

# Introduction to Computer Vision

*11. Learning from data*

UCLA – CS 188 – Fall 2019

Fabien Scalzo, Ph.D.

<b>Week 1</b>			26-Sep	<b>Introduction</b>
<b>Week 2</b>	1-Oct	<b>Basic Image Processing</b>	3-Oct	<b>Feature Extraction and Classification</b>
<b>Week 3</b>	8-Oct	<b>Feature Tracking/Optical Flow</b>	10-Oct	<b>SVD, 2D camera model, projective plane</b>
<b>Week 4</b>	15-Oct	<b>2D Image transformations, RANSAC</b>	17-Oct	<b>Euclidean geometry, rigid body motion</b>
<b>Week 5</b>	22-Oct	<b>Epipolar Geometry</b>	24-Oct	<b>3D Cameras and processing</b>
<b>Week 6</b>	29-Oct	<b>Midterm</b>	31-Oct	<b>3D Cameras and processing</b>
<b>Week 7</b>	5-Nov	<b>Learning from data</b>	7-Nov	<b>Neural Networks</b>
<b>Week 8</b>	12-Nov	<b>Deep Learning</b>	14-Nov	<b>Deep Learning</b>
<b>Week 9</b>	19-Nov	<b>Object Detection</b>	21-Nov	<b>Generative Models</b>
<b>Week 10</b>	26-Nov	<b>Guest Lecture (Nikhil Naik)</b>	28-Nov	
<b>Week 11</b>	3-Dec	<b>Applications</b>	5-Dec	<b>Recap</b>
<b>Week 12</b>	10-Dec		12-Dec	<b>Final</b>

- 
- The learning problem
  - Linear Regression & Classification
  - Loss, Risk, Empirical Risk
  - Gradient descent
    - Numerical, Analytical
    - General Algorithm
    - Types of gradient descent
  - Under & Overfitting
  - Regularization
  - Cross-validation
-

---

# The learning problem

## Goal

To find the parameters of the model  $f$  such that it can map any point from input space  $X$  to the output space  $Y$ .

$$f: X \rightarrow Y$$

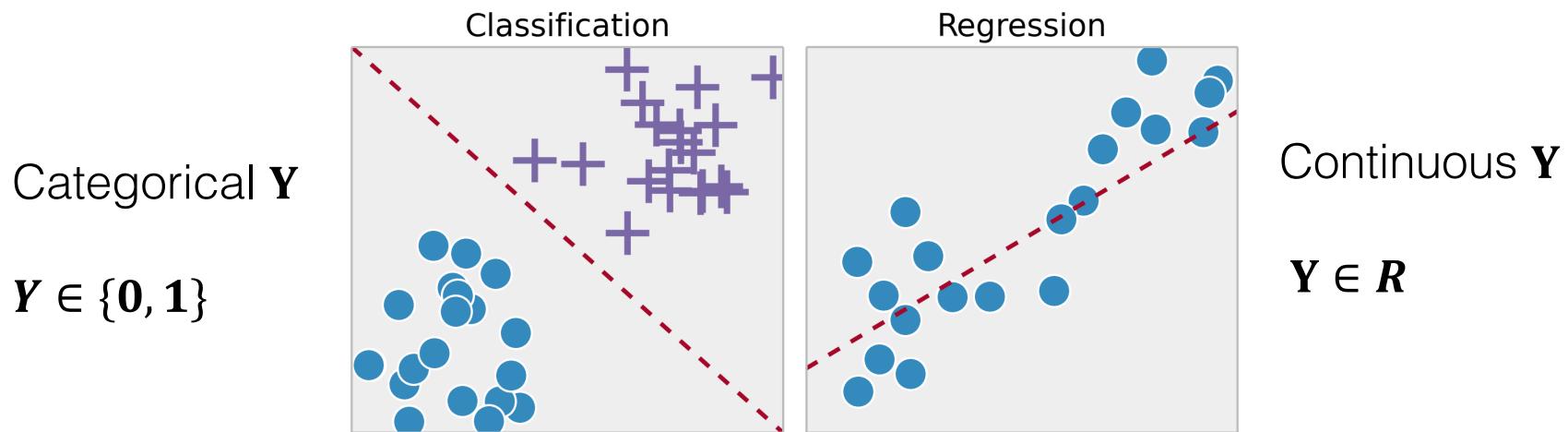
Sampled observations  $\{\mathbf{x}_i, \mathbf{y}_i\}_N$  from joint distribution  $P(x,y)$

- pairs of input variables  $x_i \in \mathbf{X}$  and output labels  $y_i \in \mathbf{Y}$
- $\{\mathbf{x}_i, \mathbf{y}_i\}_N$  denotes a set of  $N$  pairs of labeled data points
- sampled (i.i.d) from a process/distribution that is typically unknown

# The learning problem

Observed data

- pairs of input variables  $x_i \in \mathbf{X}$  and output labels  $y_i \in \mathbf{Y}$
- $\{x_i, y_i\}_N$  denotes a set of N pairs of labeled data points
- sampled (iid) from a process that is typically unknown  
-> our dataset



# Loss, Risk, and Empirical Risk

How well will the model work in practice on new data?

**Risk:**  $R(h) = \mathbf{E}[L(h(x), y)] = \int L(h(x), y) dP(x, y).$

In general, the risk  $R(h)$  cannot be computed because the distribution  $P(x, y)$  is unknown

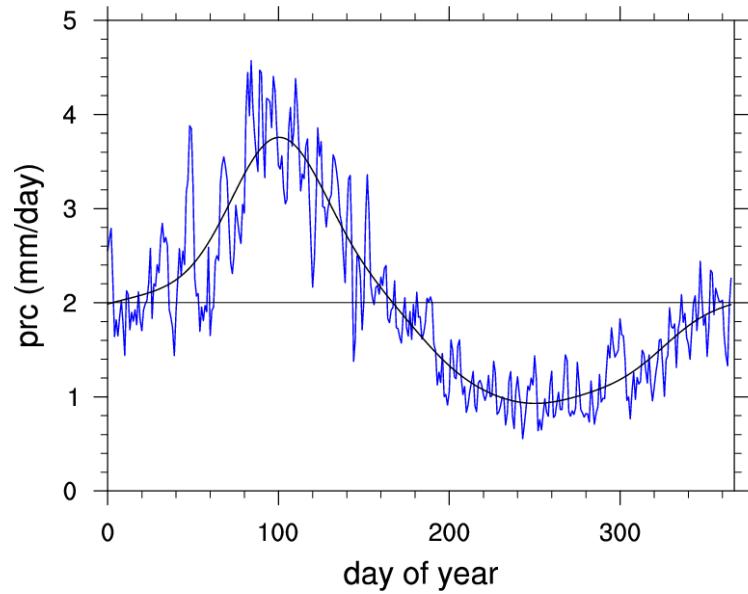
**LOSS:** measures how different the predicted output is from the groundtruth

Example:  $L(\hat{y}, y) = \begin{cases} 1 & \text{If } \hat{y} \neq y \\ 0 & \text{If } \hat{y} = y \end{cases}$ .

**Empirical Risk:** Approximation to true risk  $R(h)$ ; average loss over dataset

$$R_{\text{emp}}(h) = \frac{1}{n} \sum_{i=1}^n L(h(x_i), y_i).$$

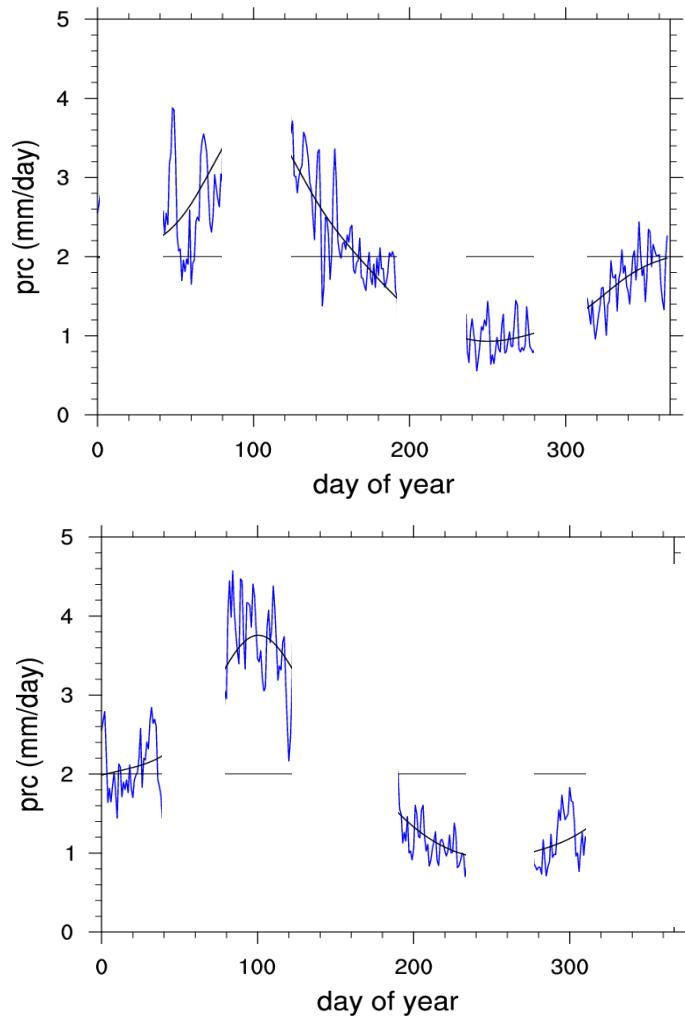
# Separate dataset into training and test set



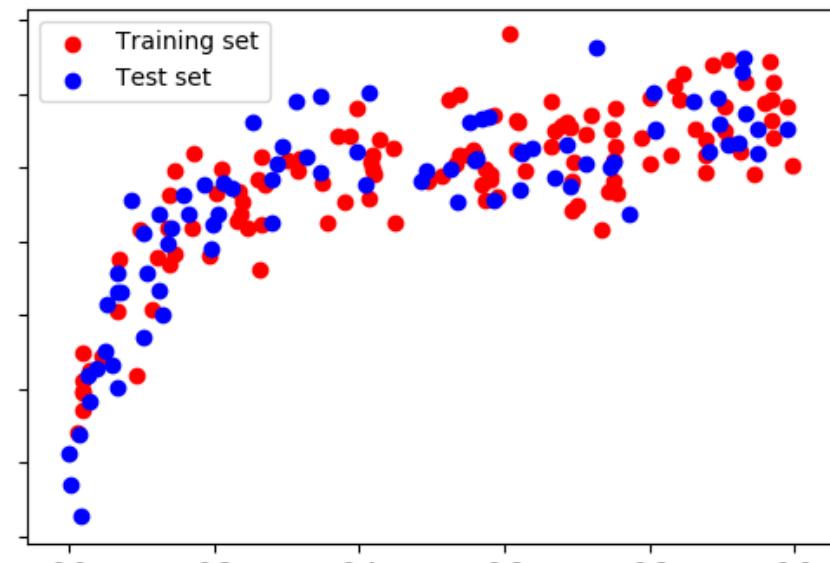
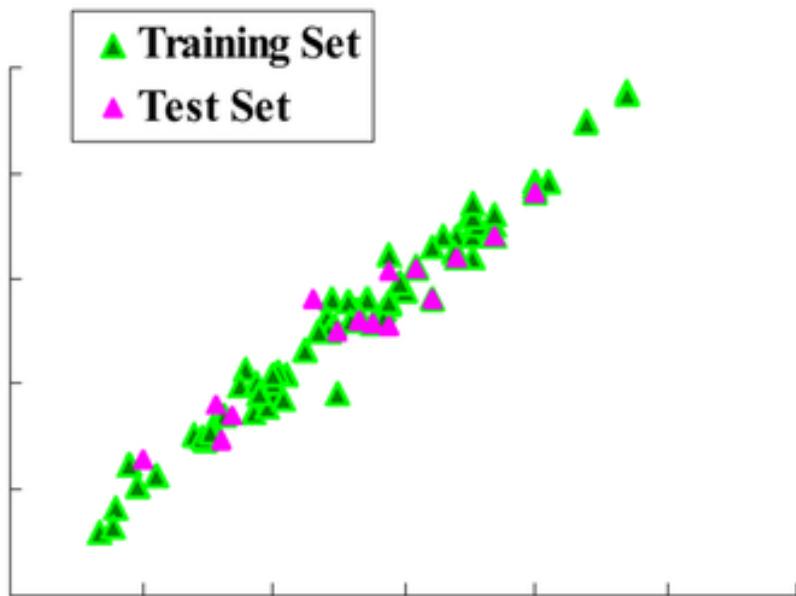
training

?

testing



# Separate dataset into training and test set

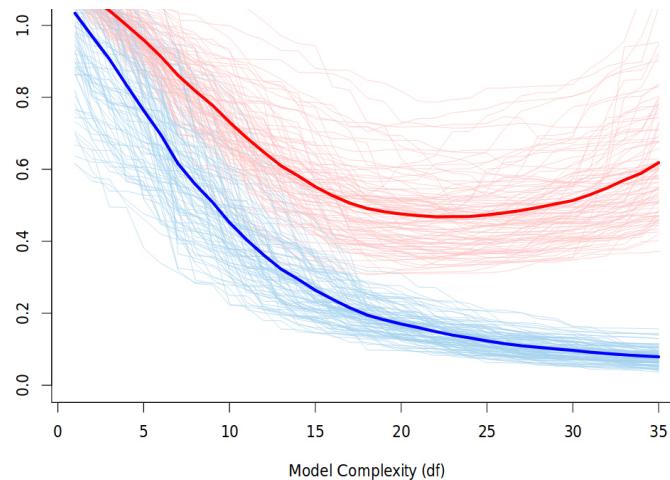
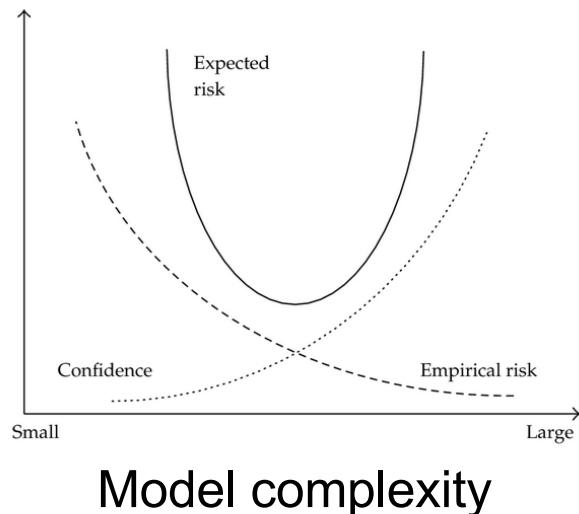


# Minimizing the Empirical Risk

$$R_{\text{emp}}(h) = \frac{1}{n} \sum_{i=1}^n L(h(x_i), y_i).$$

$$\hat{h} = \arg \min_{h \in \mathcal{H}} R_{\text{emp}}(h).$$

on **training set**  
and **test** set



# Loss functions for regression

1.Squared Loss  $(h(\mathbf{x}_i) - y_i)^2$

---

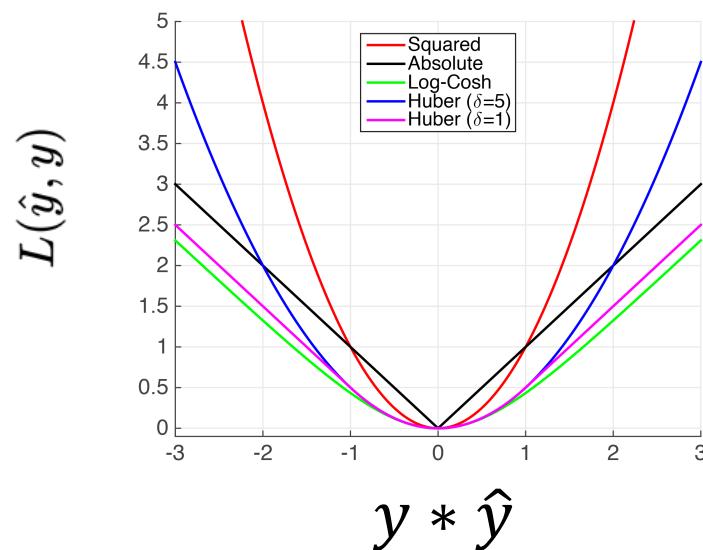
2.Absolute Loss  $|h(\mathbf{x}_i) - y_i|$

---

3.Huber Loss

- $\frac{1}{2}(h(\mathbf{x}_i) - y_i)^2$  if  $|h(\mathbf{x}_i) - y_i| < \delta$ ,
  - otherwise  $\delta(|h(\mathbf{x}_i) - y_i| - \frac{\delta}{2})$
- 

4.Log-Cosh Loss  
 $\log(\cosh(h(\mathbf{x}_i) - y_i)),$   
 $\cosh(x) = \frac{e^x + e^{-x}}{2}$



# Loss functions for classification

1.Hinge-Loss

$$\max[1 - h_{\mathbf{w}}(\mathbf{x}_i)y_i, 0]^p$$

2.Log-Loss

$$\log(1 + e^{-h_{\mathbf{w}}(\mathbf{x}_i)y_i})$$

3.Exponential Loss

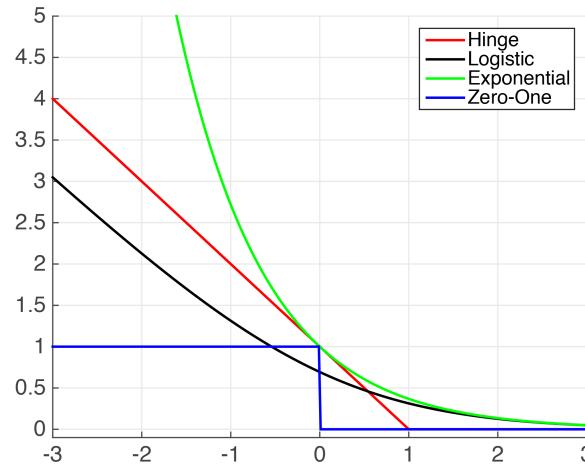
$$e^{-h_{\mathbf{w}}(\mathbf{x}_i)y_i}$$

4.Zero-One Loss

$$\delta(\text{sign}(h_{\mathbf{w}}(\mathbf{x}_i)) \neq y_i) = \begin{cases} 1 & \text{If } \hat{y} \neq y \\ 0 & \text{If } \hat{y} = y \end{cases}.$$

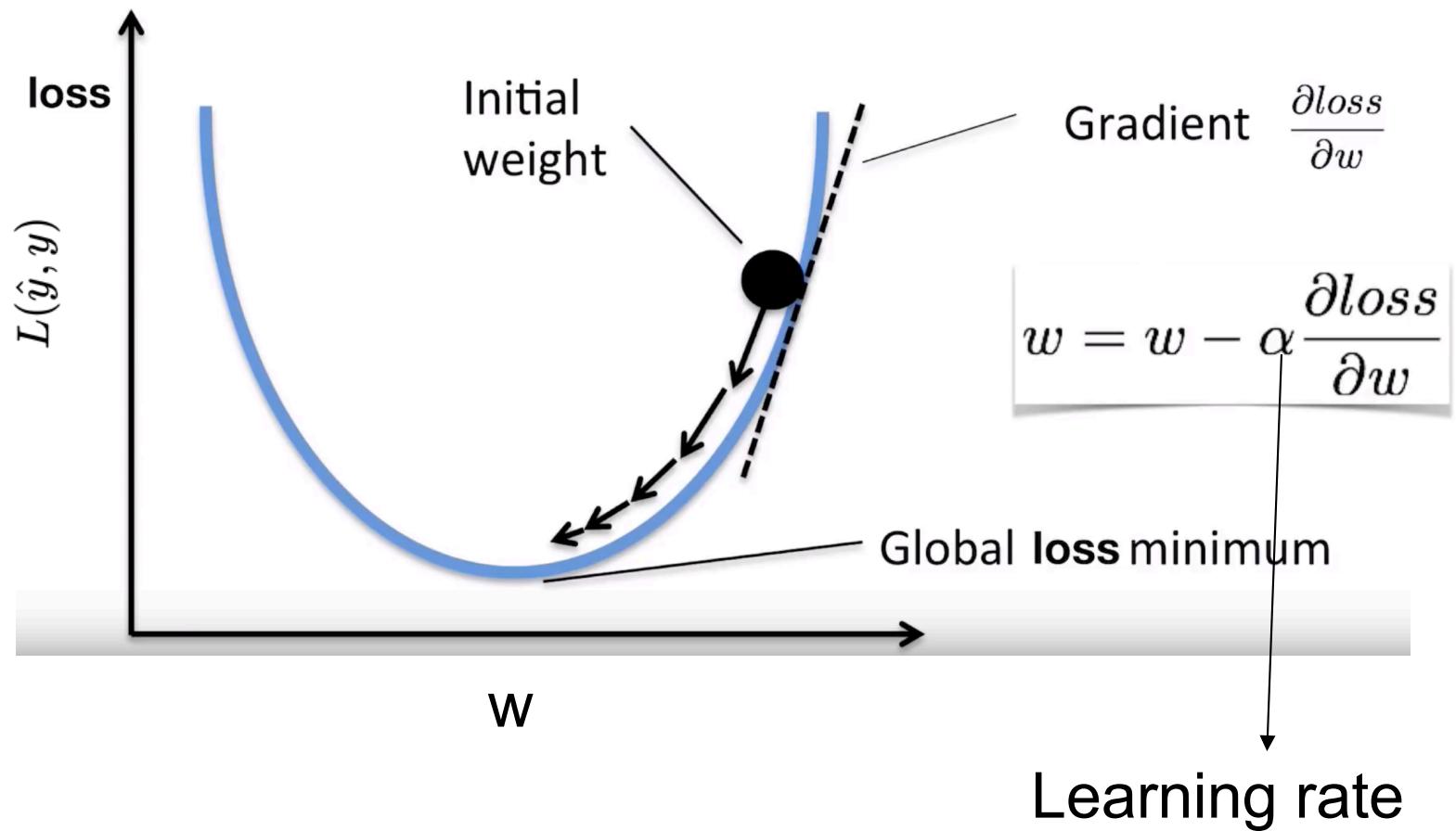
$$Y \in \{-1, +1\}$$

$$L(\hat{y}, y)$$



$$y * \hat{y}$$

# Gradient descent



$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

# Computing the Gradient

How do we compute the gradient  $\frac{\partial \text{loss}}{\partial w}$

- Numerical

Evaluate the difference of the Loss at  $L(w)$  and  $L(w + \epsilon)$

- Analytical

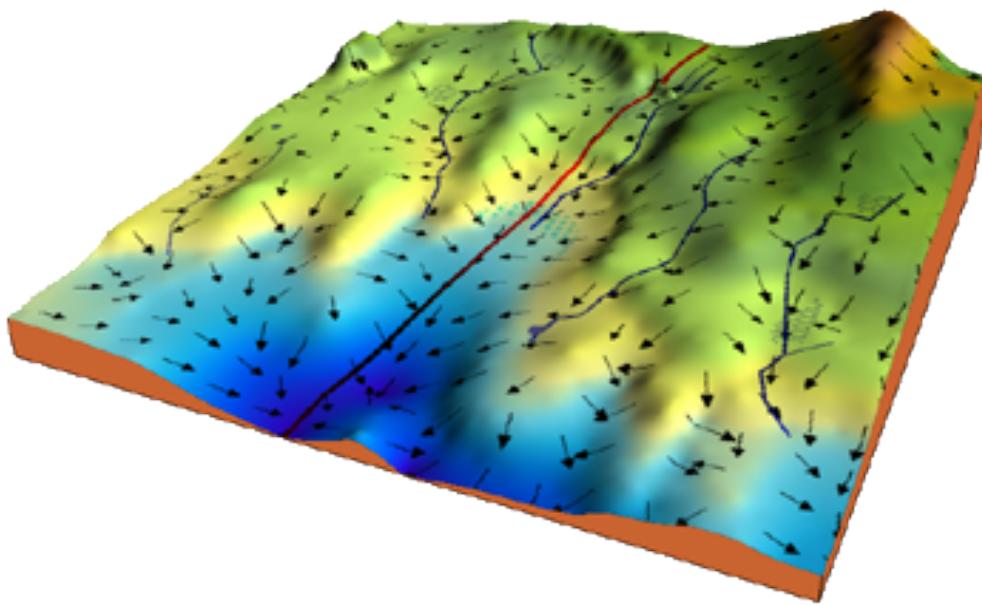
Given the cost function:

$$f(m, b) = \frac{1}{N} \sum_{i=1}^n (y_i - (mx_i + b))^2$$

The gradient can be calculated as:

$$f'(m, b) = \begin{bmatrix} \frac{df}{dm} \\ \frac{df}{db} \end{bmatrix} = \begin{bmatrix} \frac{1}{N} \sum -2x_i(y_i - (mx_i + b)) \\ \frac{1}{N} \sum -2(y_i - (mx_i + b)) \end{bmatrix}$$

# Gradient



# Gradient Descent Algorithms

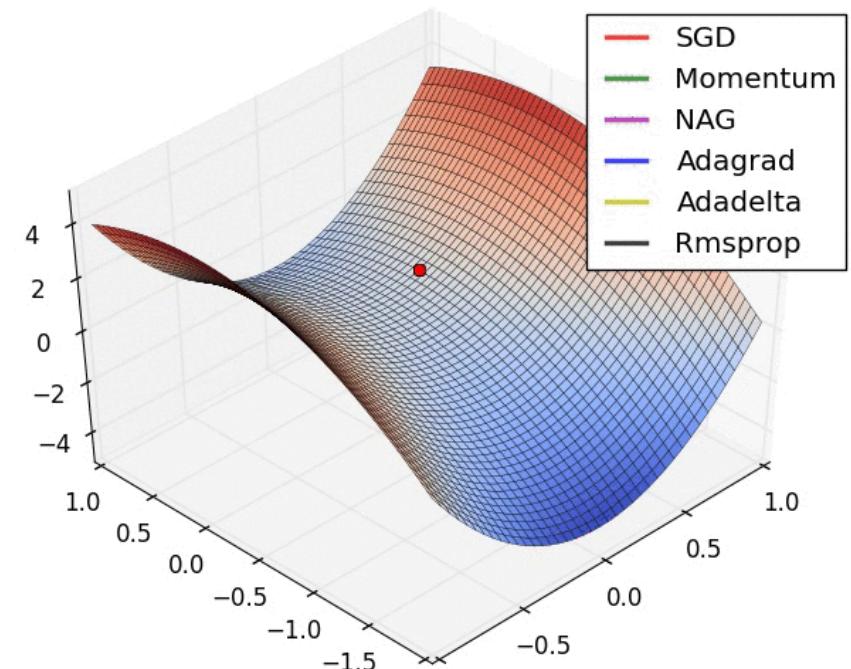
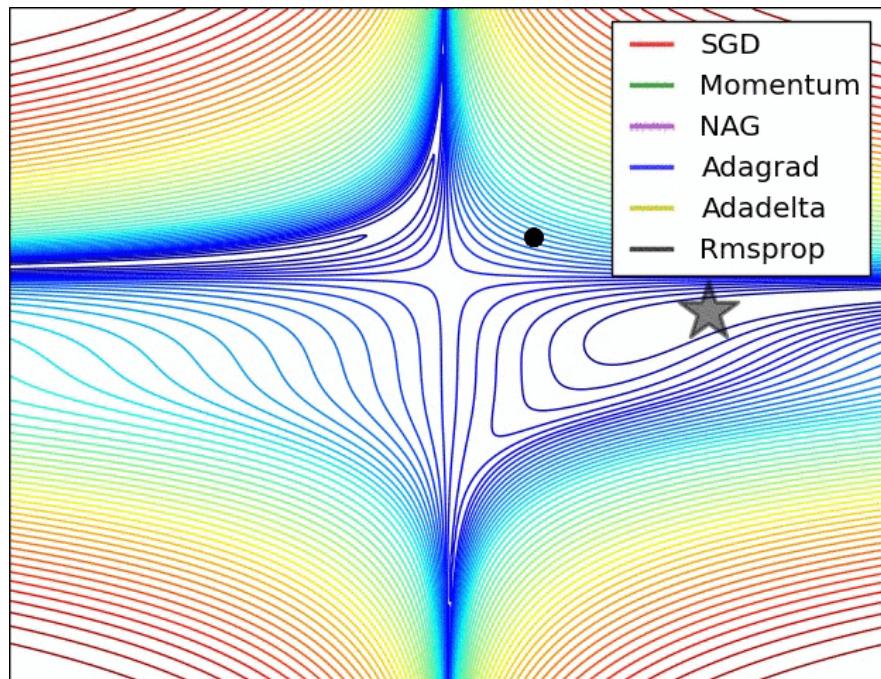
- Gradient descent variants
  - Batch gradient descent
    - Error is computed over all training samples at each iteration
  - Stochastic gradient descent
    - Error is computed over one training sample at each iteration
  - Mini-batch gradient descent
    - Error is computed over a subset of the training samples at each iteration

# Gradient Descent Algorithms

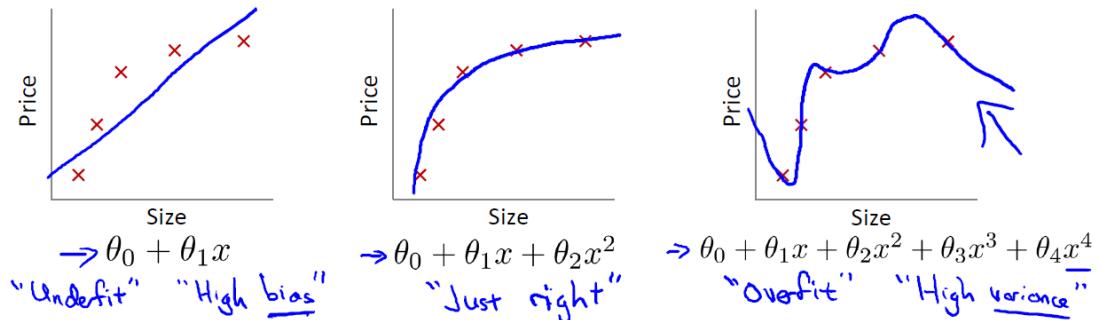
- Momentum
- Nesterov accelerated gradient
- Adagrad
- Adadelta
- RMSprop
- Adam
- AdaMax
- Nadam
- AMSGrad

Gradient Descent with Momentum  
considers the past gradients during  
the update

# Gradient Descent Algorithms



# Fitting Data



Symptoms	Underfitting	Just right	Overfitting
Regression	- High training error - Training error close to test error - High bias	- Training error slightly lower than test error	- Low training error - Training error much lower than test error - High variance
Classification			

# Regularization

The process of making the prediction function fit the training data less well in the hope that it generalizes new data better.  
Regularization penalizes large values for  $w$

$\ell_0$  complexity: number of non-zero coefficients

$\ell_1$  "lasso" complexity:  $\sum_{i=1}^d |w_i|$ , for coefficients  $w_1, \dots, w_d$

$\ell_2$  "ridge" complexity:  $\sum_{i=1}^d w_i^2$  for coefficients  $w_1, \dots, w_d$

# Regularized Empirical Risk Minimization

$$\operatorname{argmin}_{\mathbf{w}} R(\mathbf{w}) + C \sum_{i=1}^N L(y_i, \mathbf{w}^\top \mathbf{x}_i)$$

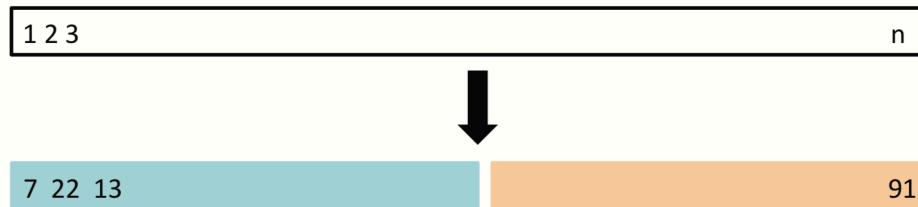
Regularization      Loss

# Cross-validation

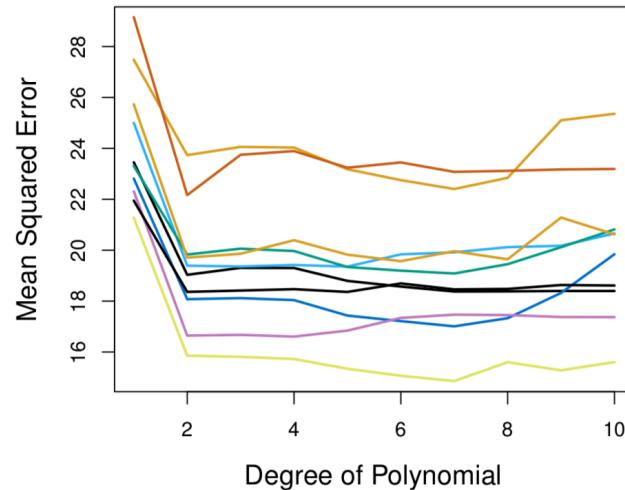
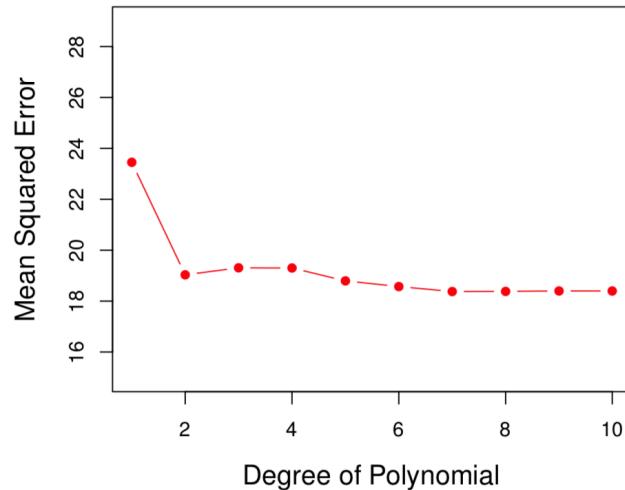
**Goal:** Estimate the test error for a supervised learning method.

**Strategy:**

- ▶ Split the data in two parts.
- ▶ Train the method in the first part.
- ▶ Compute the error on the second part.



# Cross-validation

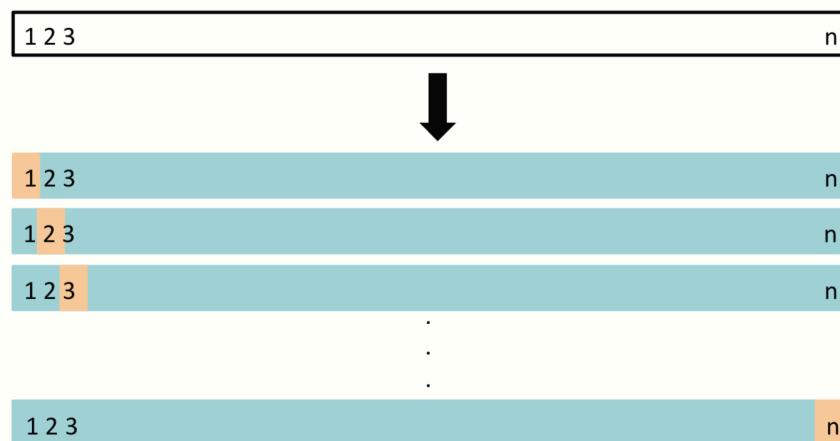


**Problems:** 1. Every split yields a different estimate of the error.  
2. Only a subset of points is used to evaluate the model.

# Cross-validation

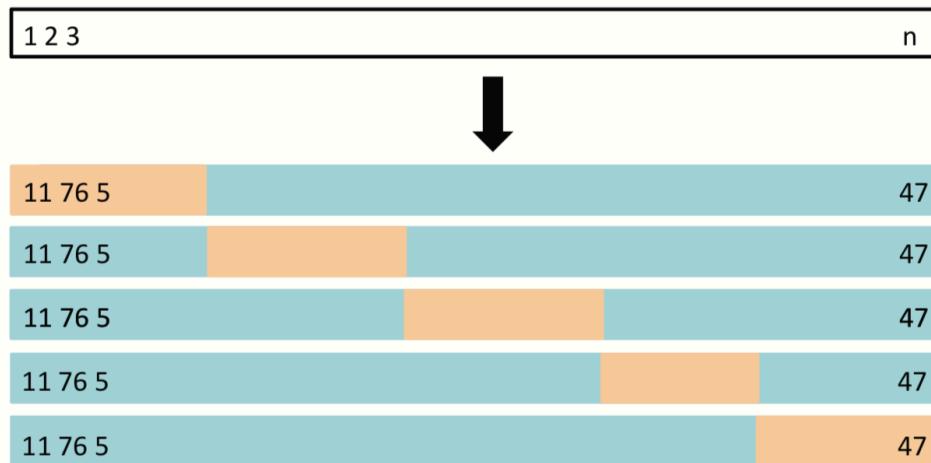
## Leave one out cross-validation

- ▶ For every  $i = 1, \dots, n$ :
  - ▶ train the model on every point except  $i$ ,
  - ▶ compute the test error on the held out point.
- ▶ Average the test errors.



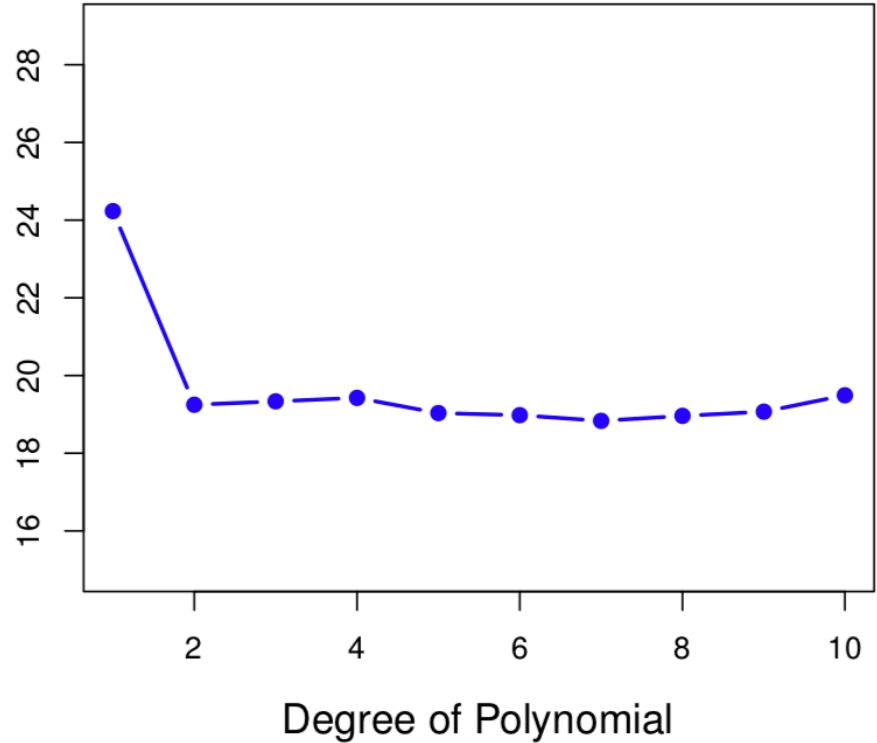
## *k*-fold cross-validation

- ▶ Split the data into  $k$  subsets or *folds*.
- ▶ For every  $i = 1, \dots, k$ :
  - ▶ train the model on every fold except the  $i$ th fold,
  - ▶ compute the test error on the  $i$ th fold.
- ▶ Average the test errors.



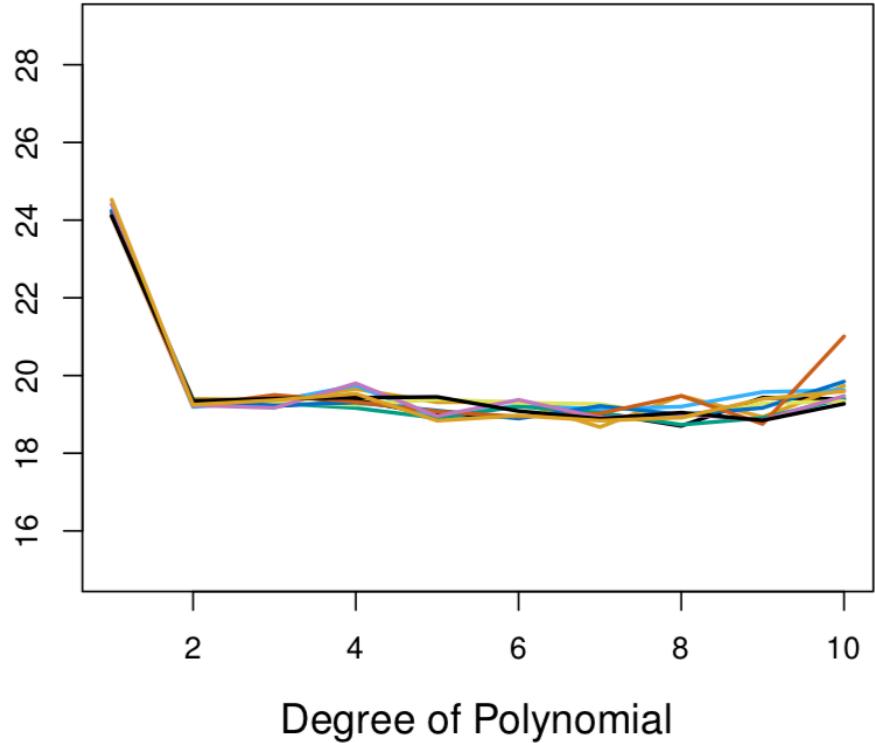
**LOOCV**

Mean Squared Error



**10-fold CV**

Mean Squared Error



# The wrong way to do cross validation

*Reading:* Section 7.10.2 of The Elements of Statistical Learning.

We want to classify 200 individuals according to whether they have cancer or not. We use logistic regression onto 1000 measurements of gene expression.

Proposed strategy:

- ▶ Using all the data, select the 20 most significant genes using  $z$ -tests.
- ▶ Estimate the test error of logistic regression with these 20 predictors via 10-fold cross validation.

---

## The wrong way to do cross validation

To see how that works, let's use the following simulated data:

- ▶ Each gene expression is standard normal and independent of all others.
- ▶ The response (cancer or not) is sampled from a coin flip — no correlation to any of the “genes”.

What should the misclassification rate be for any classification method using these predictors?

Roughly 50%.

---

## The wrong way to do cross validation

We run this simulation, and obtain a CV error rate of 3%!

Why is this?

- ▶ Since we only have 200 individuals in total, among 1000 variables, at least some will be correlated with the response.
- ▶ We do variable selection using *all the data*, so the variables we select have some correlation with the response in every subset or fold in the cross validation.

## The **right** way to do cross validation

- ▶ Divide the data into 10 folds.
- ▶ For  $i = 1, \dots, 10$ :
  - ▶ Using every fold except  $i$ , perform the variable selection and fit the model with the selected variables.
  - ▶ Compute the error on fold  $i$ .
- ▶ Average the 10 test errors obtained.

In our simulation, this produces an error estimate of close to 50%.

**Moral of the story:** Every aspect of the learning method that involves using the data — variable selection, for example — must be cross-validated.