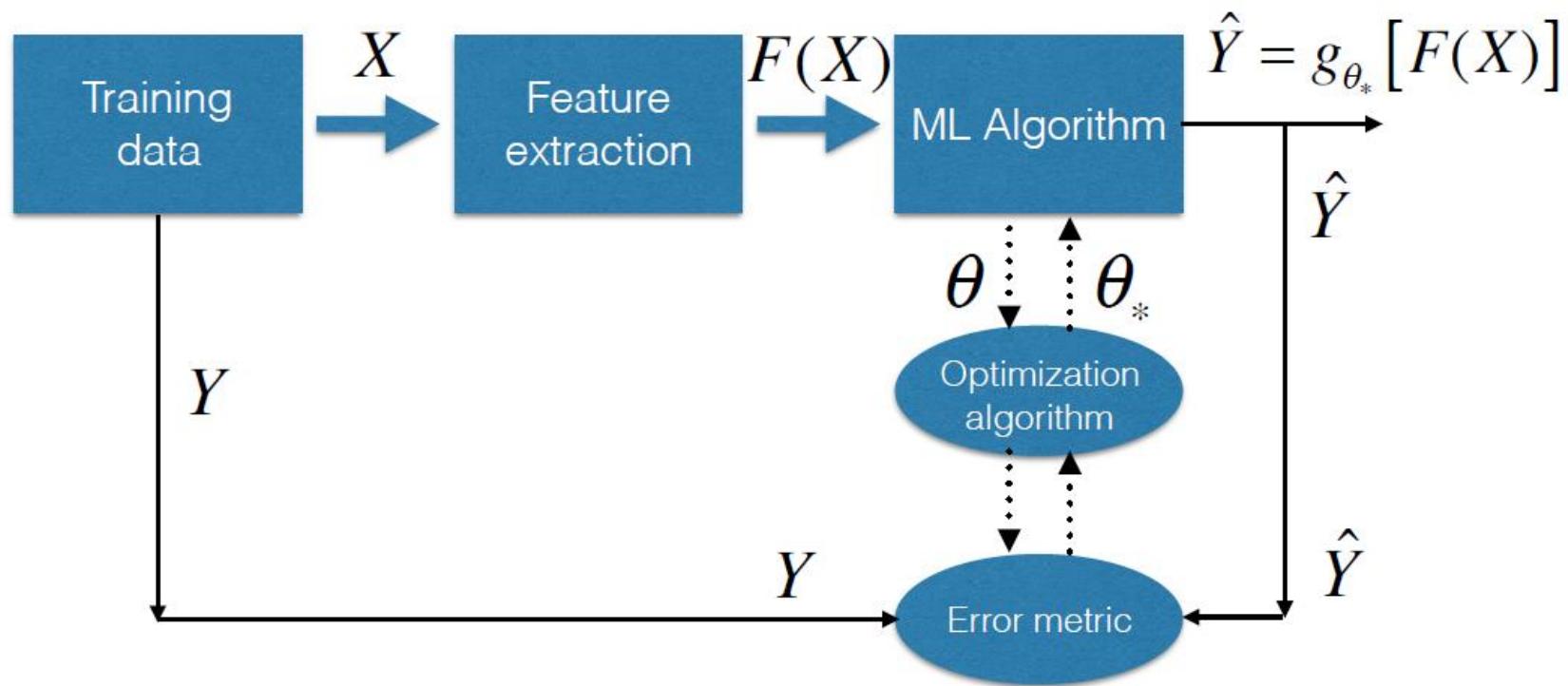


Special cases:

- $K=1$ : no test sub-set!
- $K=N$ : leave-one-out cross-validation

# Supervised Learning Diagram

General diagram for training supervised learning algorithms



# Analytical Learning in Linear Reg

# Analytical Learning in Linear Reg

**Performance measure P :** mean square error (**MSE**)  $\text{MSE}_{test}$ :

$$\text{MSE}_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} (\hat{y}_i^{test} - y_i^{test})^2 = \frac{1}{N_{test}} \left\| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right\|_2^2 = \frac{1}{N_{test}} \left\| \mathbf{X}^{test} \mathbf{W} - \mathbf{Y}^{test} \right\|_2^2$$

Parameters  $\mathbf{W} \in \mathbb{R}^D$  are found by minimizing  $\text{MSE}_{train}$

$$\text{MSE}_{train} = \frac{1}{N_{train}} \sum_{n=1}^{N_{train}} (\hat{y}_i^{train} - y_i^{train})^2 = \frac{1}{N_{train}} \left\| \hat{\mathbf{Y}}^{train} - \mathbf{Y}^{train} \right\|_2^2 = \frac{1}{N_{train}} \left\| \mathbf{X}^{train} \mathbf{W} - \mathbf{Y}^{train} \right\|_2^2$$

For optimal weights  $\mathbf{W}$ , the gradient of  $\text{MSE}_{train}$  should be 0:

$$\nabla_{\mathbf{W}} \text{MSE}_{train} = \nabla_{\mathbf{W}} \frac{1}{N_{train}} \left\| \hat{\mathbf{Y}}^{train} - \mathbf{Y}^{train} \right\|_2^2 = \frac{1}{N_{train}} \nabla_{\mathbf{W}} \left\| \mathbf{X}^{train} \mathbf{W} - \mathbf{Y}^{train} \right\|_2^2 = 0$$

$\Rightarrow$

$$\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Normal  
equation

# Analytical Learning in Linear Reg

**Performance measure P :** mean square error (**MSE**)  $MSE_{test}$ :

$$MSE_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} (\hat{y}_i^{test} - y_i^{test})^2 = \frac{1}{N_{test}} \left\| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right\|_2^2 = \frac{1}{N_{test}} \left\| \mathbf{X}^{test} \mathbf{W} - \mathbf{Y}^{test} \right\|_2^2$$

Parameters  $\mathbf{W} \in \mathbb{R}^D$  are found by minimizing  $MSE_{train}$

$$MSE_{train} = \frac{1}{N_{train}} \sum_{n=1}^{N_{train}} (\hat{y}_i^{train} - y_i^{train})^2 = \frac{1}{N_{train}} \left\| \hat{\mathbf{Y}}^{train} - \mathbf{Y}^{train} \right\|_2^2 = \frac{1}{N_{train}} \left\| \mathbf{X}^{train} \mathbf{W} - \mathbf{Y}^{train} \right\|_2^2$$

For optimal weights  $\mathbf{W}$ , the gradient of  $MSE_{train}$  should be 0:

$$\nabla_{\mathbf{W}} MSE_{train} = \nabla_{\mathbf{W}} \frac{1}{N_{train}} \left\| \hat{\mathbf{Y}}^{train} - \mathbf{Y}^{train} \right\|_2^2 = \frac{1}{N_{train}} \nabla_{\mathbf{W}} \left\| \mathbf{X}^{train} \mathbf{W} - \mathbf{Y}^{train} \right\|_2^2 = 0$$

$\Rightarrow$

$$\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Normal equation

Only works analytically for linear regression

# Analytical Learning in Linear Reg

**Performance measure P :** mean square error (**MSE**)  $MSE_{test}$ :

$$MSE_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} (\hat{y}_i^{test} - y_i^{test})^2 = \frac{1}{N_{test}} \left\| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right\|_2^2 = \frac{1}{N_{test}} \left\| \mathbf{X}^{test} \mathbf{W} - \mathbf{Y}^{test} \right\|_2^2$$

Parameters  $\mathbf{W} \in \mathbb{R}^D$  are found by minimizing  $MSE_{train}$

$$MSE_{train} = \frac{1}{N_{train}} \sum_{n=1}^{N_{train}} (\hat{y}_i^{train} - y_i^{train})^2 = \frac{1}{N_{train}} \left\| \hat{\mathbf{Y}}^{train} - \mathbf{Y}^{train} \right\|_2^2 = \frac{1}{N_{train}} \left\| \mathbf{X}^{train} \mathbf{W} - \mathbf{Y}^{train} \right\|_2^2$$

For optimal weights  $\mathbf{W}$ , the gradient of  $MSE_{train}$  should be 0:

$$\nabla_{\mathbf{W}} MSE_{train} = \nabla_{\mathbf{W}} \frac{1}{N_{train}} \left\| \hat{\mathbf{Y}}^{train} - \mathbf{Y}^{train} \right\|_2^2 = \frac{1}{N_{train}} \nabla_{\mathbf{W}} \left\| \mathbf{X}^{train} \mathbf{W} - \mathbf{Y}^{train} \right\|_2^2 = 0$$

$$\Rightarrow \boxed{\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}}$$

Normal equation

Only works analytically for linear regression

Complexity  $O(N^3)$  to  $O(N^{2.4})$

# Analytical Learning in Linear Reg

**Performance measure P :** mean square error (**MSE**)  $MSE_{test}$ :

$$MSE_{test} = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} (\hat{y}_i^{test} - y_i^{test})^2 = \frac{1}{N_{test}} \left\| \hat{\mathbf{Y}}^{test} - \mathbf{Y}^{test} \right\|_2^2 = \frac{1}{N_{test}} \left\| \mathbf{X}^{test} \mathbf{W} - \mathbf{Y}^{test} \right\|_2^2$$

Parameters  $\mathbf{W} \in \mathbb{R}^D$  are found by minimizing  $MSE_{train}$

$$MSE_{train} = \frac{1}{N_{train}} \sum_{n=1}^{N_{train}} (\hat{y}_i^{train} - y_i^{train})^2 = \frac{1}{N_{train}} \left\| \hat{\mathbf{Y}}^{train} - \mathbf{Y}^{train} \right\|_2^2 = \frac{1}{N_{train}} \left\| \mathbf{X}^{train} \mathbf{W} - \mathbf{Y}^{train} \right\|_2^2$$

For optimal weights  $\mathbf{W}$ , the gradient of  $MSE_{train}$  should be 0:

$$\nabla_{\mathbf{W}} MSE_{train} = \nabla_{\mathbf{W}} \frac{1}{N_{train}} \left\| \hat{\mathbf{Y}}^{train} - \mathbf{Y}^{train} \right\|_2^2 = \boxed{\frac{1}{N_{train}} \nabla_{\mathbf{W}} \left\| \mathbf{X}^{train} \mathbf{W} - \mathbf{Y}^{train} \right\|_2^2 = 0}$$

$$\Rightarrow \mathbf{W} = \cancel{\left( \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y}}$$

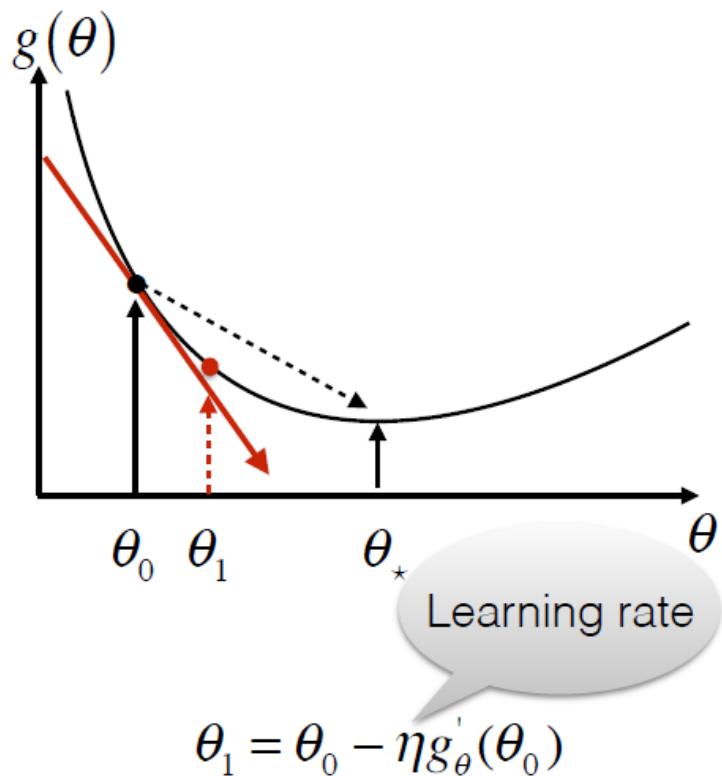
Only works analytically for linear regression

Complexity  $O(N^3)$  to  $O(N^{2.4})$

# Gradient Descent

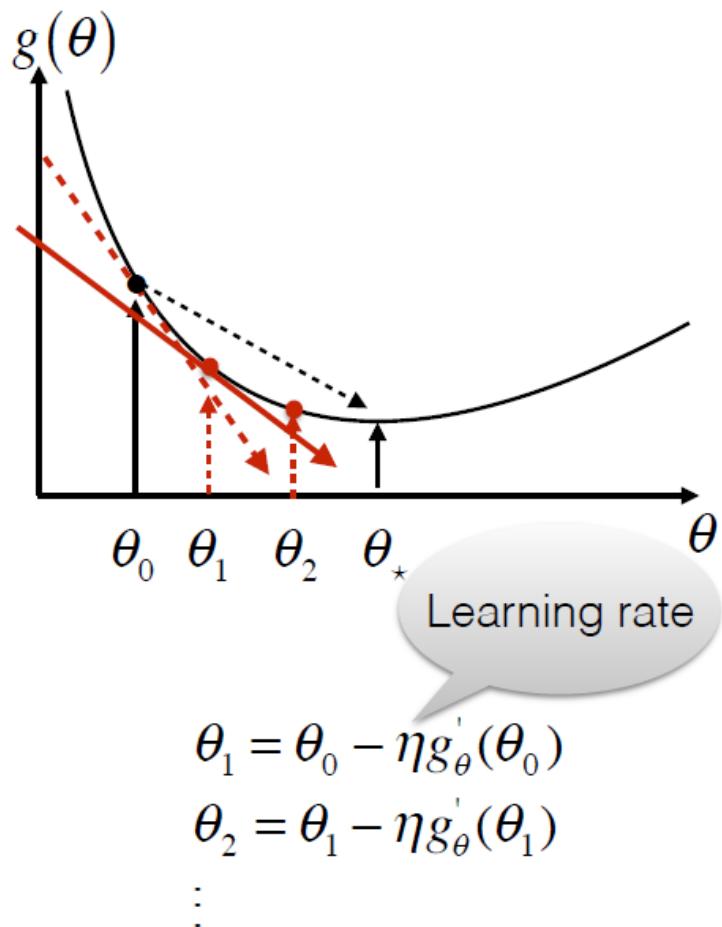
# Gradient Descent

Optimize weights **W** **gradually** by moving against the gradients of  $\text{MSE}_{\text{train}}$



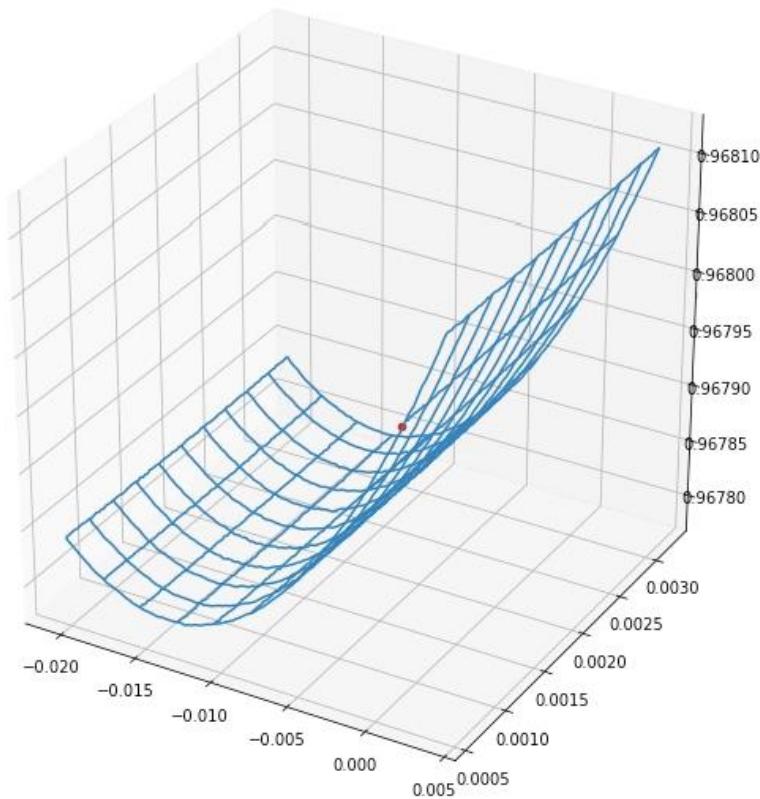
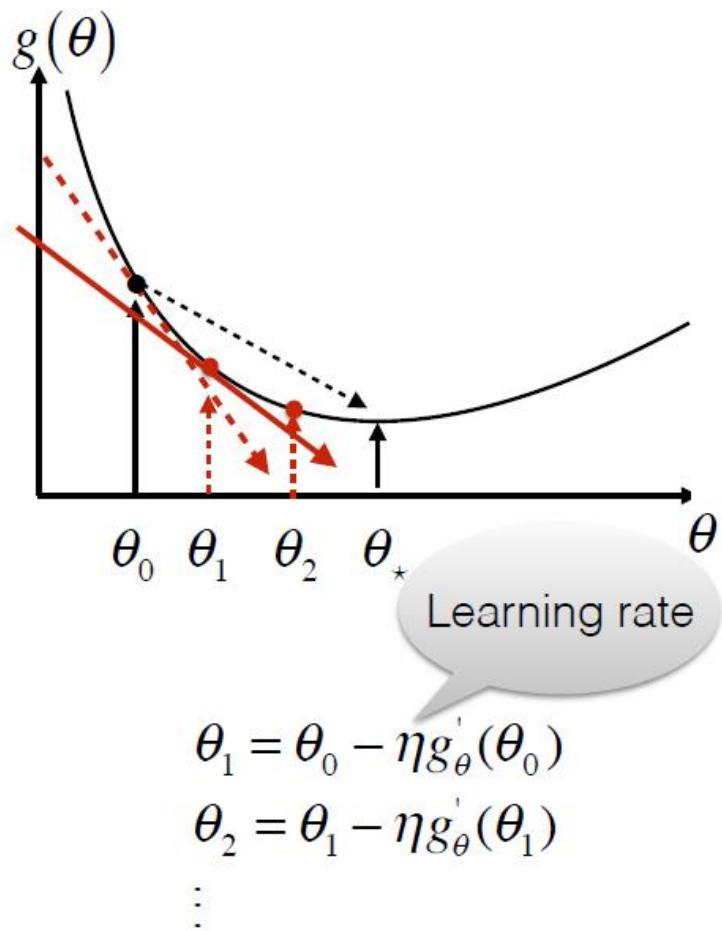
# Gradient Descent

Optimize weights **W** **gradually** by moving against the gradients of  $\text{MSE}_{\text{train}}$



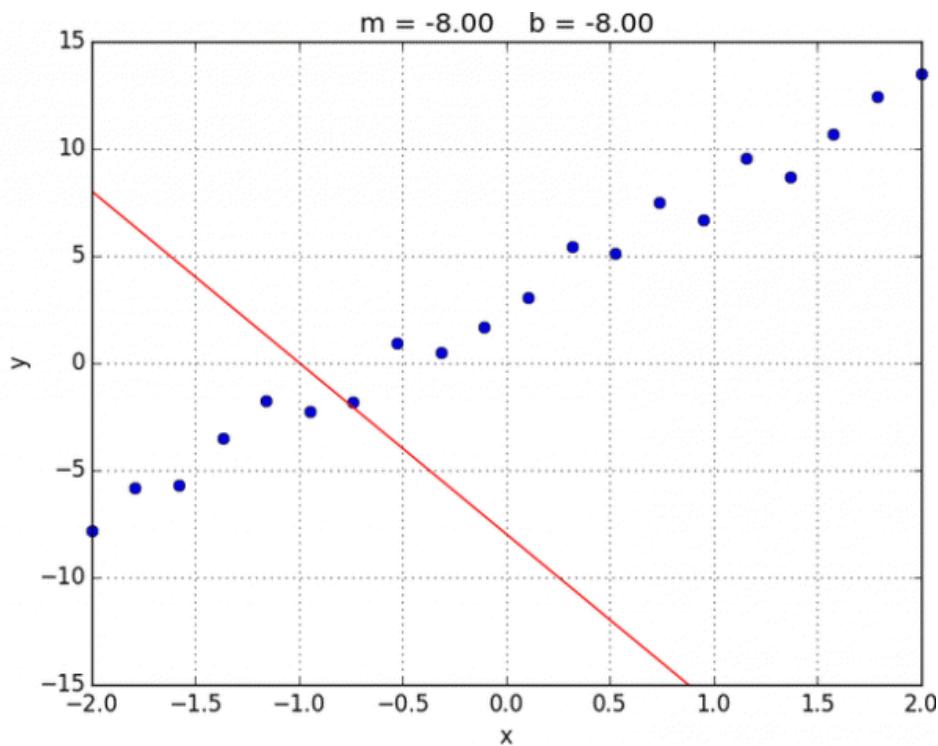
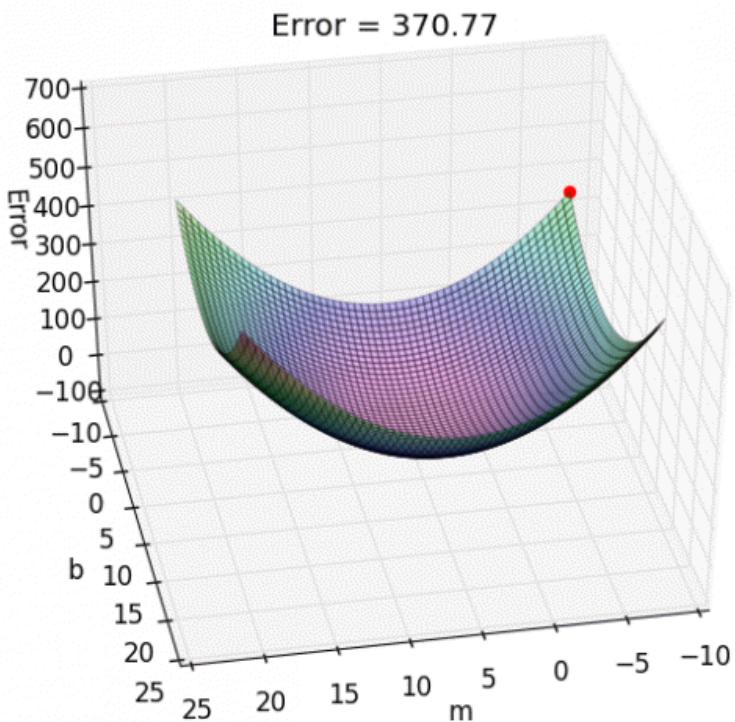
# Gradient Descent

Optimize weights **W** **gradually** by moving against the gradients of  $\text{MSE}_{\text{train}}$



# Gradient Descent

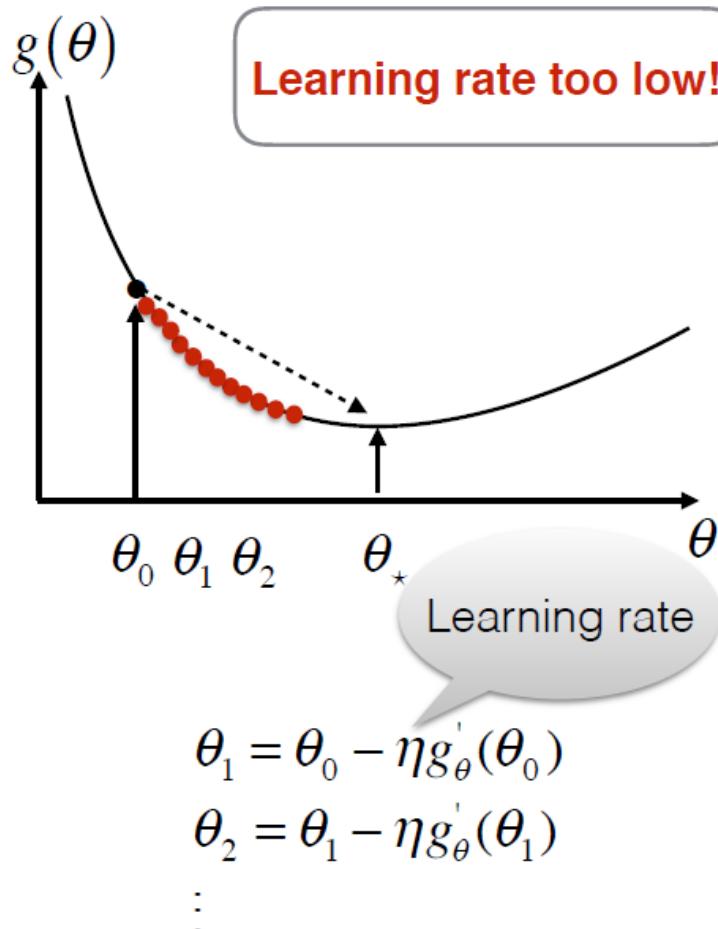
# Gradient Descent



# Choice of Learning Rate

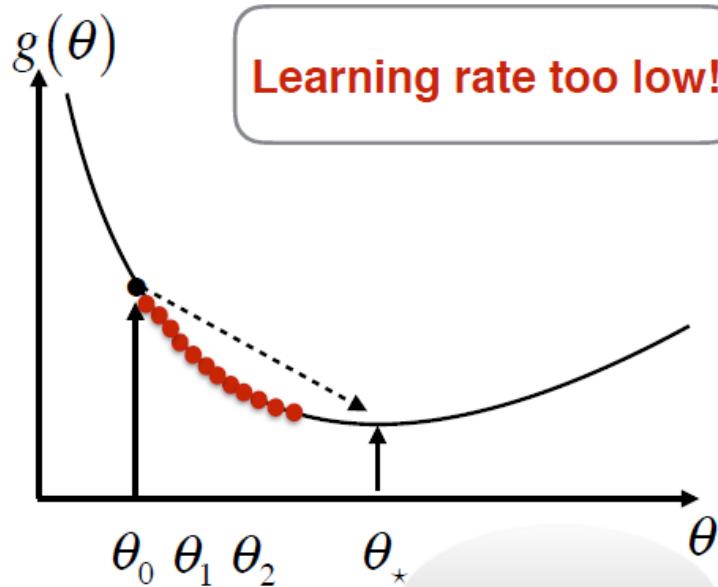
# Choice of Learning Rate

Optimize weights **W gradually** by moving against the gradients of  $\text{MSE}_{\text{train}}$



# Choice of Learning Rate

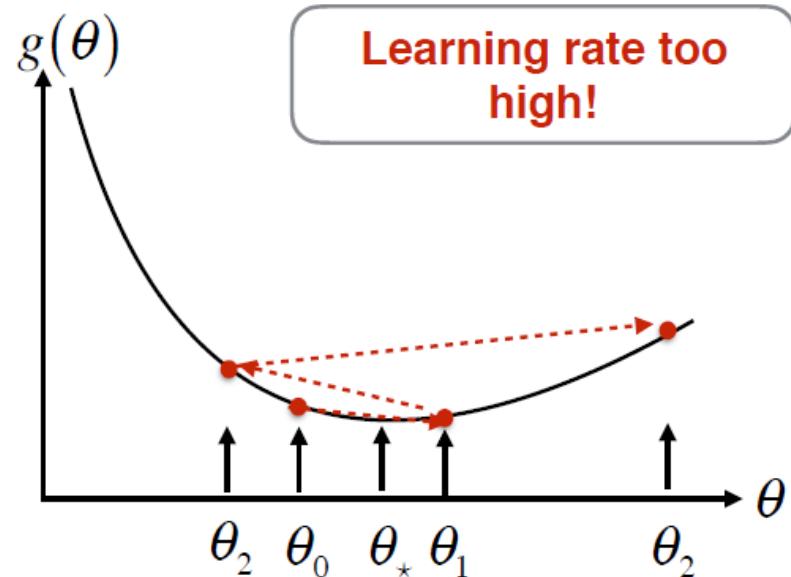
Optimize weights **W gradually** by moving against the gradients of  $\text{MSE}_{\text{train}}$



$$\theta_1 = \theta_0 - \eta g'_\theta(\theta_0)$$

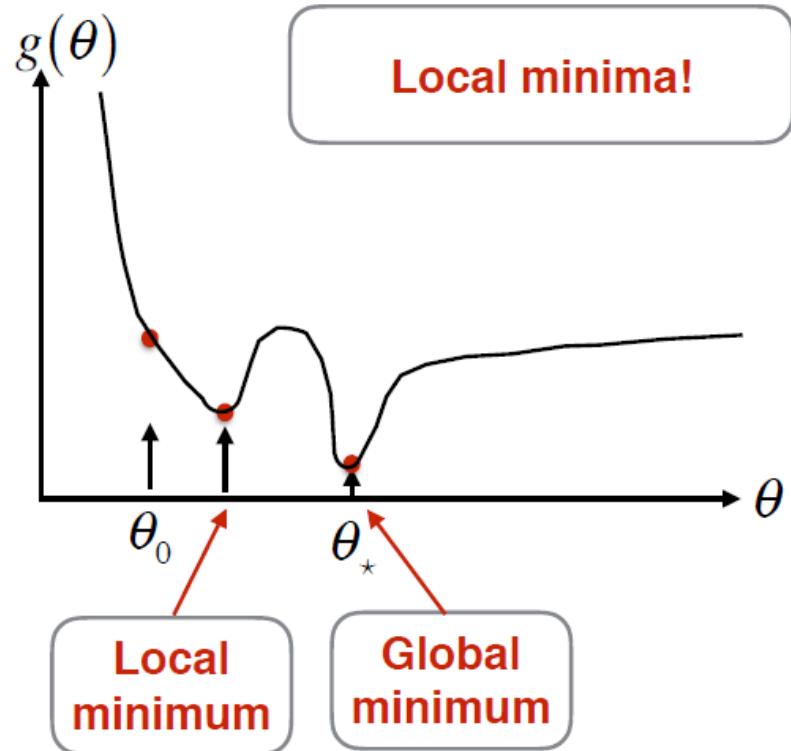
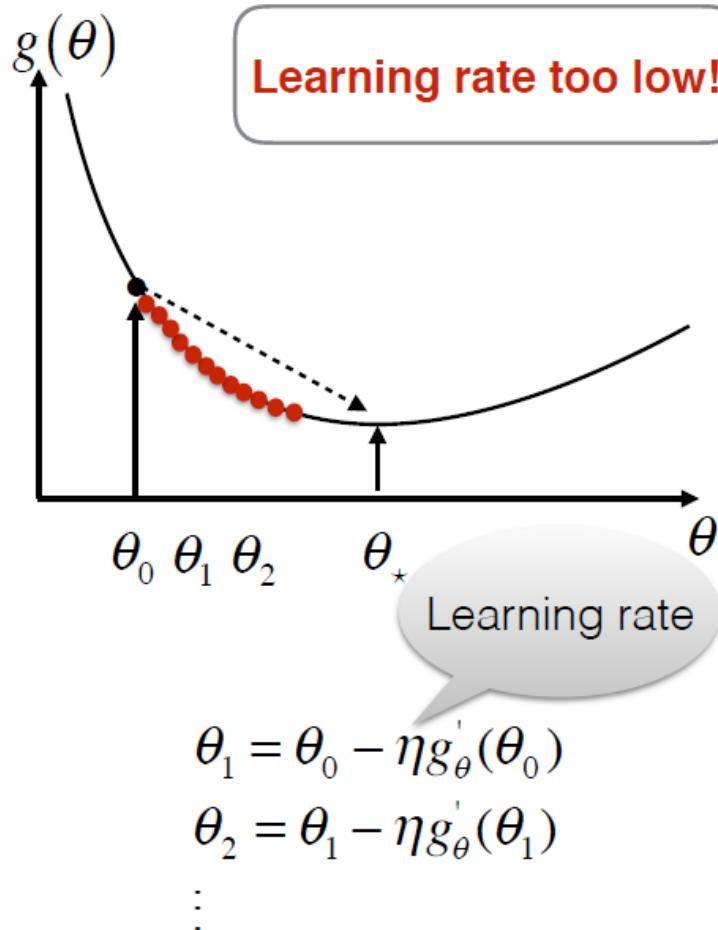
$$\theta_2 = \theta_1 - \eta g'_\theta(\theta_1)$$

⋮



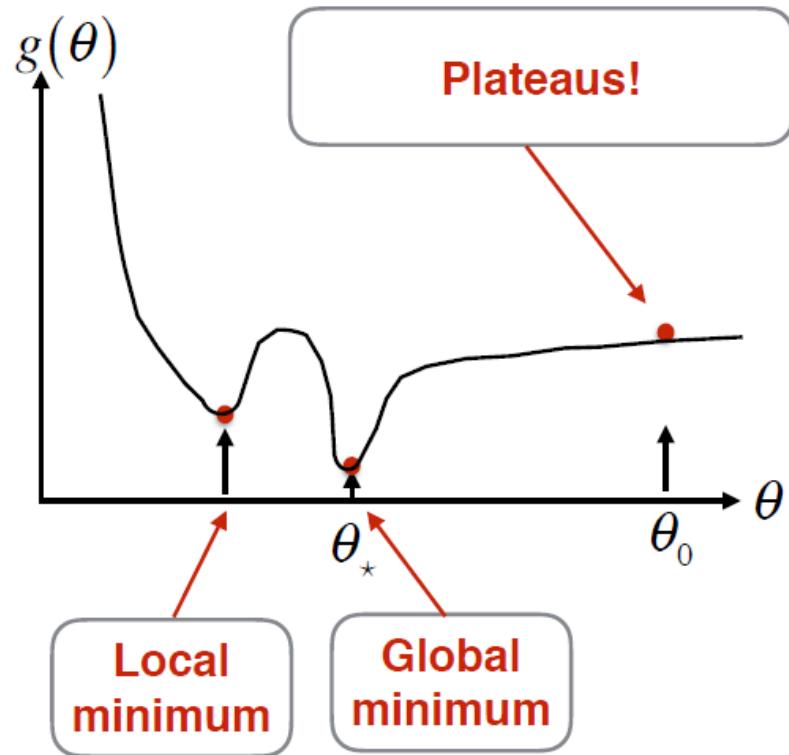
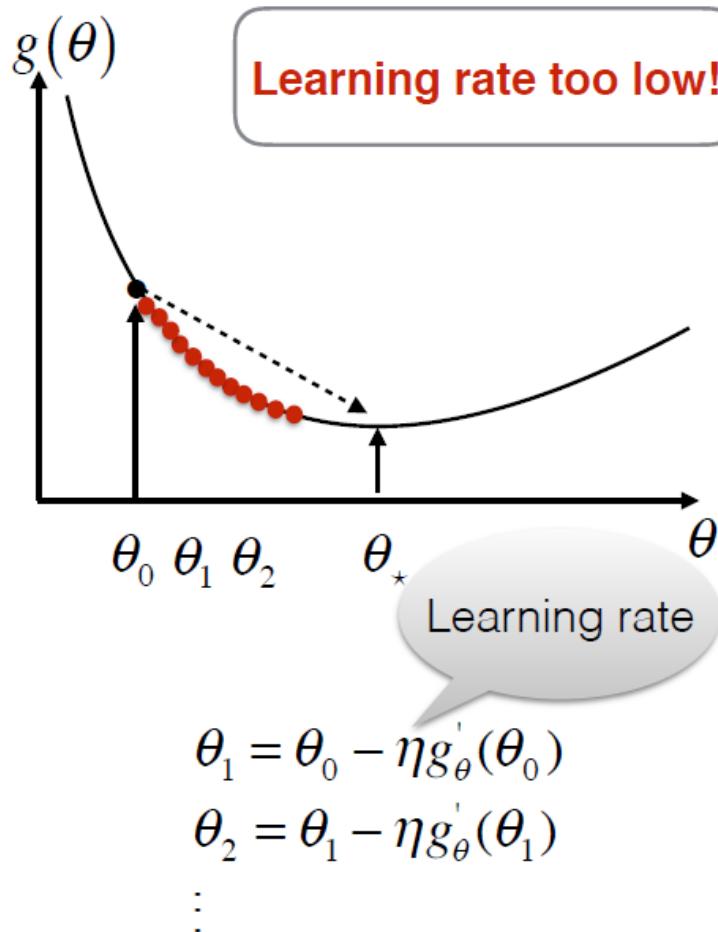
# Choice of Learning Rate

Optimize weights **W gradually** by moving against the gradients of  $\text{MSE}_{\text{train}}$



# Choice of Learning Rate

Optimize weights **W gradually** by moving against the gradients of  $\text{MSE}_{\text{train}}$



# Stochastic Gradient Descent

# Stochastic Gradient Descent

Gradient descent minimizes a train error - an approximation to the generalization error

$$E[f] = \int L(f(x), y) d P(x, y) \Rightarrow E_n[f] = \frac{1}{N} \sum_{i=1}^N L(f(x_i, w), y_i)$$

:

# Stochastic Gradient Descent

Gradient descent minimizes a train error - an approximation to the generalization error

$$E[f] = \int L(f(x), y) d P(x, y) \Rightarrow E_n[f] = \frac{1}{N} \sum_{i=1}^N L(f(x_i, w), y_i)$$

**Gradient Descent:**

$$w^{(k+1)} = w^{(k)} - \eta \frac{1}{N} \sum_{i=1}^N \nabla_w L(f(x_i, w^{(k)}), y_i)$$

:

# Stochastic Gradient Descent

Gradient descent minimizes a train error - an approximation to the generalization error

$$E[f] = \int L(f(x), y) d P(x, y) \Rightarrow E_n[f] = \frac{1}{N} \sum_{i=1}^N L(f(x_i, w), y_i)$$

Gradient Descent:

$$w^{(k+1)} = w^{(k)} - \eta \frac{1}{N} \sum_{i=1}^N \nabla_w L(f(x_i, w^{(k)}), y_i)$$

**Stochastic Gradient Descent (SGD)** with randomly selected **mini-batches**  $M_k$  of some fixed size  $N_{MB} \ll N$ :

$$w^{(k+1)} = w^{(k)} - \eta \frac{1}{N} \sum_{i \in M_k} \nabla_w L(f(x_i, w^{(k)}), y_i)$$

# Stochastic Gradient Descent

Gradient descent minimizes a train error - an approximation to the generalization error

$$E[f] = \int L(f(x), y) d P(x, y) \Rightarrow E_n[f] = \frac{1}{N} \sum_{i=1}^N L(f(x_i, w), y_i)$$

Gradient Descent:

$$w^{(k+1)} = w^{(k)} - \eta \frac{1}{N} \sum_{i=1}^N \nabla_w L(f(x_i, w^{(k)}), y_i)$$

Stochastic Gradient Descent (SGD) with randomly selected mini-batches  $M_k$  of some fixed size  $N_{MB} \ll N$ :

$$w^{(k+1)} = w^{(k)} - \eta \frac{1}{N} \sum_{i \in M_k} \nabla_w L(f(x_i, w^{(k)}), y_i)$$

As mini-batches are randomly drawn from a data-generating distribution, **SGD minimizes the generalization error plus noise:**

$$w^{(k+1)} = w^{(k)} - \eta \nabla_w E[f] + \varepsilon^{(k)}$$

# Stochastic Gradient Descent

Gradient descent minimizes a train error - an approximation to the generalization error

$$E[f] = \int L(f(x), y) d P(x, y) \Rightarrow E_n[f] = \frac{1}{N} \sum_{i=1}^N L(f(x_i, w), y_i)$$

Gradient Descent:

$$w^{(k+1)} = w^{(k)} - \eta \frac{1}{N} \sum_{i=1}^N \nabla_w L(f(x_i, w^{(k)}), y_i)$$

Stochastic Gradient Descent (SGD) with randomly selected mini-batches  $M_k$  of some fixed size  $N_{MB} \ll N$ :

$$w^{(k+1)} = w^{(k)} - \eta \frac{1}{N} \sum_{i \in M_k} \nabla_w L(f(x_i, w^{(k)}), y_i)$$

As mini-batches are randomly drawn from a data-generating distribution, SGD minimizes the generalization error plus noise:

$$w^{(k+1)} = w^{(k)} - \eta \nabla_w E[f] + \varepsilon^{(k)}$$

**On-line SGD** ( $N_{MB} = 1$ ): convergence is slow, but guaranteed under certain conditions on  $\eta$

# Stochastic Gradient Descent

Gradient descent minimizes a train error - an approximation to the generalization error

$$E[f] = \int L(f(x), y) d P(x, y) \Rightarrow E_n[f] = \frac{1}{N} \sum_{i=1}^N L(f(x_i, w), y_i)$$

Gradient Descent:

$$w^{(k+1)} = w^{(k)} - \eta \frac{1}{N} \sum_{i=1}^N \nabla_w L(f(x_i, w^{(k)}), y_i)$$

Stochastic Gradient Descent (SGD) with randomly selected mini-batches  $M_k$  of some fixed size  $N_{MB} \ll N$ :

$$w^{(k+1)} = w^{(k)} - \eta \frac{1}{N} \sum_{i \in M_k} \nabla_w L(f(x_i, w^{(k)}), y_i)$$

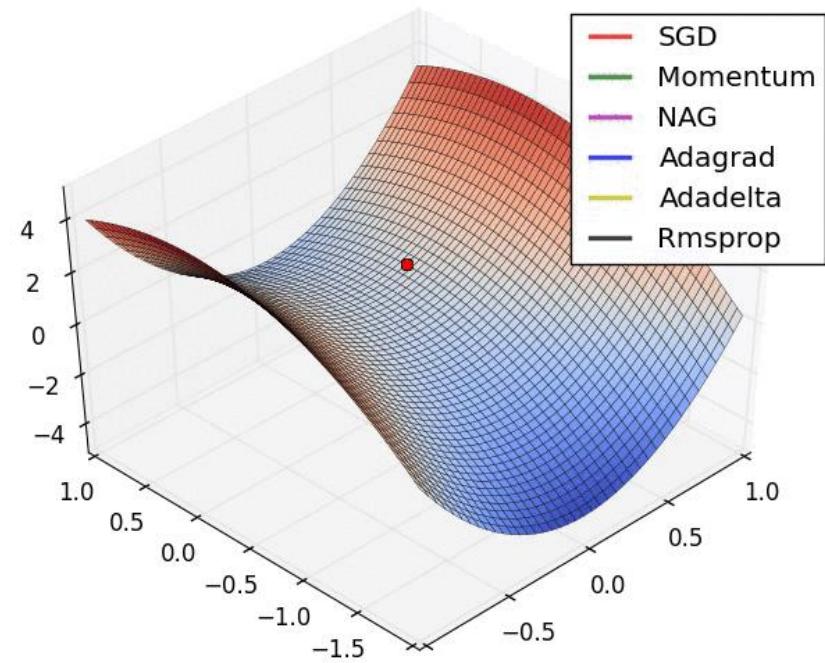
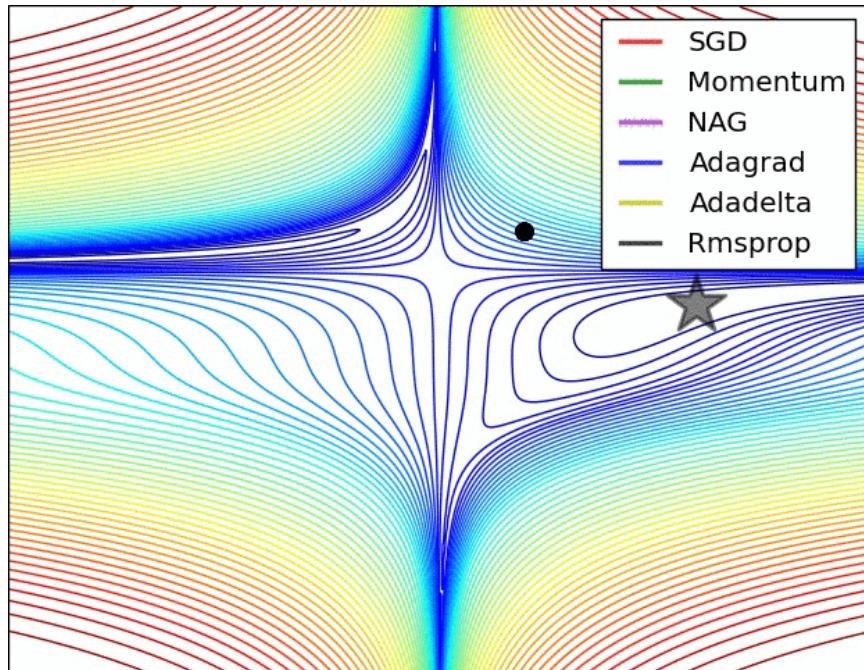
As mini-batches are randomly drawn from a data-generating distribution, SGD minimizes the generalization error plus noise:

$$w^{(k+1)} = w^{(k)} - \eta \nabla_w E[f] + \varepsilon^{(k)}$$

On-line SGD ( $N_{MB} = 1$ ): convergence is slow, but guaranteed under certain conditions on  $\eta$

**SGD is one of the most important algorithms** for neural networks, and for ML in general!

# Advanced Variants of SGD



# The Bayesian Recipe

Sum

$$p(A + B) = p(A) + p(B) - p(AB)$$

Product

$$p(AB) = p(A|B)p(B) = p(B|A)p(A)$$

Bayes' theorem

$$p(A|B) = \frac{p(B|A) p(A)}{p(B)}$$

Law of Total Probability

$$p(A|\mathcal{I}) = \sum_i p(A|B_i, \mathcal{I})p(B_i|\mathcal{I})$$

for  $\{B_i\}$   
exclusive  
and  
exhaustive

Normalisation

$$\sum_i p(B_i|\mathcal{I}) = 1$$

for  $\{B_i\}$   
exclusive  
and  
exhaustive

# ML Training for Probabilistic Models

# ML Training for Probabilistic Models

ML essentially amounts to a **model estimation** problem: What are the most probable values of parameters  $\Theta$ , given the model  $\mathbf{M} = f(\mathbf{X}, \Theta)$  and data  $\mathbf{D}$ ? That is, we need to find  $\Theta$  that maximizes  $p(\Theta | \mathbf{D}, \mathbf{M})$ . Use **Bayes' rule**: (assuming our model is a **probabilistic model!**)

$$\max_{\Theta} p(\Theta | \mathbf{D}, \mathbf{M}) = \max_{\Theta} \frac{p(\mathbf{D} | \mathbf{M}, \Theta) p(\Theta | \mathbf{M})}{p(\mathbf{D} | \mathbf{M})}$$

# ML Training for Probabilistic Models

ML essentially amounts to a **model estimation** problem: What are the most probable values of parameters  $\Theta$ , given the model  $\mathbf{M} = f(\mathbf{X}, \Theta)$  and data  $\mathbf{D}$ ? That is, we need to find  $\Theta$  that maximizes  $p(\Theta | \mathbf{D}, \mathbf{M})$ . Use **Bayes' rule**: (assuming our model is a **probabilistic model!**)

$$\begin{aligned}\max_{\Theta} p(\Theta | \mathbf{D}, \mathbf{M}) &= \max_{\Theta} \frac{p(\mathbf{D} | \mathbf{M}, \Theta) p(\Theta | \mathbf{M})}{p(\mathbf{D} | \mathbf{M})} \\ &= \max_{\Theta} \frac{p(\mathbf{D} | \mathbf{M}, \Theta) p(\Theta | \mathbf{M})}{\int p(\mathbf{D} | \mathbf{M}, \Theta) p(\Theta | \mathbf{M}) d\Theta}\end{aligned}$$

# ML Training for Probabilistic Models

ML essentially amounts to a **model estimation** problem: What are the most probable values of parameters  $\Theta$ , given the model  $\mathbf{M} = f(\mathbf{X}, \Theta)$  and data  $\mathbf{D}$ ? That is, we need to find  $\Theta$  that maximizes  $p(\Theta | \mathbf{D}, \mathbf{M})$ . Use **Bayes' rule**: (assuming our model is a **probabilistic model!**)

$$\begin{aligned}\max_{\Theta} p(\Theta | \mathbf{D}, \mathbf{M}) &= \max_{\Theta} \frac{p(\mathbf{D} | \mathbf{M}, \Theta) p(\Theta | \mathbf{M})}{p(\mathbf{D} | \mathbf{M})} \\ &= \max_{\Theta} \frac{p(\mathbf{D} | \mathbf{M}, \Theta) p(\Theta | \mathbf{M})}{\int p(\mathbf{D} | \mathbf{M}, \Theta) p(\Theta | \mathbf{M}) d\Theta} \\ &= \max_{\Theta} \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}\end{aligned}$$

# ML Training for Probabilistic Models

ML essentially amounts to a **model estimation** problem: What are the most probable values of parameters  $\Theta$ , given the model  $\mathbf{M} = f(\mathbf{X}, \Theta)$  and data  $\mathbf{D}$ ? That is, we need to find  $\Theta$  that maximizes  $p(\Theta | \mathbf{D}, \mathbf{M})$ . Use **Bayes' rule**: (assuming our model is a **probabilistic model!**)

$$\begin{aligned}\max_{\Theta} p(\Theta | \mathbf{D}, \mathbf{M}) &= \max_{\Theta} \frac{p(\mathbf{D} | \mathbf{M}, \Theta) p(\Theta | \mathbf{M})}{p(\mathbf{D} | \mathbf{M})} \\ &= \max_{\Theta} \frac{p(\mathbf{D} | \mathbf{M}, \Theta) p(\Theta | \mathbf{M})}{\int p(\mathbf{D} | \mathbf{M}, \Theta) p(\Theta | \mathbf{M}) d\Theta} \\ &= \max_{\Theta} \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}\end{aligned}$$

Measures the match between the data and prediction of the model

Measures how well the model can describe the data

Prior beliefs about which values of the model are most reasonable

# Maximum Likelihood Estimation MLE

Bayes' rule produced a general relation:

$$\max_{\Theta} p(\Theta | D, M) = \max_{\Theta} \frac{p(D | M, \Theta) p(\Theta | M)}{p(D | M)} = \max_{\Theta} \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$

**Special cases:**

- The **model is fixed**, can drop conditioning on  $M$ :  $p(D | M) \rightarrow p(D)$
- The model is fixed, and a flat prior  $p(\Theta | M) = \text{const}$  is used:  
**Maximum Likelihood Estimation (MLE)**

$$\max_{\Theta} p(\Theta | D, M) = \max_{\Theta} p(D | M, \Theta)$$

# Why MLE

## Maximum Likelihood Estimation (MLE):

$$\max_{\Theta} p(\Theta | \mathbf{D}, \mathbf{M}) = \max_{\Theta} p(\mathbf{D} | \mathbf{M}, \Theta)$$

### Properties:

- **Consistent** estimator (i.e.  $\hat{\Theta}_{N \rightarrow \infty} \rightarrow \Theta$ ) if:
  - The true distribution  $p_{\text{data}}(\mathbf{x})$  lies within the family  $\mathcal{F}_{\Theta}$
  - There is a unique value of  $\Theta$  that corresponds to  $p_{\text{data}}(\mathbf{x})$
- Asymptotically (as  $N \rightarrow \infty$ ), has the **lowest possible MSE** among all consistent estimators (Cramer-Rao lower bound, 1945-1946)
- For **finite**  $N$ , biased estimators (e.g. a regularized MLE) can be more efficient (i.e. reach the same level of generalization error with a smaller number of samples  $N$ )

# Maximum A-Posteriori Estimation

Bayes' rule produced a general relation:

$$\max_{\Theta} p(\Theta | D, M) = \max_{\Theta} \frac{p(D | M, \Theta) p(\Theta | M)}{p(D | M)} = \max_{\Theta} \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$

**Special cases:**

- The **model is fixed**, can drop conditioning on  $M$ :  $p(D | M) \rightarrow p(D)$
- The model is fixed, and a flat prior  $p(\Theta | M) = \text{const}$  is used:  
Maximum Likelihood Estimation (MLE)

$$\max_{\Theta} p(\Theta | D, M) = \max_{\Theta} p(D | M, \Theta)$$

- The model is fixed, and a non-flat prior  $p(\Theta | M)$  is used:  
**Maximum A-Posteriori (MAP) estimation**

$$\max_{\Theta} p(\Theta | D, M) = \max_{\Theta} p(D | M, \Theta) p(\Theta | M)$$

# MLE and Least Squares Loss

# MLE and Least Squares Loss

Assume we want to fit  $N$  noisy measurements of quantity  $y_n$  as a function of variable  $\mathbf{x}_n$  :  $y_n = f(\mathbf{x}_n, \Theta) + \varepsilon_n$

Assuming that errors  $\varepsilon_n$  are i.i.d. and have a Gaussian distribution with variance  $\sigma_n^2$ , the likelihood is

$$p(\mathbf{D} | \mathbf{M}, \Theta) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma_n^2}} e^{-\frac{(y_n - f(\mathbf{x}_n, \Theta))^2}{2\sigma_n^2}}$$

# MLE and Least Squares Loss

Assume we want to fit  $N$  noisy measurements of quantity  $y_n$  as a function of variable  $\mathbf{x}_n$  :  $y_n = f(\mathbf{x}_n, \Theta) + \varepsilon_n$

Assuming that errors  $\varepsilon_n$  are i.i.d. and have a Gaussian distribution with variance  $\sigma_n^2$ , the likelihood is

$$p(\mathbf{D} | \mathbf{M}, \Theta) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma_n^2}} e^{-\frac{(y_n - f(\mathbf{x}_n, \Theta))^2}{2\sigma_n^2}}$$

Maximization of  $p(\mathbf{D} | \mathbf{M}, \Theta)$  is equivalent to minimization of **negative log-likelihood (NLL)**:

$$\min_{\Theta} (-\log p(\mathbf{D} | \mathbf{M}, \Theta)) = \min_{\Theta} \sum_{n=1}^N \frac{(y_n - f(\mathbf{x}_n, \Theta))^2}{2\sigma_n^2} + \frac{1}{2} \log(2\pi\sigma_n^2)$$

# MLE and Least Squares Loss

Assume we want to fit  $N$  noisy measurements of quantity  $y_n$  as a function of variable  $\mathbf{x}_n$  :  $y_n = f(\mathbf{x}_n, \Theta) + \varepsilon_n$

Assuming that errors  $\varepsilon_n$  are i.i.d. and have a Gaussian distribution with variance  $\sigma_n^2$ , the likelihood is

$$p(\mathbf{D} | \mathbf{M}, \Theta) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma_n^2}} e^{-\frac{(y_n - f(\mathbf{x}_n, \Theta))^2}{2\sigma_n^2}}$$

Maximization of  $p(\mathbf{D} | \mathbf{M}, \Theta)$  is equivalent to minimization of **negative log-likelihood (NLL)**:

$$\min_{\Theta} (-\log p(\mathbf{D} | \mathbf{M}, \Theta)) = \min_{\Theta} \sum_{n=1}^N \frac{(y_n - f(\mathbf{x}_n, \Theta))^2}{2\sigma_n^2} + \frac{1}{2} \log(2\pi\sigma_n^2)$$

When variances  $\sigma_n^2$  are constant, this is equivalent to a **minimum least squares error (MSE)** function

$$\min_{\Theta} \sum_{n=1}^N (y_n - f(\mathbf{x}_n, \Theta))^2$$

Becomes linear regression if  
 $f(\mathbf{x}_n, \Theta) = \sum \Theta_k x_{nk} = \Theta^T \mathbf{x}$

# Kullback-Leibler (KL) Divergence

Another interpretation of MLE : minimization of the KL divergence between the model distribution  $p_{\text{model}}(\mathbf{x}, \Theta)$  and the true distribution  $p_{\text{data}}(\mathbf{x})$ :

$$\begin{aligned} D_{KL}(p_{\text{data}} \parallel p_{\text{model}}) &= \mathbb{E}_{x \sim p_{\text{data}}} \left[ \log \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{model}}(\mathbf{x})} \right] \\ &= \mathbb{E}_{x \sim p_{\text{data}}} [\log p_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x})] \\ &= \underbrace{\mathbb{E}_{x \sim p_{\text{data}}} [\log p_{\text{data}}(\mathbf{x})]}_{\text{independent of } p_{\text{model}}} - \mathbb{E}_{x \sim p_{\text{data}}} [\log p_{\text{model}}(\mathbf{x})] \\ &\simeq -\mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x})] + \dots = -\sum \log p_{\text{model}}(\mathbf{x}_n, \Theta) \end{aligned}$$

KL divergence between two distributions  $p_1(\mathbf{x})$  and  $p_2(\mathbf{x})$  measures their dissimilarity:

$$D_{KL}(p_1 \parallel p_2) \geq 0$$

$$D_{KL}(p_1 \parallel p_2) = 0 \quad \text{iff } p_1(\mathbf{x}) = p_2(\mathbf{x})$$

KL divergence is widely used in ML!

# MAP and Regularization

# MAP and Regularization

Assume we want to fit  $N$  noisy measurements of quantity  $y_n$  as a function of variable  $\mathbf{X}_n$ :

$$\max_{\Theta} p(\Theta | \mathbf{D}, \mathbf{M}) = \max_{\Theta} p(\mathbf{D} | \mathbf{M}, \Theta) p(\Theta | \mathbf{M})$$

Negative log-likelihood for our problem  $y_n = f(\mathbf{X}_n, \Theta) + \varepsilon_n$ :

$$\min_{\Theta} (-\log p(\mathbf{D} | \mathbf{M}, \Theta)) = \min_{\Theta} \sum_{n=1}^N \frac{(y_n - f(\mathbf{x}_n, \Theta))^2}{2\sigma_n^2} - \log p(\Theta | \mathbf{M}) + \dots$$



# MAP and Regularization

Recall the definition of the Maximum A-Posteriori (MAP) Estimation of model parameters:

$$\max_{\Theta} p(\Theta | \mathbf{D}, \mathbf{M}) = \max_{\Theta} p(\mathbf{D} | \mathbf{M}, \Theta) p(\Theta | \mathbf{M})$$

Negative log-likelihood for our problem  $y_n = f(\mathbf{x}_n, \Theta) + \varepsilon_n$  :

$$\min_{\Theta} (-\log p(\mathbf{D} | \mathbf{M}, \Theta)) = \min_{\Theta} \sum_{n=1}^N \frac{(y_n - f(\mathbf{x}_n, \Theta))^2}{2\sigma_n^2} - \log p(\Theta | \mathbf{M}) + \dots$$



Examples:

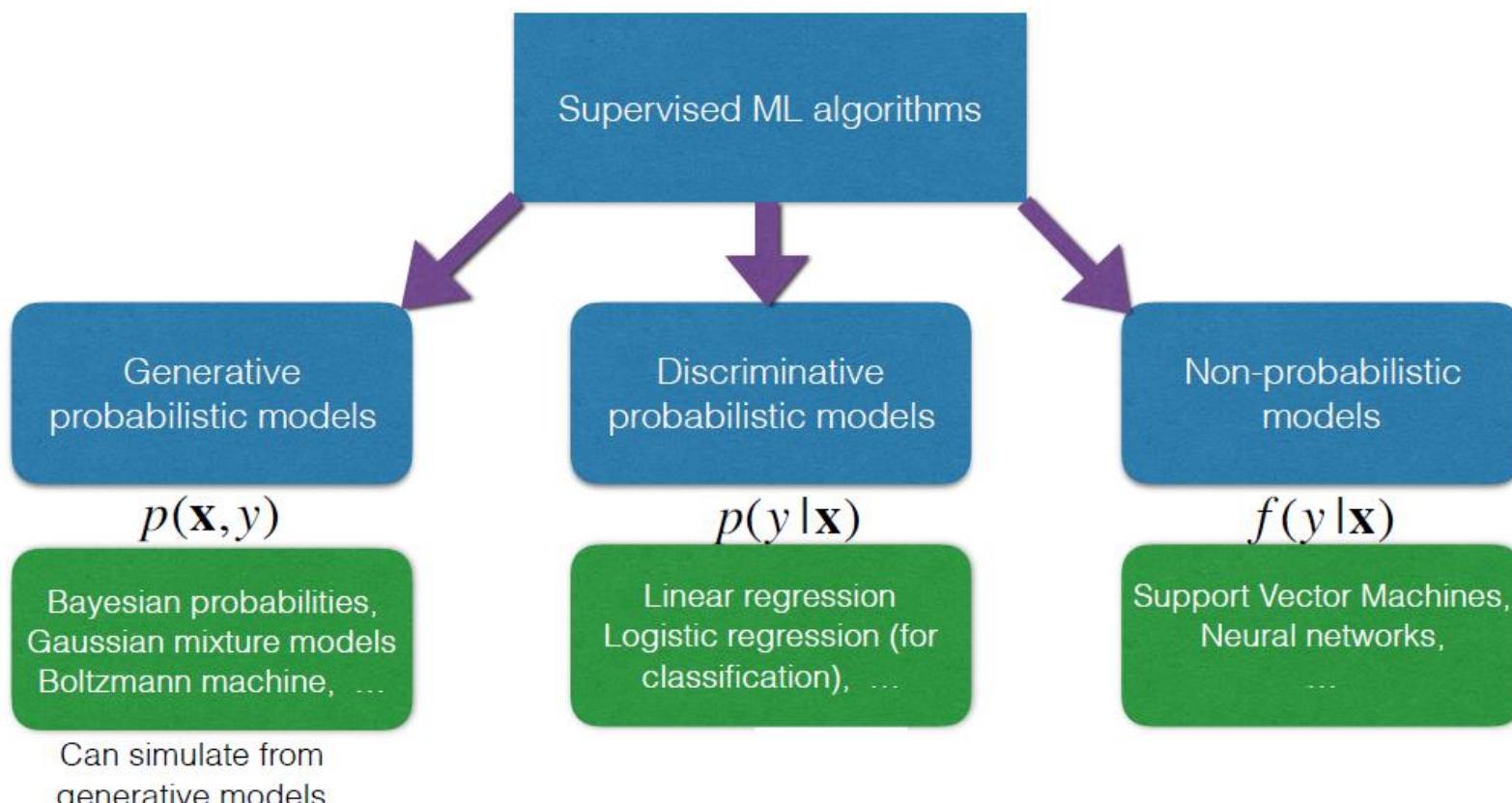
- Gaussian prior  $\sim e^{-\lambda(\theta-\theta_0)^2} \Rightarrow L_2$  regularization  $\lambda \|\Theta - \Theta_0\|_2$
- Laplace prior  $\sim e^{-\alpha|\theta-\theta_0|} \Rightarrow L_1$  regularization  $\lambda \|\Theta - \Theta_0\|_1$

# Supervised Learning Algorithms

Most, but not all, supervised ML algorithms amount to **estimating a probability distribution**  $p(y|x)$

Example: Linear Regression is equivalent to a **discriminative probabilistic model**

$$p(y|x) = \mathcal{N}(y; \Theta^T x; \sigma^2)$$



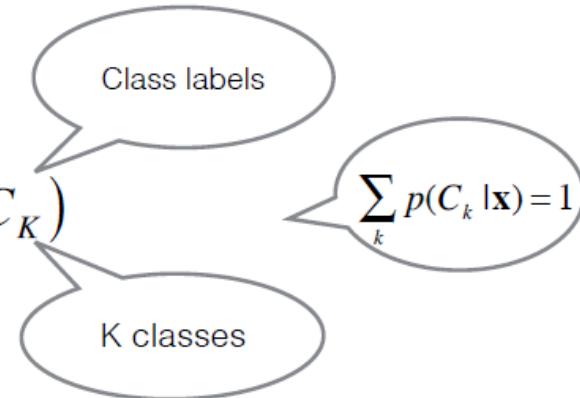
# Probabilistic Classification Models

# Probabilistic Classification Models

For probabilistic **regression** models, we estimate a probability distribution  $p(y | \mathbf{x})$ ,  $y \in \mathbb{R}$

Probabilistic **classification** models:

$$p(y | \mathbf{x}), \quad y \in (C_1, \dots, C_K)$$

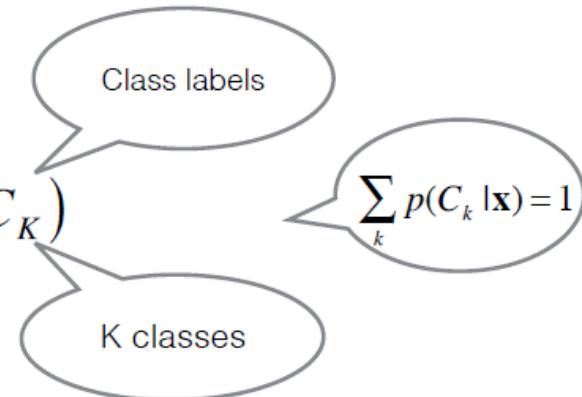


# Probabilistic Classification Models

For probabilistic **regression** models, we estimate a probability distribution  $p(y | \mathbf{x})$ ,  $y \in \mathbb{R}$

Probabilistic **classification** models:

$$p(y | \mathbf{x}), \quad y \in (C_1, \dots, C_K)$$



**Binary classification** ( $K=2$ ):

Here  $y = \{0,1\}$ , or  $y = \{-1,1\}$ , or  $y = \{\text{low risk}, \text{high risk}\}$ , e.t.c.  
Only one unknown: the “positive class” probability

$$p(y = 1 | \mathbf{x})$$

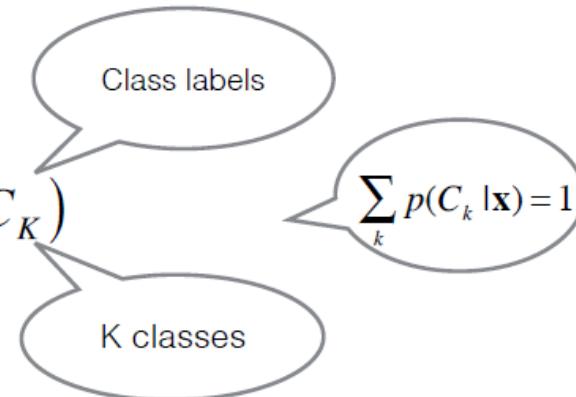
with  $p(y = 0 | \mathbf{x}) = 1 - p(y = 1 | \mathbf{x})$

# Probabilistic Classification Models

For probabilistic **regression** models, we estimate a probability distribution  $p(y | \mathbf{x})$ ,  $y \in \mathbb{R}$

Probabilistic **classification** models:

$$p(y | \mathbf{x}), \quad y \in (C_1, \dots, C_K)$$



**Binary classification** ( $K=2$ ):

Here  $y = \{0,1\}$ , or  $y = \{-1,1\}$ , or  $y = \{\text{low risk, high risk}\}$ , e.t.c.  
Only one unknown: the “positive class” probability

$$p(y = 1 | \mathbf{x})$$

with  $p(y = 0 | \mathbf{x}) = 1 - p(y = 1 | \mathbf{x})$

Examples of binary classification: bank failures, mortgage defaults, credit card fraud, anti-money laundering, price direction prediction

# Logistic Regression

# Logistic Regression

Recall that a Linear Regression model fits a linear function

$$y = \Theta^T \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^N, \boxed{y \in \mathbb{R}}$$

This is *not* on its own suitable to describe probabilities (as they should be numbers between 0 and 1).

# Logistic Regression

Recall that a Linear Regression model fits a linear function

$$y = \Theta^T \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^N, \boxed{y \in \mathbb{R}}$$

This is *not* on its own suitable to describe probabilities (as they should be numbers between 0 and 1). **But!**

Let us apply a “squashing” function to  $y$ :

$$y \rightarrow h(y), \quad h: \mathbb{R} \rightarrow [0,1]$$

Nonlinear function

# Logistic Regression

Recall that a Linear Regression model fits a linear function

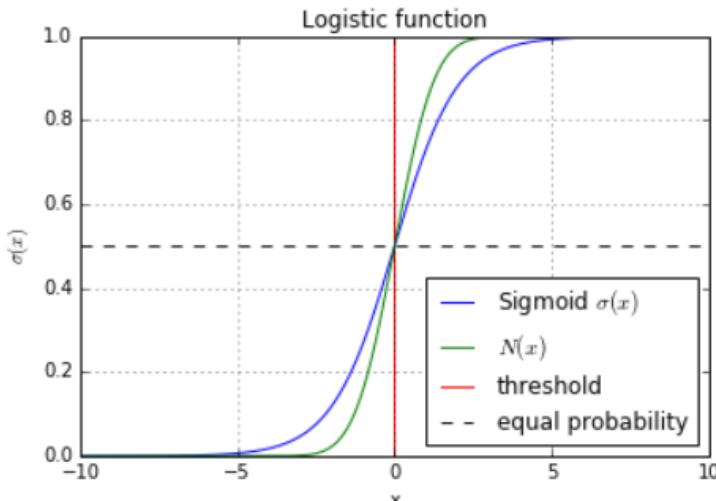
$$y = \Theta^T \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^N, \quad y \in \mathbb{R}$$

This is *not* on its own suitable to describe probabilities (as they should be numbers between 0 and 1). But!

Let us apply a “squashing” function to  $y$ :

$$y \rightarrow h(y), \quad h: \mathbb{R} \rightarrow [0,1]$$

**Logistic Regression:** take  $h(y) = \sigma(y) = \frac{1}{1 + \exp(-y)} \Rightarrow$



$$p(y=1 | \mathbf{x}) = \sigma(\theta^T \mathbf{x}) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})}$$

Can also be derived from other models

# MLE for Logistic Regression

**Logistic Regression:**

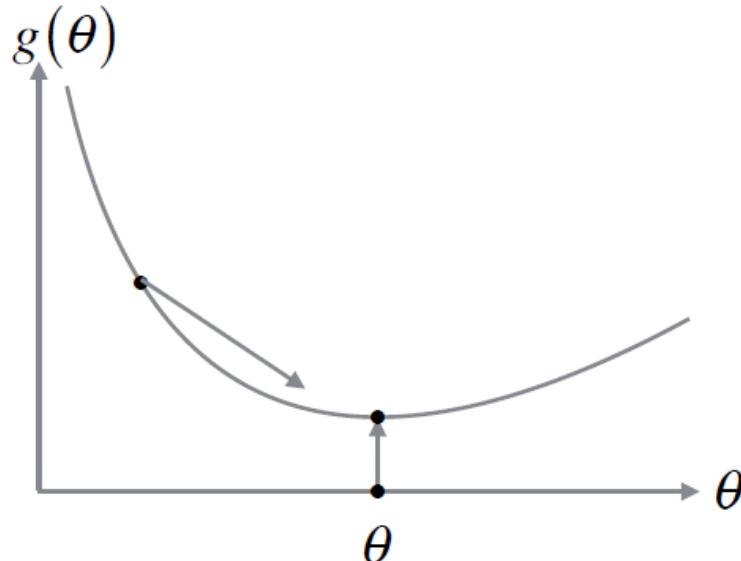
$$p_n(\theta) = p(y_n = 1 | \mathbf{x}_n) = \sigma(\theta^T \mathbf{x}_n) = \frac{1}{1 + \exp(-\theta^T \mathbf{x}_n)}$$

Likelihood:

$$p(\mathbf{D} | \mathbf{M}, \Theta) = \prod_{n=1}^N p_n(\theta)^{y_n} (1 - p_n(\theta))^{1-y_n}$$

observed  
values {0,1}

Negative LL:  $-\log p(\mathbf{D} | \mathbf{M}, \Theta) = \sum_{n=1}^N [y_n \log p_n(\theta) + (1 - y_n) \log(1 - p_n(\theta))]$



Can be  
minimized by gradient  
descent

Convex function of  $\theta$

Unique solution

# Regression and Classification Demo

